

# Package ‘tidybulk’

April 6, 2021

**Type** Package

**Title** Brings transcriptomics to the tidyverse

**Version** 1.2.1

**Description** This is a collection of utility functions that allow to perform exploration of and calculations to RNA sequencing data, in a modular, pipe-friendly and tidy fashion.

**License** GPL-3

**Depends** R (>= 4.0.0)

**Imports** tibble,  
readr,  
dplyr,  
magrittr,  
tidyr,  
stringr,  
rlang,  
purrr,  
preprocessCore,  
stats,  
parallel,  
utils,  
lifecycle,  
scales,  
SummarizedExperiment,  
methods

**Suggests** BiocStyle,  
testthat,  
vctrs,  
AnnotationDbi,  
BiocManager,  
Rsubread,  
e1071,  
edgeR,  
limma,  
org.Hs.eg.db,  
org.Mm.eg.db,  
sva,  
GGally,  
knitr,

qpdf,  
 covr,  
 Seurat,  
 KernSmooth,  
 Rtsne,  
 S4Vectors,  
 ggplot2,  
 widyr,  
 clusterProfiler,  
 msigdb,  
 DESeq2,  
 broom,  
 survival,  
 boot,  
 betareg,  
 tidyHeatmap,  
 pasilla,  
 ggrepel,  
 devtools,  
 functional

**VignetteBuilder** knitr

**RdMacros** lifecycle

**Biarch** true

**biocViews** AssayDomain, Infrastructure, RNASeq, DifferentialExpression, GeneExpression, Normalization, Clustering, QualityControl, Sequencing, Transcription, Transcriptomics

**Encoding** UTF-8

**LazyData** true

**RoxygenNote** 7.1.1

**git\_url** <https://git.bioconductor.org/packages/tidybulk>

**git\_branch** RELEASE\_3\_12

**git\_last\_commit** 15d7391

**git\_last\_commit\_date** 2021-04-06

**Date/Publication** 2021-04-06

## R topics documented:

adjust_abundance . . . . .	4
aggregate_duplicates . . . . .	6
arrange . . . . .	8
as_matrix . . . . .	10
bind . . . . .	10
breast_tcga_mini . . . . .	11
cluster_elements . . . . .	12
counts . . . . .	14
counts_ensembl . . . . .	15
counts_mini . . . . .	15
deconvolve_cellularity . . . . .	15
describe_transcript . . . . .	18

distinct . . . . .	18
ensembl_symbol_mapping . . . . .	19
ensembl_to_symbol . . . . .	19
fill_missing_abundance . . . . .	20
filter . . . . .	22
flybaseIDs . . . . .	23
full_join . . . . .	24
get_bibliography . . . . .	24
group_by . . . . .	25
identify_abundant . . . . .	26
impute_missing_abundance . . . . .	28
inner_join . . . . .	31
keep_abundant . . . . .	31
keep_variable . . . . .	34
left_join . . . . .	36
log10_reverse_trans . . . . .	37
logit_trans . . . . .	38
mutate . . . . .	38
parse_formula_survival . . . . .	40
pivot_sample . . . . .	40
pivot_transcript . . . . .	41
reduce_dimensions . . . . .	42
remove_redundancy . . . . .	46
rename . . . . .	49
right_join . . . . .	50
rotate_dimensions . . . . .	50
rowwise . . . . .	53
scale_abundance . . . . .	54
se . . . . .	56
se_mini . . . . .	57
summarise . . . . .	57
symbol_to_entrez . . . . .	58
test_deseq2_df . . . . .	59
test_differential_abundance . . . . .	59
test_differential_cellularity . . . . .	63
test_gene_enrichment . . . . .	66
test_gene_overrepresentation . . . . .	69
tidybulk . . . . .	71
tidybulk_SAM_BAM . . . . .	72
unnest . . . . .	73
vignette_manuscript_signature_boxplot . . . . .	74
vignette_manuscript_signature_tsne . . . . .	75
vignette_manuscript_signature_tsne2 . . . . .	75
X_cibersort . . . . .	75

---

adjust\_abundance

Adjust transcript abundance for unwanted variation

---

## Description

adjust\_abundance() takes as input a 'tbl' formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>| and returns a 'tbl' with an additional adjusted abundance column. This method uses scaled counts if present.

## Usage

```
adjust_abundance(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'spec_tbl_df'
adjust_abundance(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'tbl_df'
adjust_abundance(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'tidybulk'
adjust_abundance(
  .data,
  .formula,
  .sample = NULL,
```

```

        .transcript = NULL,
        .abundance = NULL,
        log_transform = TRUE,
        action = "add",
        ...
    )

## S4 method for signature 'SummarizedExperiment'
adjust_abundance(
    .data,
    .formula,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    log_transform = TRUE,
    action = "add",
    ...
)

## S4 method for signature 'RangedSummarizedExperiment'
adjust_abundance(
    .data,
    .formula,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    log_transform = TRUE,
    action = "add",
    ...
)

```

### Arguments

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.formula</code>	A formula with no response variable, representing the desired linear model where the first covariate is the factor of interest and the second covariate is the unwanted variation (of the kind $\sim$ factor_of_intrest + batch)
<code>.sample</code>	The name of the sample column
<code>.transcript</code>	The name of the transcript/gene column
<code>.abundance</code>	The name of the transcript/gene abundance column
<code>log_transform</code>	A boolean, whether the value should be log-transformed (e.g., TRUE for RNA sequencing data)
<code>action</code>	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).
<code>...</code>	Further parameters passed to the function <code>sva::ComBat</code>

### Details

#### [Maturing]

This function adjusts the abundance for (known) unwanted variation. At the moment just an unwanted covariate is allowed at a time using Combat (DOI: 10.1093/bioinformatics/bts034)

Underlying method: `sva::ComBat(data, batch = my_batch, mod = design, prior.plots = FALSE, ...)`

### Value

A 'tbl' with additional columns for the adjusted counts as '<COUNT COLUMN>\_adjusted'  
 A 'tbl' with additional columns for the adjusted counts as '<COUNT COLUMN>\_adjusted'  
 A 'tbl' with additional columns for the adjusted counts as '<COUNT COLUMN>\_adjusted'  
 A 'tbl' with additional columns for the adjusted counts as '<COUNT COLUMN>\_adjusted'  
 A 'SummarizedExperiment' object  
 A 'SummarizedExperiment' object

### Examples

```
cm = tidybulk::counts_mini
cm$batch = 0
cm$batch[cm$sample %in% c("SRR1740035", "SRR1740043")] = 1

res =
  cm %>%
  tidybulk(sample, transcript, count) %>%
  identify_abundant() %>%
  adjust_abundance( ~ condition + batch )
```

---

aggregate_duplicates	<i>Aggregates multiple counts from the same samples (e.g., from isoforms), concatenates other character columns, and averages other numeric columns</i>
----------------------	---

---

### Description

`aggregate_duplicates()` takes as input a 'tbl' formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>| and returns a 'tbl' with aggregated transcripts that were duplicated.

### Usage

```
aggregate_duplicates(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  aggregation_function = sum,
  keep_integer = TRUE
)

## S4 method for signature 'spec_tbl_df'
aggregate_duplicates(
```

```

    .data,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    aggregation_function = sum,
    keep_integer = TRUE
)

## S4 method for signature 'tbl_df'
aggregate_duplicates(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  aggregation_function = sum,
  keep_integer = TRUE
)

## S4 method for signature 'tidybulk'
aggregate_duplicates(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  aggregation_function = sum,
  keep_integer = TRUE
)

## S4 method for signature 'SummarizedExperiment'
aggregate_duplicates(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  aggregation_function = sum,
  keep_integer = TRUE
)

## S4 method for signature 'RangedSummarizedExperiment'
aggregate_duplicates(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  aggregation_function = sum,
  keep_integer = TRUE
)

```

## Arguments

.data	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
.sample	The name of the sample column

.transcript      The name of the transcript/gene column  
.abundance        The name of the transcript/gene abundance column  
aggregation\_function      A function for counts aggregation (e.g., sum, median, or mean)  
keep\_integer      A boolean. Whether to force the aggregated counts to integer

## Details

### [Maturing]

This function aggregates duplicated transcripts (e.g., isoforms, ensembl). For example, we often have to convert ensembl symbols to gene/transcript symbol, but in doing so we have to deal with duplicates. 'aggregate\_duplicates' takes a tibble and column names (as symbols; for 'sample', 'transcript' and 'count') as arguments and returns a tibble with aggregate transcript with the same name. All the rest of the column are appended, and factors and boolean are appended as characters.

Underlying custom method: `data filter(n_aggr > 1) group_by(!.sample,!.transcript) dplyr::mutate(!.abundance := !.abundance`

## Value

A 'tbl' object with aggregated transcript abundance and annotation  
A 'tbl' object with aggregated transcript abundance and annotation  
A 'tbl' object with aggregated transcript abundance and annotation  
A 'tbl' object with aggregated transcript abundance and annotation  
A 'SummarizedExperiment' object  
A 'SummarizedExperiment' object

## Examples

```
aggregate_duplicates(
  tidybulk::counts_mini,
  sample,
  transcript,
  `count`,
  aggregation_function = sum
)
```

---

arrange

drplyr-methods

---

## Description

'arrange()' order the rows of a data frame rows by the values of selected columns.

Unlike other dplyr verbs, 'arrange()' largely ignores grouping; you need to explicit mention grouping variables (or use 'by\_group = TRUE') in order to group by them, and functions of variables are evaluated once per data frame, not once per group.



**Usage**

```

arrange(.data, ..., .by_group = FALSE)

## Default S3 method:
arrange(.data, ..., .by_group = FALSE)

bind_rows(..., .id = NULL)

bind_cols(..., .id = NULL)

ungroup(x, ...)

```

**Arguments**

<code>.data</code>	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from <code>dbplyr</code> or <code>dtplyr</code> ). See <i>*Methods*</i> , below, for more details.
<code>...</code>	<[‘tidy-eval’][dplyr_tidy_eval]> Variables, or functions or variables. Use <code>[desc()]</code> to sort a variable in descending order.
<code>.by_group</code>	If ‘TRUE’, will sort first by grouping variable. Applies to grouped data frames only.
<code>.id</code>	Data frame identifier. When ‘.id’ is supplied, a new column of identifiers is created to link each row to its original data frame. The labels are taken from the named arguments to ‘ <code>bind_rows()</code> ’. When a list of data frames is supplied, the labels are taken from the names of the list. If no names are found a numeric sequence is used instead.
<code>x</code>	A <code>[tbl()]</code>

**Details**

**## Locales** The sort order for character vectors will depend on the collating sequence of the locale in use: see `[locales()]`.

**## Missing values** Unlike base sorting with ‘`sort()`’, ‘NA’ are: \* always sorted to the end for local data, even when wrapped with ‘`desc()`’. \* treated differently for remote data, depending on the backend.

**Value**

A tibble Arrange rows by column values

An object of the same type as ‘.data’.

\* All rows appear in the output, but (usually) in a different place. \* Columns are not modified. \* Groups are not modified. \* Data frame attributes are preserved.

**Methods**

This function is a *\*\*generic\*\**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages:

See Also

Other single table verbs: `filter()`, `mutate()`, `rename()`, `summarise()`

Examples

```
`%>%` = magrittr::`%>%`  
arrange(mtcars, cyl, disp)
```

---

as_matrix	<i>Get matrix from tibble</i>
-----------	-------------------------------

---

Description

Get matrix from tibble

Usage

```
as_matrix(tbl, rownames = NULL, do_check = TRUE)
```

Arguments

- tbl                   A tibble
- rownames            A character string of the rownames
- do\_check            A boolean

Value

A matrix

Examples

```
as_matrix(head(dplyr::select(tidybulk::counts_mini, transcript, count)), rownames=transcript)
```

---

bind	<i>Efficiently bind multiple data frames by row and column</i>
------	--

---

Description

This is an efficient implementation of the common pattern of ‘do.call(rbind, dfs)’ or ‘do.call(cbind, dfs)’ for binding many data frames into one.

**Arguments**

- `...` Data frames to combine.  
 Each argument can either be a data frame, a list that could be a data frame, or a list of data frames.  
 When row-binding, columns are matched by name, and any missing columns will be filled with NA.  
 When column-binding, rows are matched by position, so all data frames must have the same number of rows. To match by value, not position, see [mutate-joins].
- `.id` Data frame identifier.  
 When `'id'` is supplied, a new column of identifiers is created to link each row to its original data frame. The labels are taken from the named arguments to `'bind_rows()'`. When a list of data frames is supplied, the labels are taken from the names of the list. If no names are found a numeric sequence is used instead.

**Details**

The output of `'bind_rows()'` will contain a column if that column appears in any of the inputs.

**Value**

`'bind_rows()'` and `'bind_cols()'` return the same type as the first input, either a data frame, `'tbl_df'`, or `'grouped_df'`.

**Examples**

```
`%>%` = magrittr::`%>%`
one <- mtcars[1:4, ]
two <- mtcars[11:14, ]

# You can supply data frames as arguments:
bind_rows(one, two)
```

---

breast_tcga_mini	<i>Data set</i>
------------------	-----------------

---

**Description**

Data set

**Usage**

```
breast_tcga_mini
```

**Format**

An object of class `tbl_df` (inherits from `tbl`, `data.frame`) with 125500 rows and 5 columns.

---

cluster_elements	<i>Get clusters of elements (e.g., samples or transcripts)</i>
------------------	--

---

### Description

cluster\_elements() takes as input a 'tbl' formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>| and identify clusters in the data.

### Usage

```
cluster_elements(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'spec_tbl_df'
cluster_elements(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'tbl_df'
cluster_elements(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'tidybulk'
cluster_elements(
  .data,
```

```

    .element = NULL,
    .feature = NULL,
    .abundance = NULL,
    method,
    of_samples = TRUE,
    log_transform = TRUE,
    action = "add",
    ...
)

## S4 method for signature 'SummarizedExperiment'
cluster_elements(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  log_transform = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'RangedSummarizedExperiment'
cluster_elements(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  log_transform = TRUE,
  action = "add",
  ...
)

```

## Arguments

.data	A 'tbl' formatted as  <SAMPLE> <TRANSCRIPT> <COUNT> <...>
.element	The name of the element column (normally samples).
.feature	The name of the feature column (normally transcripts/genes)
.abundance	The name of the column including the numerical value the clustering is based on (normally transcript abundance)
method	A character string. The cluster algorithm to use, at the moment k-means is the only algorithm included.
of_samples	A boolean. In case the input is a tidybulk object, it indicates Whether the element column will be sample or transcript column
log_transform	A boolean, whether the value should be log-transformed (e.g., TRUE for RNA sequencing data)
action	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).

... Further parameters passed to the function kmeans

## Details

### [Maturing]

identifies clusters in the data, normally of samples. This function returns a tibble with additional columns for the cluster annotation. At the moment only k-means (DOI: 10.2307/2346830) and SNN clustering (DOI:10.1016/j.cell.2019.05.031) is supported, the plan is to introduce more clustering methods.

Underlying method for kmeans `do.call(kmeans(.data, iter.max = 1000, ...))`

Underlying method for SNN `Seurat::CreateSeuratObject() Seurat::ScaleData(display.progress = TRUE, num.cores = 4, do.par = TRUE) Seurat::FindVariableFeatures(selection.method = "vst") Seurat::RunPCA(npcs = 30) Seurat::FindNeighbors() Seurat::FindClusters(method = "igraph", ...)`

## Value

A tibble object with additional columns with cluster labels

A tibble object with additional columns with cluster labels

A tibble object with additional columns with cluster labels

A tibble object with additional columns with cluster labels

A 'SummarizedExperiment' object

A 'SummarizedExperiment' object

## Examples

```
cluster_elements(tidybulk::counts_mini, sample, transcript, count, centers = 2, method="kmeans")
```

---

counts	<i>Example data set</i>
--------	-------------------------

---

## Description

Example data set

## Usage

```
counts
```

## Format

An object of class `tibble_df` (inherits from `tibble`, `data.frame`) with 408624 rows and 8 columns.

---

counts_ensembl	<i>Counts with ensembl annotation</i>
----------------	---------------------------------------

---

**Description**

Counts with ensembl annotation

**Usage**

```
counts_ensembl
```

**Format**

An object of class `tbl_df` (inherits from `tbl`, `data.frame`) with 119 rows and 6 columns.

---

counts_mini	<i>Example data set reduced</i>
-------------	---------------------------------

---

**Description**

Example data set reduced

**Usage**

```
counts_mini
```

**Format**

An object of class `spec_tbl_df` (inherits from `tbl_df`, `tbl`, `data.frame`) with 2635 rows and 6 columns.

---

deconvolve_cellularity	<i>Get cell type proportions from samples</i>
------------------------	---

---

**Description**

`deconvolve_cellularity()` takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a 'tbl' with the estimated cell type abundance for each sample

**Usage**

```
deconvolve_cellularity(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  reference = X_cibersort,  
  method = "cibersort",  
  action = "add",  
  ...  
)  
  
## S4 method for signature 'spec_tbl_df'  
deconvolve_cellularity(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  reference = X_cibersort,  
  method = "cibersort",  
  action = "add",  
  ...  
)  
  
## S4 method for signature 'tbl_df'  
deconvolve_cellularity(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  reference = X_cibersort,  
  method = "cibersort",  
  action = "add",  
  ...  
)  
  
## S4 method for signature 'tidybulk'  
deconvolve_cellularity(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  reference = X_cibersort,  
  method = "cibersort",  
  action = "add",  
  ...  
)  
  
## S4 method for signature 'SummarizedExperiment'  
deconvolve_cellularity(  
  .data,  
  .sample = NULL,
```



```

        .transcript = NULL,
        .abundance = NULL,
        reference = X_cibersort,
        method = "cibersort",
        action = "add",
        ...
    )

## S4 method for signature 'RangedSummarizedExperiment'
deconvolve_cellularity(
    .data,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    reference = X_cibersort,
    method = "cibersort",
    action = "add",
    ...
)

```

### Arguments

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.sample</code>	The name of the sample column
<code>.transcript</code>	The name of the transcript/gene column
<code>.abundance</code>	The name of the transcript/gene abundance column
<code>reference</code>	A data frame. The transcript/cell_type data frame of integer transcript abundance
<code>method</code>	A character string. The method to be used. At the moment Cibersort (default) and llr (linear least squares regression) are available.
<code>action</code>	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).
<code>...</code>	Further parameters passed to the function Cibersort

### Details

#### [Maturing]

This function infers the cell type composition of our samples (with the algorithm Cibersort; Newman et al., 10.1038/nmeth.3337).

Underlying method: CIBERSORT(Y = data, X = reference, ...)

### Value

A 'tbl' object including additional columns for each cell type estimated  
 A 'tbl' object including additional columns for each cell type estimated  
 A 'tbl' object including additional columns for each cell type estimated  
 A 'tbl' object including additional columns for each cell type estimated  
 A 'SummarizedExperiment' object  
 A 'SummarizedExperiment' object

Examples

```
# Subsetting for time efficiency
deconvolve_cellularity(filter(tidybulk::counts, sample=="SRR1740034"), sample, transcript, `count`, cores = 1)
```

---

describe_transcript	<i>Get DESCRIPTION from gene SYMBOL for Human and Mouse</i>
---------------------	---

---

Description

Get DESCRIPTION from gene SYMBOL for Human and Mouse

Usage

```
describe_transcript(.data, .transcript = NULL)
```

Arguments

- .data           A tt or tbl object.
- .transcript    A character. The name of the gene symbol column.

Value

A tbl

Examples

```
describe_transcript(tidybulk::counts_mini, .transcript = transcript)
```

---

distinct	<i>distinct</i>
----------	-----------------

---

Description

distinct

Usage

```
distinct(.data, ..., .keep_all = FALSE)
```

Arguments

- .data           A tbl. (See dplyr)
- ...            Data frames to combine (See dplyr)
- .keep\_all      If TRUE, keep all variables in .data. If a combination of ... is not distinct, this keeps the first row of values. (See dplyr)

**Value**

A tt object

**Examples**

```
distinct(tidybulk::counts_mini)
```

---

```
ensembl_symbol_mapping
```

*Data set*

---

**Description**

Data set

**Usage**

```
ensembl_symbol_mapping
```

**Format**

An object of class `spec_tbl_df` (inherits from `tbl_df`, `tbl`, `data.frame`) with 291249 rows and 3 columns.

---

```
ensembl_to_symbol
```

*Add transcript symbol column from ensembl id for human and mouse data*

---

**Description**

`ensembl_to_symbol()` takes as input a 'tbl' formatted as |<SAMPLE>|<ENSEMBL\_ID>|<COUNT>|<...>| and returns a 'tbl' with the additional transcript symbol column

**Usage**

```
ensembl_to_symbol(.data, .ensembl, action = "add")
```

```
## S4 method for signature 'spec_tbl_df'
ensembl_to_symbol(.data, .ensembl, action = "add")
```

```
## S4 method for signature 'tbl_df'
ensembl_to_symbol(.data, .ensembl, action = "add")
```

```
## S4 method for signature 'tidybulk'
ensembl_to_symbol(.data, .ensembl, action = "add")
```

**Arguments**

.data	A 'tbl' formatted as   <SAMPLE>   <ENSEMBL_ID>   <COUNT>   <...>
.ensembl	A character string. The column that is represents ensembl gene id
action	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).

**Details****[Questioning]**

This is useful since different resources use ensembl IDs while others use gene symbol IDs. At the moment this work for human (genes and transcripts) and mouse (genes) data.

**Value**

A 'tbl' object including additional columns for transcript symbol  
 A 'tbl' object including additional columns for transcript symbol  
 A 'tbl' object including additional columns for transcript symbol  
 A 'tbl' object including additional columns for transcript symbol

**Examples**

```
ensembl_to_symbol(tidybulk::counts_ensembl, ens)
```

---

```
fill_missing_abundance
```

*Fill transcript abundance if missing from sample-transcript pairs*

---

**Description**

fill\_missing\_abundance() takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a 'tbl' with new observations

**Usage**

```
fill_missing_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  fill_with
)

## S4 method for signature 'spec_tbl_df'
fill_missing_abundance(
  .data,
```

```

    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    fill_with
  )

## S4 method for signature 'tbl_df'
fill_missing_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  fill_with
)

## S4 method for signature 'tidybulk'
fill_missing_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  fill_with
)

## S4 method for signature 'SummarizedExperiment'
fill_missing_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  fill_with
)

## S4 method for signature 'RangedSummarizedExperiment'
fill_missing_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  fill_with
)

```

### Arguments

.data	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
.sample	The name of the sample column
.transcript	The name of the transcript column
.abundance	The name of the transcript abundance column
fill_with	A numerical abundance with which fill the missing data points

## Details

### [Maturing]

This function fills the abundance of missing sample-transcript pair using the median of the sample group defined by the formula

## Value

A 'tbl' non-sparse abundance  
 A 'tbl' with filled abundance  
 A 'tbl' with filled abundance  
 A 'tbl' with filled abundance  
 A 'SummarizedExperiment' object  
 A 'SummarizedExperiment' object

## Examples

```
fill_missing_abundance(tidybulk::counts_mini, sample, transcript, count, fill_with = 0)
```

---

filter	<i>Subset rows using column values</i>
--------	--

---

## Description

'filter()' retains the rows where the conditions you provide a 'TRUE'. Note that, unlike base subsetting with '[', rows where the condition evaluates to 'NA' are dropped.

## Usage

```
filter(.data, ..., .preserve = FALSE)
```

## Arguments

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>*Methods*</i> , below, for more details.
...	<['tidy-eval'] [dplyr_tidy_eval]> Logical predicates defined in terms of the variables in '.data'. Multiple conditions are combined with '&'. Only rows where the condition evaluates to 'TRUE' are kept.
.preserve	when 'FALSE' (the default), the grouping structure is recalculated based on the resulting data, otherwise it is kept as is.

## Details

dplyr is not yet smart enough to optimise filtering optimisation on grouped datasets that don't need grouped calculations. For this reason, filtering is often considerably faster on [ungroup()]ed data.

**Value**

An object of the same type as `‘.data’`.

\* Rows are a subset of the input, but appear in the same order. \* Columns are not modified. \* The number of groups may be reduced (if `‘.preserve’` is not `‘TRUE’`). \* Data frame attributes are preserved.

**Useful filter functions**

\* `['==', ['>'], ['>=']` etc \* `['&'], ['|'], [';'], [xor()]` \* `[is.na()]` \* `[between()], [near()]`

**Grouped tibbles**

Because filtering expressions are computed within groups, they may yield different results on grouped tibbles. This will be the case as soon as an aggregating, lagging, or ranking function is involved. Compare this ungrouped filtering:

The former keeps rows with `‘mass’` greater than the global average whereas the latter keeps rows with `‘mass’` greater than the gender average.

**Methods**

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages:

**See Also**

`[filter_all()], [filter_if()]` and `[filter_at()]`.

Other single table verbs: `arrange()`, `mutate()`, `rename()`, `summarise()`

**Examples**

```
# Learn more in ?dplyr_tidy_eval
```

---

flybaseIDs

*flybaseIDs*

---

**Description**

flybaseIDs

**Usage**

flybaseIDs

**Format**

An object of class character of length 14599.

---

full_join	<i>Full join datasets</i>
-----------	---------------------------

---

**Description**

Full join datasets

**Usage**

```
full_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ...)
```

**Arguments**

x	tbls to join. (See dplyr)
y	tbls to join. (See dplyr)
by	A character vector of variables to join by. (See dplyr)
copy	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. (See dplyr)
suffix	If there are non-joined duplicate variables in x and y, these suffixes will be added to the output to disambiguate them. Should be a character vector of length 2. (See dplyr)
...	Data frames to combine (See dplyr)

**Value**

A tt object

**Examples**

```
`%>%` = magrittr::`%>%`
annotation = tidybulk::counts %>% distinct(sample) %>% mutate(source = "AU")
tidybulk::counts %>% full_join(annotation)
```

---

get_bibliography	<i>Produces the bibliography list of your workflow</i>
------------------	--

---

**Description**

get\_bibliography() takes as input a 'tidybulk'

**Usage**

```
get_bibliography(.data)

## S4 method for signature 'tidybulk'
get_bibliography(.data)
```



**Arguments**

.data                      A 'tidybulk' tibble

**Details****[Maturing]**

This methods returns the bibliography list of your workflow from the internals of a tidybulk tibble (attr(., "internals"))

**Value**

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

**Examples**

```
# Define tidybulk tibble
df = tidybulk(tidybulk::counts_mini, sample, transcript, count)

get_bibliography(df)
```

---

group\_by

*Group by one or more variables*


---

**Description**

Most data operations are done on groups defined by variables. 'group\_by()' takes an existing tbl and converts it into a grouped tbl where operations are performed "by group". 'ungroup()' removes grouping.

**Usage**

```
group_by(.data, ..., .add = FALSE, .drop = group_by_drop_default(.data))
```

**Arguments**

.data	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>*Methods*</i> , below, for more details.
...	In 'group_by()', variables or computations to group by. In 'ungroup()', variables to remove from the grouping.
.add	When 'FALSE', the default, 'group_by()' will override existing groups. To add to the existing groups, use '.add = TRUE'. This argument was previously called 'add', but that prevented creating a new grouping variable called 'add', and conflicts with our naming conventions.
.drop	When '.drop = TRUE', empty groups are dropped. See [group_by_drop_default()] for what the default value is for this argument.

**Value**

A [grouped data frame][grouped\_df()], unless the combination of ‘...’ and ‘add’ yields a non empty set of grouping columns, a regular (ungrouped) data frame otherwise.

**Methods**

These function are **generic**s, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

**Examples**

```
`%>%` = magrittr::`%>%`
by_cyl <- mtcars %>% group_by(cyl)
```

---

identify_abundant	<i>find abundant transcripts</i>
-------------------	----------------------------------

---

**Description**

identify\_abundant() takes as input a ‘tbl’ formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>| and returns a ‘tbl’ with additional columns for the statistics from the hypothesis test.

**Usage**

```
identify_abundant(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  factor_of_interest = NULL,
  minimum_counts = 10,
  minimum_proportion = 0.7
)

## S4 method for signature 'spec_tbl_df'
identify_abundant(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  factor_of_interest = NULL,
  minimum_counts = 10,
  minimum_proportion = 0.7
)

## S4 method for signature 'tbl_df'
identify_abundant(
  .data,
```

```

        .sample = NULL,
        .transcript = NULL,
        .abundance = NULL,
        factor_of_interest = NULL,
        minimum_counts = 10,
        minimum_proportion = 0.7
    )

## S4 method for signature 'tidybulk'
identify_abundant(
    .data,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    factor_of_interest = NULL,
    minimum_counts = 10,
    minimum_proportion = 0.7
)

## S4 method for signature 'SummarizedExperiment'
identify_abundant(
    .data,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    factor_of_interest = NULL,
    minimum_counts = 10,
    minimum_proportion = 0.7
)

## S4 method for signature 'RangedSummarizedExperiment'
identify_abundant(
    .data,
    .sample = NULL,
    .transcript = NULL,
    .abundance = NULL,
    factor_of_interest = NULL,
    minimum_counts = 10,
    minimum_proportion = 0.7
)

```

### Arguments

.data	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
.sample	The name of the sample column
.transcript	The name of the transcript/gene column
.abundance	The name of the transcript/gene abundance column
factor_of_interest	The name of the column of the factor of interest. This is used for defining sample groups for the filtering process. It uses the filterByExpr function from edgeR.
minimum_counts	A real positive number. It is the threshold of count per million that is used to filter transcripts/genes out from the scaling procedure.

minimum\_proportion

A real positive number between 0 and 1. It is the threshold of proportion of samples for each transcripts/genes that have to be characterised by a cmp bigger than the threshold to be included for scaling procedure.

## Details

### [Maturing]

At the moment this function uses edgeR (DOI: 10.1093/bioinformatics/btp616)

Underlying method: `edgeR::filterByExpr( data, min.count = minimum_counts, group = string_factor_of_interest, min.prop = minimum_proportion )`

## Value

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'SummarizedExperiment' object

A 'SummarizedExperiment' object

## Examples

```
identify_abundant(
  tidybulk::counts_mini,
  sample,
  transcript,
  `count`
)
```

---

impute\_missing\_abundance

*impute transcript abundance if missing from sample-transcript pairs*

---

## Description

`impute_missing_abundance()` takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a 'tbl' with an additional adjusted abundance column. This method uses scaled counts if present.

**Usage**

```
impute_missing_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL  
)  
  
## S4 method for signature 'spec_tbl_df'  
impute_missing_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL  
)  
  
## S4 method for signature 'tbl_df'  
impute_missing_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL  
)  
  
## S4 method for signature 'tidybulk'  
impute_missing_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL  
)  
  
## S4 method for signature 'SummarizedExperiment'  
impute_missing_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL  
)  
  
## S4 method for signature 'RangedSummarizedExperiment'  
impute_missing_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL  
)
```

)

## Arguments

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.formula</code>	A formula with no response variable, representing the desired linear model where the first covariate is the factor of interest and the second covariate is the unwanted variation (of the kind <code>~ factor_of_intrest + batch</code> )
<code>.sample</code>	The name of the sample column
<code>.transcript</code>	The name of the transcript/gene column
<code>.abundance</code>	The name of the transcript/gene abundance column

## Details

### [Maturing]

This function imputes the abundance of missing sample-transcript pair using the median of the sample group defined by the formula

## Value

- A 'tbl' non-sparse abundance
- A 'tbl' with imputed abundance
- A 'tbl' with imputed abundance
- A 'tbl' with imputed abundance
- A 'SummarizedExperiment' object
- A 'SummarizedExperiment' object

## Examples

```
res =
  impute_missing_abundance(
    tidybulk::counts_mini,
    ~ condition,
    .sample = sample,
    .transcript = transcript,
    .abundance = count
  )
```

---

inner_join	<i>Inner join datasets</i>
------------	----------------------------

---

**Description**

Inner join datasets

**Usage**

```
inner_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ...)
```

**Arguments**

x	tbls to join. (See dplyr)
y	tbls to join. (See dplyr)
by	A character vector of variables to join by. (See dplyr)
copy	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. (See dplyr)
suffix	If there are non-joined duplicate variables in x and y, these suffixes will be added to the output to disambiguate them. Should be a character vector of length 2. (See dplyr)
...	Data frames to combine (See dplyr)

**Value**

A tt object

**Examples**

```
`%>%` = magrittr::`%>%`
annotation = tidybulk::counts %>% distinct(sample) %>% mutate(source = "AU")
tidybulk::counts %>% inner_join(annotation)
```

---

keep_abundant	<i>Keep abundant transcripts</i>
---------------	----------------------------------

---

**Description**

keep\_abundant() takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a 'tbl' with additional columns for the statistics from the hypothesis test.

**Usage**

```
keep_abundant(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  factor_of_interest = NULL,  
  minimum_counts = 10,  
  minimum_proportion = 0.7  
)  
  
## S4 method for signature 'spec_tbl_df'  
keep_abundant(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  factor_of_interest = NULL,  
  minimum_counts = 10,  
  minimum_proportion = 0.7  
)  
  
## S4 method for signature 'tbl_df'  
keep_abundant(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  factor_of_interest = NULL,  
  minimum_counts = 10,  
  minimum_proportion = 0.7  
)  
  
## S4 method for signature 'tidybulk'  
keep_abundant(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  factor_of_interest = NULL,  
  minimum_counts = 10,  
  minimum_proportion = 0.7  
)  
  
## S4 method for signature 'SummarizedExperiment'  
keep_abundant(  
  .data,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  factor_of_interest = NULL,  
  minimum_counts = 10,
```



```

    minimum_proportion = 0.7
)

## S4 method for signature 'RangedSummarizedExperiment'
keep_abundant(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  factor_of_interest = NULL,
  minimum_counts = 10,
  minimum_proportion = 0.7
)

```

### Arguments

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.sample</code>	The name of the sample column
<code>.transcript</code>	The name of the transcript/gene column
<code>.abundance</code>	The name of the transcript/gene abundance column
<code>factor_of_interest</code>	The name of the column of the factor of interest. This is used for defining sample groups for the filtering process. It uses the <code>filterByExpr</code> function from <code>edgeR</code> .
<code>minimum_counts</code>	A real positive number. It is the threshold of count per million that is used to filter transcripts/genes out from the scaling procedure.
<code>minimum_proportion</code>	A real positive number between 0 and 1. It is the threshold of proportion of samples for each transcripts/genes that have to be characterised by a <code>cmp</code> bigger than the threshold to be included for scaling procedure.

### Details

#### [Questioning]

At the moment this function uses `edgeR` (DOI: 10.1093/bioinformatics/btp616)

Underlying method: `edgeR::filterByExpr( data, min.count = minimum_counts, group = string_factor_of_interest, min.prop = minimum_proportion )`

### Value

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'SummarizedExperiment' object

A 'SummarizedExperiment' object

## Examples

```
keep_abundant(
  tidybulk::counts_mini,
    sample,
    transcript,
    `count`
)
```

---

keep_variable	<i>Keep variable transcripts</i>
---------------	----------------------------------

---

## Description

keep\_variable() takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a 'tbl' with additional columns for the statistics from the hypothesis test.

## Usage

```
keep_variable(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  top = 500,
  log_transform = TRUE
)
```

```
## S4 method for signature 'spec_tbl_df'
```

```
keep_variable(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  top = 500,
  log_transform = TRUE
)
```

```
## S4 method for signature 'tbl_df'
```

```
keep_variable(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  top = 500,
  log_transform = TRUE
)
```

```
## S4 method for signature 'tidybulk'
keep_variable(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  top = 500,
  log_transform = TRUE
)

## S4 method for signature 'SummarizedExperiment'
keep_variable(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  top = 500,
  log_transform = TRUE
)

## S4 method for signature 'RangedSummarizedExperiment'
keep_variable(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  top = 500,
  log_transform = TRUE
)
```

### Arguments

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.sample</code>	The name of the sample column
<code>.transcript</code>	The name of the transcript/gene column
<code>.abundance</code>	The name of the transcript/gene abundance column
<code>top</code>	Integer. Number of top transcript to consider
<code>log_transform</code>	A boolean, whether the value should be log-transformed (e.g., TRUE for RNA sequencing data)

### Details

#### [Maturing]

At the moment this function uses edgeR <https://doi.org/10.1093/bioinformatics/btp616>

### Value

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

Underlying method: `s <- rowMeans((x - rowMeans(x)) ^ 2)` `o <- order(s, decreasing = TRUE)` `x <- x[o[1L:top], , drop = FALSE]` `variable_transcripts = rownames(x)`

A ‘tbl’ with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A ‘tbl’ with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A ‘tbl’ with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A ‘SummarizedExperiment’ object

A ‘SummarizedExperiment’ object

Examples

```
keep_variable(  
  tidybulk::counts_mini,  
    sample,  
    transcript,  
    `count`,  
    top = 500  
)
```

---

left_join	Left join datasets
-----------	--------------------

---

Description

Left join datasets

Usage

```
left_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ...)
```

Arguments

x	tbls to join. (See dplyr)
y	tbls to join. (See dplyr)
by	A character vector of variables to join by. (See dplyr)
copy	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. (See dplyr)
suffix	If there are non-joined duplicate variables in x and y, these suffixes will be added to the output to disambiguate them. Should be a character vector of length 2. (See dplyr)
...	Data frames to combine (See dplyr)

**Value**

A tt object

**Examples**

```
`%>%` = magrittr::`%>%`
annotation = tidybulk::counts %>% distinct(sample) %>% mutate(source = "AU")
tidybulk::counts %>% left_join(annotation)
```

---

log10_reverse_trans	<i>log10_reverse_trans</i>
---------------------	----------------------------

---

**Description**

it perform log scaling and reverse the axis. Useful to plot negative log probabilities. To not be used directly but with ggplot (e.g. `scale_y_continuous(trans = "log10_reverse")`)

**Usage**

```
log10_reverse_trans()
```

**Details**

**[Maturing]**

**Value**

A scales object

**Examples**

```
library(ggplot2)
library(tibble)

tibble(pvalue = c(0.001, 0.05, 0.1), fold_change = 1:3) %>%
  ggplot(aes(fold_change, pvalue)) +
  geom_point() +
  scale_y_continuous(trans = "log10_reverse")
```

---

logit_trans	<i>logit scale</i>
-------------	--------------------

---

**Description**

it perform logit scaling with right axis formatting. To not be used directly but with ggplot (e.g. `scale_y_continuous(trans = "log10_reverse")`)

**Usage**

```
logit_trans()
```

**Details**

**[Maturing]**

**Value**

A scales object

**Examples**

```
library(ggplot2)
library(tibble)

tibble(pvalue = c(0.001, 0.05, 0.1), fold_change = 1:3) %>%
  ggplot(aes(fold_change , pvalue)) +
  geom_point() +
  scale_y_continuous(trans = "log10_reverse")
```

---

mutate	<i>Create, modify, and delete columns</i>
--------	---

---

**Description**

‘mutate()’ adds new variables and preserves existing ones; ‘transmute()’ adds new variables and drops existing ones. New variables overwrite existing variables of the same name. Variables can be removed by setting their value to ‘NULL’.

**Usage**

```
mutate(.data, ...)
```

## Arguments

- `.data` A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from `dbplyr` or `dtplyr`). See *\*Methods\**, below, for more details.
- `...` `<[‘tidy-eval’][dplyr_tidy_eval]>` Name-value pairs. The name gives the name of the column in the output.  
The value can be:
  - \* A vector of length 1, which will be recycled to the correct length. \* A vector the same length as the current group (or the whole data frame if ungrouped).
  - \* ‘NULL’, to remove the column. \* A data frame or tibble, to create multiple columns in the output.

## Value

An object of the same type as ‘.data’.

For ‘mutate()’:

\* Rows are not affected. \* Existing columns will be preserved unless explicitly modified. \* New columns will be added to the right of existing columns. \* Columns given value ‘NULL’ will be removed. \* Groups will be recomputed if a grouping variable is mutated. \* Data frame attributes are preserved.

For ‘transmute()’:

\* Rows are not affected. \* Apart from grouping variables, existing columns will be removed unless explicitly kept. \* Column order matches order of expressions. \* Groups will be recomputed if a grouping variable is mutated. \* Data frame attributes are preserved.

## Useful mutate functions

- \* [‘+’], [‘-’], [log()], etc., for their usual mathematical meanings
- \* [lead()], [lag()]
- \* [dense\_rank()], [min\_rank()], [percent\_rank()], [row\_number()], [cume\_dist()], [ntile()]
- \* [cumsum()], [cummean()], [cummin()], [cummax()], [cumany()], [cumall()]
- \* [na\_if()], [coalesce()]
- \* [if\_else()], [recode()], [case\_when()]

## Grouped tibbles

Because mutating expressions are computed within groups, they may yield different results on grouped tibbles. This will be the case as soon as an aggregating, lagging, or ranking function is involved. Compare this ungrouped mutate:

With the grouped equivalent:

The former normalises ‘mass’ by the global average whereas the latter normalises by the averages within gender levels.

## Methods

These functions are *\*\*generic\*\**s, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

Methods available in currently loaded packages:

**See Also**

Other single table verbs: [arrange\(\)](#), [filter\(\)](#), [rename\(\)](#), [summarise\(\)](#)

**Examples**

```
`%>%` = magrittr::`%>%`  
# Newly created variables are available immediately  
mtcars %>% as_tibble() %>% mutate(  
  cyl2 = cyl * 2,  
  cyl4 = cyl2 * 2  
)
```

---

parse_formula_survival	<i>Formula parser with survival</i>
------------------------	-------------------------------------

---

**Description**

Formula parser with survival

**Usage**

```
parse_formula_survival(fm)
```

**Arguments**

fm                      A formula

**Value**

A character vector

---

pivot_sample	<i>Extract sample-wise information</i>
--------------	--

---

**Description**

`pivot_sample()` takes as input a ‘tbl’ formatted as | <SAMPLE> | <ENSEMBL\_ID> | <COUNT> | <...> | and returns a ‘tbl’ with only sample-related columns

**Usage**

```
pivot_sample(.data, .sample = NULL)  
  
## S4 method for signature 'spec_tbl_df'  
pivot_sample(.data, .sample = NULL)  
  
## S4 method for signature 'tbl_df'  
pivot_sample(.data, .sample = NULL)  
  
## S4 method for signature 'tidybulk'  
pivot_sample(.data, .sample = NULL)
```



**Arguments**

`.data` A 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> |

`.sample` The name of the sample column

**Details****[Maturing]**

This function extracts only sample-related information for downstream analysis (e.g., visualisation). It is disruptive in the sense that it cannot be passed anymore to tidybulk function.

**Value**

A 'tbl' object  
 A 'tbl' object  
 A 'tbl' object  
 A 'tbl' object

**Examples**

```

pivot_sample(
  tidybulk::counts_mini,
  .sample = sample
)

```

---

pivot_transcript	<i>Extract transcript-wise information</i>
------------------	--

---

**Description**

pivot\_transcript() takes as input a 'tbl' formatted as | <SAMPLE> | <ENSEMBL\_ID> | <COUNT> | <...> | and returns a 'tbl' with only sample-related columns

**Usage**

```

pivot_transcript(.data, .transcript = NULL)

## S4 method for signature 'spec_tbl_df'
pivot_transcript(.data, .transcript = NULL)

## S4 method for signature 'tbl_df'
pivot_transcript(.data, .transcript = NULL)

## S4 method for signature 'tidybulk'
pivot_transcript(.data, .transcript = NULL)

```

**Arguments**

`.data` A 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> |  
`.transcript` The name of the transcript column

**Details****[Maturing]**

This function extracts only transcript-related information for downstream analysis (e.g., visualisation). It is disruptive in the sense that it cannot be passed anymore to tidybulk function.

**Value**

A 'tbl' object  
 A 'tbl' object  
 A 'tbl' object  
 A 'tbl' object

**Examples**

```
pivot_transcript(
  tidybulk::counts_mini,
  .transcript = transcript
)
```

---

reduce_dimensions	<i>Dimension reduction of the transcript abundance data</i>
-------------------	---

---

**Description**

reduce\_dimensions() takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and calculates the reduced dimensional space of the transcript abundance.

**Usage**

```
reduce_dimensions(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  .dims = 2,
  top = 500,
  of_samples = TRUE,
  log_transform = TRUE,
  scale = TRUE,
  action = "add",
```

```

    ...
)

## S4 method for signature 'spec_tbl_df'
reduce_dimensions(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  .dims = 2,
  top = 500,
  of_samples = TRUE,
  log_transform = TRUE,
  scale = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'tbl_df'
reduce_dimensions(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  .dims = 2,
  top = 500,
  of_samples = TRUE,
  log_transform = TRUE,
  scale = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'tidybulk'
reduce_dimensions(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  .dims = 2,
  top = 500,
  of_samples = TRUE,
  log_transform = TRUE,
  scale = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'SummarizedExperiment'
```

```

reduce_dimensions(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  .dims = 2,
  top = 500,
  of_samples = TRUE,
  log_transform = TRUE,
  scale = TRUE,
  action = "add",
  ...
)

## S4 method for signature 'RangedSummarizedExperiment'
reduce_dimensions(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  .dims = 2,
  top = 500,
  of_samples = TRUE,
  log_transform = TRUE,
  scale = TRUE,
  action = "add",
  ...
)

```

### Arguments

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.element</code>	The name of the element column (normally samples).
<code>.feature</code>	The name of the feature column (normally transcripts/genes)
<code>.abundance</code>	The name of the column including the numerical value the clustering is based on (normally transcript abundance)
<code>method</code>	A character string. The dimension reduction algorithm to use (PCA, MDS, tSNE).
<code>.dims</code>	An integer. The number of dimensions your are interested in (e.g., 4 for returning the first four principal components).
<code>top</code>	An integer. How many top genes to select for dimensionality reduction
<code>of_samples</code>	A boolean. In case the input is a tidybulk object, it indicates Whether the element column will be sample or transcript column
<code>log_transform</code>	A boolean, whether the value should be log-transformed (e.g., TRUE for RNA sequencing data)
<code>scale</code>	A boolean for method="PCA", this will be passed to the 'prcomp' function. It is not included in the ... argument because although the default for 'prcomp' if FALSE, it is advisable to set it as TRUE.

action	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).
...	Further parameters passed to the function prcomp if you choose method="PCA" or Rtsne if you choose method="tSNE"

## Details

### [Maturing]

This function reduces the dimensions of the transcript abundances. It can use multi-dimensional scaling (MDS; DOI.org/10.1186/gb-2010-11-3-r25), principal component analysis (PCA), or tSNE (Jesse Krijthe et al. 2018)

Underlying method for PCA: prcomp(scale = scale, ...)

Underlying method for MDS: limma::plotMDS(ndim = .dims, plot = FALSE, top = top)

Underlying method for tSNE: Rtsne::Rtsne(data, ...)

## Value

A tbl object with additional columns for the reduced dimensions

A tbl object with additional columns for the reduced dimensions

A tbl object with additional columns for the reduced dimensions

A tbl object with additional columns for the reduced dimensions

A ‘SummarizedExperiment’ object

A ‘SummarizedExperiment’ object

## Examples

```
counts.MDS =
  tidybulk::counts_mini %>%
  tidybulk(sample, transcript, count) %>%
  identify_abundant() %>%
  reduce_dimensions( method="MDS", .dims = 3)
```

```
counts.PCA =
  tidybulk::counts_mini %>%
  tidybulk(sample, transcript, count) %>%
  identify_abundant() %>%
  reduce_dimensions(method="PCA", .dims = 3)
```

---

remove_redundancy	<i>Drop redundant elements (e.g., samples) for which feature (e.g., transcript/gene) abundances are correlated</i>
-------------------	--

---

## Description

remove\_redundancy() takes as input a 'tbl' formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>| for correlation method or |<DIMENSION 1>|<DIMENSION 2>|<...>| for reduced\_dimensions method, and returns a 'tbl' with dropped elements (e.g., samples).

## Usage

```
remove_redundancy(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  correlation_threshold = 0.9,
  top = Inf,
  log_transform = FALSE,
  Dim_a_column,
  Dim_b_column
)

## S4 method for signature 'spec_tbl_df'
remove_redundancy(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  correlation_threshold = 0.9,
  top = Inf,
  log_transform = FALSE,
  Dim_a_column = NULL,
  Dim_b_column = NULL
)

## S4 method for signature 'tbl_df'
remove_redundancy(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  correlation_threshold = 0.9,
  top = Inf,
```

```

    log_transform = FALSE,
    Dim_a_column = NULL,
    Dim_b_column = NULL
  )

## S4 method for signature 'tidybulk'
remove_redundancy(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  correlation_threshold = 0.9,
  top = Inf,
  log_transform = FALSE,
  Dim_a_column = NULL,
  Dim_b_column = NULL
)

## S4 method for signature 'SummarizedExperiment'
remove_redundancy(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  correlation_threshold = 0.9,
  top = Inf,
  log_transform = FALSE,
  Dim_a_column = NULL,
  Dim_b_column = NULL
)

## S4 method for signature 'RangedSummarizedExperiment'
remove_redundancy(
  .data,
  .element = NULL,
  .feature = NULL,
  .abundance = NULL,
  method,
  of_samples = TRUE,
  correlation_threshold = 0.9,
  top = Inf,
  log_transform = FALSE,
  Dim_a_column = NULL,
  Dim_b_column = NULL
)

```

### Arguments

`.data` A 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> |

<code>.element</code>	The name of the element column (normally samples).
<code>.feature</code>	The name of the feature column (normally transcripts/genes)
<code>.abundance</code>	The name of the column including the numerical value the clustering is based on (normally transcript abundance)
<code>method</code>	A character string. The cluster algorithm to use, at the moment k-means is the only algorithm included.
<code>of_samples</code>	A boolean. In case the input is a tidybulk object, it indicates Whether the element column will be sample or transcript column
<code>correlation_threshold</code>	A real number between 0 and 1. For correlation based calculation.
<code>top</code>	An integer. How many top genes to select for correlation based method
<code>log_transform</code>	A boolean, whether the value should be log-transformed (e.g., TRUE for RNA sequencing data)
<code>Dim_a_column</code>	A character string. For reduced_dimension based calculation. The column of one principal component
<code>Dim_b_column</code>	A character string. For reduced_dimension based calculation. The column of another principal component

### Value

A tbl object with with dropped redundant elements (e.g., samples).

A tbl object with with dropped redundant elements (e.g., samples).

A tbl object with with dropped redundant elements (e.g., samples).

A tbl object with with dropped redundant elements (e.g., samples).

A ‘SummarizedExperiment’ object

A ‘SummarizedExperiment’ object

### Examples

```
tidybulk::counts_mini %>%
tidybulk(sample, transcript, count) %>%
identify_abundant() %>%
  remove_redundancy(
    .element = sample,
    .feature = transcript,
    .abundance = count,
    method = "correlation"
  )

counts.MDS =
tidybulk::counts_mini %>%
tidybulk(sample, transcript, count) %>%
identify_abundant() %>%
  reduce_dimensions( method="MDS", .dims = 3)

remove_redundancy(
counts.MDS,
Dim_a_column = `Dim1`,
```



```
Dim_b_column = `Dim2`,
  .element = sample,
  method = "reduced_dimensions"
)
```

---

 rename

*Rename columns*


---

## Description

Rename individual variables using ‘new\_name = old\_name’ syntax.

## Usage

```
rename(.data, ...)
```

## Arguments

<code>.data</code>	A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from <code>dbplyr</code> or <code>dtplyr</code> ). See <i>*Methods*</i> , below, for more details.
<code>...</code>	<[‘tidy-select’][ <code>dplyr_tidy_select</code> ]> Use ‘new_name = old_name’ to rename selected variables.

## Value

An object of the same type as ‘.data’. \* Rows are not affected. \* Column names are changed; column order is preserved \* Data frame attributes are preserved. \* Groups are updated to reflect new names.

## Scoped selection and renaming

Use the three scoped variants (`[rename_all()]`, `[rename_if()]`, `[rename_at()]`) to renaming a set of variables with a function.

## Methods

This function is a *\*\*generic\*\**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages:

## See Also

Other single table verbs: `arrange()`, `filter()`, `mutate()`, `summarise()`

## Examples

```
`%>%` = magrittr::`%>%`
iris <- as_tibble(iris) # so it prints a little nicer
rename(iris, petal_length = Petal.Length)
```

---

right_join	<i>Right join datasets</i>
------------	----------------------------

---

### Description

Right join datasets

### Usage

```
right_join(x, y, by = NULL, copy = FALSE, suffix = c(".x", ".y"), ...)
```

### Arguments

x	tbls to join. (See dplyr)
y	tbls to join. (See dplyr)
by	A character vector of variables to join by. (See dplyr)
copy	If x and y are not from the same data source, and copy is TRUE, then y will be copied into the same src as x. (See dplyr)
suffix	If there are non-joined duplicate variables in x and y, these suffixes will be added to the output to disambiguate them. Should be a character vector of length 2. (See dplyr)
...	Data frames to combine (See dplyr)

### Value

A tt object

### Examples

```
`%>%` = magrittr::`%>%`
annotation = tidybulk::counts %>% distinct(sample) %>% mutate(source = "AU")
tidybulk::counts %>% right_join(annotation)
```

---

rotate_dimensions	<i>Rotate two dimensions (e.g., principal components) of an arbitrary angle</i>
-------------------	---

---

### Description

rotate\_dimensions() takes as input a 'tbl' formatted as | <DIMENSION 1> | <DIMENSION 2> | <...> | and calculates the rotated dimensional space of the transcript abundance.

**Usage**

```
rotate_dimensions(  
  .data,  
  dimension_1_column,  
  dimension_2_column,  
  rotation_degrees,  
  .element = NULL,  
  of_samples = TRUE,  
  dimension_1_column_rotated = NULL,  
  dimension_2_column_rotated = NULL,  
  action = "add"  
)  
  
## S4 method for signature 'spec_tbl_df'  
rotate_dimensions(  
  .data,  
  dimension_1_column,  
  dimension_2_column,  
  rotation_degrees,  
  .element = NULL,  
  of_samples = TRUE,  
  dimension_1_column_rotated = NULL,  
  dimension_2_column_rotated = NULL,  
  action = "add"  
)  
  
## S4 method for signature 'tbl_df'  
rotate_dimensions(  
  .data,  
  dimension_1_column,  
  dimension_2_column,  
  rotation_degrees,  
  .element = NULL,  
  of_samples = TRUE,  
  dimension_1_column_rotated = NULL,  
  dimension_2_column_rotated = NULL,  
  action = "add"  
)  
  
## S4 method for signature 'tidybulk'  
rotate_dimensions(  
  .data,  
  dimension_1_column,  
  dimension_2_column,  
  rotation_degrees,  
  .element = NULL,  
  of_samples = TRUE,  
  dimension_1_column_rotated = NULL,  
  dimension_2_column_rotated = NULL,  
  action = "add"  
)
```

```
## S4 method for signature 'SummarizedExperiment'
rotate_dimensions(
  .data,
  dimension_1_column,
  dimension_2_column,
  rotation_degrees,
  .element = NULL,
  of_samples = TRUE,
  dimension_1_column_rotated = NULL,
  dimension_2_column_rotated = NULL,
  action = "add"
)

## S4 method for signature 'RangedSummarizedExperiment'
rotate_dimensions(
  .data,
  dimension_1_column,
  dimension_2_column,
  rotation_degrees,
  .element = NULL,
  of_samples = TRUE,
  dimension_1_column_rotated = NULL,
  dimension_2_column_rotated = NULL,
  action = "add"
)
```

## Arguments

.data	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
dimension_1_column	A character string. The column of the dimension 1
dimension_2_column	A character string. The column of the dimension 2
rotation_degrees	A real number between 0 and 360
.element	The name of the element column (normally samples).
of_samples	A boolean. In case the input is a tidybulk object, it indicates Whether the element column will be sample or transcript column
dimension_1_column_rotated	A character string. The column of the rotated dimension 1 (optional)
dimension_2_column_rotated	A character string. The column of the rotated dimension 2 (optional)
action	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).

## Details

### [Maturing]

This function to rotate two dimensions such as the reduced dimensions.

Underlying custom method:  $\text{rotation} = \text{function}(m, d) \text{ // } r = \text{the angle // } m \text{ data matrix } r = d * \pi / 180 \text{ ((dplyr::bind\_rows( c('1' = cos(r), '2' = -sin(r)), c('1' = sin(r), '2' = cos(r)) )$

**Value**

A tbl object with additional columns for the reduced dimensions. additional columns for the rotated dimensions. The rotated dimensions will be added to the original data set as ‘<NAME OF DIMENSION> rotated <ANGLE>’ by default, or as specified in the input arguments.

A tbl object with additional columns for the reduced dimensions. additional columns for the rotated dimensions. The rotated dimensions will be added to the original data set as ‘<NAME OF DIMENSION> rotated <ANGLE>’ by default, or as specified in the input arguments.

A tbl object with additional columns for the reduced dimensions. additional columns for the rotated dimensions. The rotated dimensions will be added to the original data set as ‘<NAME OF DIMENSION> rotated <ANGLE>’ by default, or as specified in the input arguments.

A tbl object with additional columns for the reduced dimensions. additional columns for the rotated dimensions. The rotated dimensions will be added to the original data set as ‘<NAME OF DIMENSION> rotated <ANGLE>’ by default, or as specified in the input arguments.

A ‘SummarizedExperiment’ object

A ‘SummarizedExperiment’ object

**Examples**

```
counts.MDS =
  tidybulk::counts_mini %>%
  tidybulk(sample, transcript, count) %>%
  identify_abundant() %>%
  reduce_dimensions( method="MDS", .dims = 3)

counts.MDS.rotated = rotate_dimensions(counts.MDS, `Dim1`, `Dim2`, rotation_degrees = 45, .element = sample)
```

---

rowwise

*Group input by rows*


---

**Description**

See [this repository](<https://github.com/jennybc/row-oriented-workflows>) for alternative ways to perform row-wise operations.

**Usage**

```
rowwise(.data)
```

**Arguments**

.data                      Input data frame.

**Details**

‘rowwise()’ is used for the results of [do()] when you create list-variables. It is also useful to support arbitrary complex operations that need to be applied to each row.

Currently, rowwise grouping only works with data frames. Its main impact is to allow you to work with list-variables in [summarise()] and [mutate()] without having to use [[1]]. This makes ‘summarise()’ on a rowwise tbl effectively equivalent to [plyr::ldply()].

**Value**

A 'tbl'

A 'tbl'

**Examples**

```
`%>%` = magrittr::`%>%`
df <- expand.grid(x = 1:3, y = 3:1)
df_done <- df %>% rowwise() %>% do(i = seq(.$x, .$y))
```

scale\_abundance

*Scale the counts of transcripts/genes***Description**

scale\_abundance() takes as input a 'tbl' formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>| and Scales transcript abundance compensating for sequencing depth (e.g., with TMM algorithm, Robinson and Oshlack doi.org/10.1186/gb-2010-11-3-r25).

**Usage**

```
scale_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "TMM",
  reference_sample = NULL,
  action = "add",
  reference_selection_function = NULL
)

## S4 method for signature 'spec_tbl_df'
scale_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "TMM",
  reference_sample = NULL,
  action = "add",
  reference_selection_function = NULL
)

## S4 method for signature 'tbl_df'
scale_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
```

```

    method = "TMM",
    reference_sample = NULL,
    action = "add",
    reference_selection_function = NULL
)

## S4 method for signature 'tidybulk'
scale_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "TMM",
  reference_sample = NULL,
  action = "add",
  reference_selection_function = NULL
)

## S4 method for signature 'SummarizedExperiment'
scale_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "TMM",
  reference_sample = NULL,
  action = "add",
  reference_selection_function = NULL
)

## S4 method for signature 'RangedSummarizedExperiment'
scale_abundance(
  .data,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "TMM",
  reference_sample = NULL,
  action = "add",
  reference_selection_function = NULL
)

```

### Arguments

.data	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
.sample	The name of the sample column
.transcript	The name of the transcript/gene column
.abundance	The name of the transcript/gene abundance column
method	A character string. The scaling method passed to the back-end function (i.e., edgeR::calcNormFactors; "TMM", "TMMwsp", "RLE", "upperquartile")

reference_sample	A character string. The name of the reference sample. If NULL the sample with highest total read count will be selected as reference.
action	A character string between "add" (default) and "only". "add" joins the new information to the input tbl (default), "only" return a non-redundant tbl with the just new information.
reference_selection_function	DEPRECATED. please use reference_sample.

Details

[Maturing]

Scales transcript abundance compensating for sequencing depth (e.g., with TMM algorithm, Robinson and Oshlack doi.org/10.1186/gb-2010-11-3-r25). Lowly transcribed transcripts/genes (defined with minimum\_counts and minimum\_proportion parameters) are filtered out from the scaling procedure. The scaling inference is then applied back to all unfiltered data.

Underlying method `edgeR::calcNormFactors(.data, method = c("TMM", "TMMwsp", "RLE", "upperquartile"))`

Value

- A tbl object with additional columns with scaled data as ‘<NAME OF COUNT COLUMN>\_scaled’
- A tbl object with additional columns with scaled data as ‘<NAME OF COUNT COLUMN>\_scaled’
- A tbl object with additional columns with scaled data as ‘<NAME OF COUNT COLUMN>\_scaled’
- A tbl object with additional columns with scaled data as ‘<NAME OF COUNT COLUMN>\_scaled’
- A ‘SummarizedExperiment’ object
- A ‘SummarizedExperiment’ object

Examples

```
tidybulk::counts_mini %>%
  tidybulk(sample, transcript, count) %>%
  identify_abundant() %>%
  scale_abundance()
```

---

se	<i>SummarizedExperiment</i>
----	-----------------------------

---

Description

SummarizedExperiment

Usage

se



**Format**

An object of class RangedSummarizedExperiment with 100 rows and 8 columns.

---

se_mini	<i>SummarizedExperiment mini for vignette</i>
---------	---

---

**Description**

SummarizedExperiment mini for vignette

**Usage**

se\_mini

**Format**

An object of class SummarizedExperiment with 527 rows and 5 columns.

---

summarise	<i>Summarise each group to fewer rows</i>
-----------	---

---

**Description**

‘summarise()’ creates a new data frame. It will have one (or more) rows for each combination of grouping variables; if there are no grouping variables, the output will have a single row summarising all observations in the input. It will contain one column for each grouping variable and one column for each of the summary statistics that you have specified.

‘summarise()’ and ‘summarize()’ are synonyms.

**Usage**

summarise(.data, ...)

**Arguments**

- |       |   |
|-------|---|
| .data | A data frame, data frame extension (e.g. a tibble), or a lazy data frame (e.g. from dbplyr or dtplyr). See <i>*Methods*</i> , below, for more details.  |
| ...   | <[‘tidy-eval’][dplyr_tidy_eval]> Name-value pairs of summary functions. The name will be the name of the variable in the result.<br>The value can be:<br>* A vector of length 1, e.g. ‘min(x)’, ‘n()’, or ‘sum(is.na(y))’. * A vector of length ‘n’, e.g. ‘quantile()’. * A data frame, to add multiple columns from a single expression. |

**Value**

An object `_usually_` of the same type as `‘.data’`.

\* The rows come from the underlying `‘group_keys()’`. \* The columns are a combination of the grouping keys and the summary expressions that you provide. \* If `‘x’` is grouped by more than one variable, the output will be another `[grouped_df]` with the right-most group removed. \* If `‘x’` is grouped by one variable, or is not grouped, the output will be a `[tibble]`. \* Data frame attributes are **not** preserved, because `‘summarise()’` fundamentally creates a new data frame.

**Useful functions**

\* Center: `[mean()]`, `[median()]` \* Spread: `[sd()]`, `[IQR()]`, `[mad()]` \* Range: `[min()]`, `[max()]`, `[quantile()]` \* Position: `[first()]`, `[last()]`, `[nth()]`, \* Count: `[n()]`, `[n_distinct()]` \* Logical: `[any()]`, `[all()]`

**Backend variations**

The data frame backend supports creating a variable and using it in the same summary. This means that previously created summary variables can be further transformed or combined within the summary, as in `[mutate()]`. However, it also means that summary variables with the same names as previous variables overwrite them, making those variables unavailable to later summary variables.

This behaviour may not be supported in other backends. To avoid unexpected results, consider using new names for your summary variables, especially when creating multiple summaries.

**Methods**

This function is a **generic**, which means that packages can provide implementations (methods) for other classes. See the documentation of individual methods for extra arguments and differences in behaviour.

The following methods are currently available in loaded packages:

**See Also**

Other single table verbs: `arrange()`, `filter()`, `mutate()`, `rename()`

**Examples**

```
`%>%` = magrittr::`%>%`
# A summary applied to ungrouped tbl returns a single row
mtcars %>%
  summarise(mean = mean(disp))
```

---

symbol\_to\_entrez

*Get ENTREZ id from gene SYMBOL*

---

**Description**

Get ENTREZ id from gene SYMBOL

**Usage**

```
symbol_to_entrez(.data, .transcript = NULL, .sample = NULL)
```

**Arguments**

.data            A tt or tbl object.

.transcript     A character. The name of the gene symbol column.

.sample         The name of the sample column

**Value**

A tbl

**Examples**

```
symbol_to_entrez(tidybulk::counts_mini, .transcript = transcript, .sample = sample)
```

---

test_deseq2_df	<i>SummarizedExperiment mini for vignette</i>
----------------	---

---

**Description**

SummarizedExperiment mini for vignette

**Usage**

```
test_deseq2_df
```

**Format**

An object of class DESeqDataSet with 9921 rows and 7 columns.

---

test_differential_abundance	<i>Perform differential transcription testing using edgeR QLT, edgeR LR, limma-voom or DESeq2</i>
-----------------------------	---

---

**Description**

test\_differential\_abundance() takes as input a ‘tbl’ formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a ‘tbl’ with additional columns for the statistics from the hypothesis test.

**Usage**

```
test_differential_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  .contrasts = NULL,  
  method = "edgeR_quasi_likelihood",  
  scaling_method = "TMM",  
  omit_contrast_in_colnames = FALSE,  
  prefix = "",  
  action = "add",  
  significance_threshold = NULL,  
  fill_missing_values = NULL  
)  
  
## S4 method for signature 'spec_tbl_df'  
test_differential_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  .contrasts = NULL,  
  method = "edgeR_quasi_likelihood",  
  scaling_method = "TMM",  
  omit_contrast_in_colnames = FALSE,  
  prefix = "",  
  action = "add",  
  significance_threshold = NULL,  
  fill_missing_values = NULL  
)  
  
## S4 method for signature 'tbl_df'  
test_differential_abundance(  
  .data,  
  .formula,  
  .sample = NULL,  
  .transcript = NULL,  
  .abundance = NULL,  
  .contrasts = NULL,  
  method = "edgeR_quasi_likelihood",  
  scaling_method = "TMM",  
  omit_contrast_in_colnames = FALSE,  
  prefix = "",  
  action = "add",  
  significance_threshold = NULL,  
  fill_missing_values = NULL  
)  
  
## S4 method for signature 'tidybulk'
```

```

test_differential_abundance(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  .contrasts = NULL,
  method = "edgeR_quasi_likelihood",
  scaling_method = "TMM",
  omit_contrast_in_colnames = FALSE,
  prefix = "",
  action = "add",
  significance_threshold = NULL,
  fill_missing_values = NULL
)

## S4 method for signature 'SummarizedExperiment'
test_differential_abundance(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  .contrasts = NULL,
  method = "edgeR_quasi_likelihood",
  scaling_method = "TMM",
  omit_contrast_in_colnames = FALSE,
  prefix = "",
  action = "add",
  significance_threshold = NULL,
  fill_missing_values = NULL
)

## S4 method for signature 'RangedSummarizedExperiment'
test_differential_abundance(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  .contrasts = NULL,
  method = "edgeR_quasi_likelihood",
  scaling_method = "TMM",
  omit_contrast_in_colnames = FALSE,
  prefix = "",
  action = "add",
  significance_threshold = NULL,
  fill_missing_values = NULL
)

```

### Arguments

`.data` A 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> |

<code>.formula</code>	A formula with no response variable, representing the desired linear model
<code>.sample</code>	The name of the sample column
<code>.transcript</code>	The name of the transcript/gene column
<code>.abundance</code>	The name of the transcript/gene abundance column
<code>.contrasts</code>	This parameter takes the shape of the contrast parameter of the method of choice. For edgeR and limma-voom is a character vector. For DESeq2 is a list including a character vectors of length three. If contrasts are not present the first covariate is the one the model is tested against (e.g., <code>~ factor_of_interest</code> )
<code>method</code>	A string character. Either "edgeR_quasi_likelihood" (i.e., QLF), "edgeR_likelihood_ratio" (i.e., LRT), "DESeq2", "limma_voom"
<code>scaling_method</code>	A character string. The scaling method passed to the back-end function (i.e., <code>edgeR::calcNormFactors</code> ; "TMM", "TMMwsp", "RLE", "upperquartile")
<code>omit_contrast_in_colnames</code>	If just one contrast is specified you can choose to omit the contrast label in the colnames.
<code>prefix</code>	A character string. The prefix you would like to add to the result columns. It is useful if you want to compare several methods.
<code>action</code>	A character string. Whether to join the new information to the input tbl (add), or just get the non-redundant tbl with the new information (get).
<code>significance_threshold</code>	A real between 0 and 1 (usually 0.05).
<code>fill_missing_values</code>	A boolean. Whether to fill missing sample/transcript values with the median of the transcript. This is rarely needed.

## Details

### [Maturing]

This function provides the option to use edgeR <https://doi.org/10.1093/bioinformatics/btp616>, limma-voom <https://doi.org/10.1186/gb-2014-15-2-r29>, or DESeq2 <https://doi.org/10.1186/s13059-014-0550-8> to perform the testing. All methods use raw counts, irrespective of if `scale_abundance` or `adjust_abundance` have been calculated, therefore it is essential to add covariates such as batch effects (if applicable) in the formula.

Underlying method for edgeR framework: `.data`

```
# Filter keep_abundant( factor_of_interest = !(as.symbol(parse_formula(.formula)[1])), minimum_counts
= minimum_counts, minimum_proportion = minimum_proportion )
```

```
# Format select(!.transcript,!sample,!abundance) spread(!.sample,!abundance) as_matrix(rownames
= !.transcript)
```

```
# edgeR edgeR::DGEList(counts = .) edgeR::calcNormFactors(method = scaling_method) edgeR::estimateDisp(design)
```

```
# Fit edgeR::glmQLFit(design) edgeR::glmQLFTest(coef = 2, contrast = my_contrasts) // or glmLRT
according to choice
```

Underlying method for DESeq2 framework: `keep_abundant( factor_of_interest = !(as.symbol(parse_formula(.formula)[1]))`  
`minimum_counts = minimum_counts, minimum_proportion = minimum_proportion )`

```
# DESeq2 DESeq2::DESeqDataSet( design = .formula) DESeq2::DESeq() DESeq2::results()
```

**Value**

A 'tbl' with additional columns for the statistics from the test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'SummarizedExperiment' object

A 'SummarizedExperiment' object

**Examples**

```
tidybulk::counts_mini %>%
tidybulk(sample, transcript, count) %>%
identify_abundant() %>%
test_differential_abundance( ~ condition )

# The function `test_differential_abundance` operates with contrasts too

tidybulk::counts_mini %>%
tidybulk(sample, transcript, count) %>%
identify_abundant() %>%
test_differential_abundance(
  ~ 0 + condition,
  .contrasts = c( "conditionTRUE - conditionFALSE" )
)
```

---

test\_differential\_cellularity

*Add differential tissue composition information to a tbl*

---

**Description**

test\_differential\_cellularity() takes as input a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> | and returns a 'tbl' with additional columns for the statistics from the hypothesis test.

**Usage**

```
test_differential_cellularity(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
```

```

    method = "cibersort",
    reference = X_cibersort,
    significance_threshold = 0.05,
    ...
)

## S4 method for signature 'spec_tbl_df'
test_differential_cellularity(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "cibersort",
  reference = X_cibersort,
  significance_threshold = 0.05,
  ...
)

## S4 method for signature 'tbl_df'
test_differential_cellularity(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "cibersort",
  reference = X_cibersort,
  significance_threshold = 0.05,
  ...
)

## S4 method for signature 'tidybulk'
test_differential_cellularity(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "cibersort",
  reference = X_cibersort,
  significance_threshold = 0.05,
  ...
)

## S4 method for signature 'SummarizedExperiment'
test_differential_cellularity(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,

```



```

    method = "cibersort",
    reference = X_cibersort,
    significance_threshold = 0.05,
    ...
)

## S4 method for signature 'RangedSummarizedExperiment'
test_differential_cellularity(
  .data,
  .formula,
  .sample = NULL,
  .transcript = NULL,
  .abundance = NULL,
  method = "cibersort",
  reference = X_cibersort,
  significance_threshold = 0.05,
  ...
)

```

## Arguments

.data	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
.formula	A formula with no response variable, representing the desired linear model. If censored regression is desired (coxph) the formula should be of the form \"survival::Surv(y, dead) ~ x\"
.sample	The name of the sample column
.transcript	The name of the transcript/gene column
.abundance	The name of the transcript/gene abundance column
method	A string character. Either \"cibersort\" or \"llsr\"
reference	A data frame. The transcript/cell_type data frame of integer transcript abundance
significance_threshold	A real between 0 and 1 (usually 0.05).
...	Further parameters passed to the method deconvolve_cellularity

## Details

### [Maturing]

This routine applies a deconvolution method (e.g., Cibersort; DOI: 10.1038/nmeth.3337) and passes the proportions inferred into a generalised linear model (DOI:dx.doi.org/10.1007/s11749-010-0189-z) or a cox regression model (ISBN: 978-1-4757-3294-8)

Underlying method for the generalised linear model: `data deconvolve_cellularity( !!.sample, !!.transcript, !!.abundance, method=method, reference = reference, action="get", ... ) [..] betareg::betareg(.my_formula, .)`

Underlying method for the cox regression: `data deconvolve_cellularity( !!.sample, !!.transcript, !!.abundance, method=method, reference = reference, action="get", ... ) [..] mutate(.proportion_0_corrected = .proportion_0_corrected survival::coxph(.my_formula, .)`

**Value**

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'tbl' with additional columns for the statistics from the hypothesis test (e.g., log fold change, p-value and false discovery rate).

A 'SummarizedExperiment' object

A 'SummarizedExperiment' object

**Examples**

```
test_differential_cellularity(
  tidybulk::counts_mini,
  ~ condition,
  sample,
  transcript,
  count,
  cores = 1
)
```

---

test\_gene\_enrichment    *analyse gene enrichment with EGSEA*

---

**Description**

test\_gene\_enrichment() takes as input a 'tbl' formatted as | <SAMPLE> | <ENSEMBL\_ID> | <COUNT> | <...> | and returns a 'tbl' with the additional transcript symbol column

**Usage**

```
test_gene_enrichment(
  .data,
  .formula,
  .sample = NULL,
  .entrez,
  .abundance = NULL,
  .contrasts = NULL,
  method = c("camera", "roast", "safe", "gage", "padog", "globaltest", "ora"),
  species,
  cores = 10
)
```

```
## S4 method for signature 'spec_tbl_df'
test_gene_enrichment(
  .data,
  .formula,
  .sample = NULL,
  .entrez,
  .abundance = NULL,
  .contrasts = NULL,
  method = c("camera", "roast", "safe", "gage", "padog", "globaltest", "ora"),
  species,
  cores = 10
)

## S4 method for signature 'tbl_df'
test_gene_enrichment(
  .data,
  .formula,
  .sample = NULL,
  .entrez,
  .abundance = NULL,
  .contrasts = NULL,
  method = c("camera", "roast", "safe", "gage", "padog", "globaltest", "ora"),
  species,
  cores = 10
)

## S4 method for signature 'tidybulk'
test_gene_enrichment(
  .data,
  .formula,
  .sample = NULL,
  .entrez,
  .abundance = NULL,
  .contrasts = NULL,
  method = c("camera", "roast", "safe", "gage", "padog", "globaltest", "ora"),
  species,
  cores = 10
)
```

### Arguments

.data	A 'tbl' formatted as  <SAMPLE> <TRANSCRIPT> <COUNT> <...>
.formula	A formula with no response variable, representing the desired linear model
.sample	The name of the sample column
.entrez	The ENTREZ ID of the transcripts/genes
.abundance	The name of the transcript/gene abundance column
.contrasts	= NULL,
method	A character vector. The methods to be included in the ensembl. Type EGSEA::egsea.base() to see the supported GSE methods.
species	A character. For example, human or mouse
cores	An integer. The number of cores available

## Details

### [Maturing]

This wrapper execute ensemble gene enrichment analyses of the dataset using EGSEA (DOI:0.12688/f1000research.1254

```
dge = data_keep_abundant( factor_of_interest = !!as.symbol(parse_formula(.formula)[[1]]), !!sample, !!entrez, !!abundance )
```

```
# Make sure transcript names are adjacent [...] as_matrix(rownames = !!entrez) edgeR::DGEList(counts = .)
```

```
idx = buildIdx(entrezIDs = rownames(dge), species = species)
```

```
dge
```

```
# Calculate weights limma::voom(design, plot = FALSE)
```

```
# Execute EGSEA egsea( contrasts = my_contrasts, baseGSEAs = method, sort.by = "med.rank", num.threads = cores, report = FALSE )
```

## Value

A 'tbl' object

A 'tbl' object

A 'tbl' object

A 'tbl' object

## Examples

```
## Not run:
```

```
df_entrez = symbol_to_entrez(tidybulk::counts_mini, .transcript = transcript, .sample = sample)
df_entrez = aggregate_duplicates(df_entrez, aggregation_function = sum, .sample = sample, .transcript = entrez)
```

```
library("EGSEA")
```

```
test_gene_enrichment(
  df_entrez,
  ~ condition,
  .sample = sample,
  .entrez = entrez,
  .abundance = count,
  method = c("roast", "safe", "gage", "padog", "globaltest", "ora"),
  species="human",
  cores = 2
)
```

```
## End(Not run)
```

---

`test_gene_overrepresentation`*analyse gene over-representation with GSEA*

---

## Description

`test_gene_overrepresentation()` takes as input a 'tbl' formatted as | <SAMPLE> | <ENSEMBL\_ID> | <COUNT> | <...> | and returns a 'tbl' with the GSEA statistics

## Usage

```
test_gene_overrepresentation(  
  .data,  
  .sample = NULL,  
  .entrez,  
  .do_test,  
  species,  
  gene_set = NULL  
)  
  
## S4 method for signature 'spec_tbl_df'  
test_gene_overrepresentation(  
  .data,  
  .sample = NULL,  
  .entrez,  
  .do_test,  
  species,  
  gene_set = NULL  
)  
  
## S4 method for signature 'tbl_df'  
test_gene_overrepresentation(  
  .data,  
  .sample = NULL,  
  .entrez,  
  .do_test,  
  species,  
  gene_set = NULL  
)  
  
## S4 method for signature 'tidybulk'  
test_gene_overrepresentation(  
  .data,  
  .sample = NULL,  
  .entrez,  
  .do_test,  
  species,  
  gene_set = NULL  
)
```

**Arguments**

<code>.data</code>	A 'tbl' formatted as   <SAMPLE>   <TRANSCRIPT>   <COUNT>   <...>
<code>.sample</code>	The name of the sample column
<code>.entrez</code>	The ENTREZ ID of the transcripts/genes
<code>.do_test</code>	A boolean column name symbol. It indicates the transcript to check
<code>species</code>	A character. For example, human or mouse. MSigDB uses the latin species names (e.g., <code>"Mus musculus"</code> , <code>"Homo sapiens"</code> )
<code>gene_set</code>	A character vector. The subset of MSigDB datasets you want to test against (e.g. <code>"C2"</code> ). If NULL all gene sets are used (suggested). This argument was added to avoid time overflow of the examples.

**Details****[Maturing]**

This wrapper execute gene enrichment analyses of the dataset using a list of transcripts and GSEA. This wrapper uses clusterProfiler (DOI: [doi.org/10.1089/omi.2011.0118](https://doi.org/10.1089/omi.2011.0118)) on the back-end.

Undelying method: `msigdbr::msigdbr(species = species) nest(data = -gs_cat) mutate(test = map(data, ~ clusterProfiler::enricher( my_entrez_rank, TERM2GENE=.x pvalueCutoff = 1 ) ))`

**Value**

A 'tbl' object  
 A 'tbl' object  
 A 'tbl' object  
 A 'tbl' object

**Examples**

```
df_entrez = symbol_to_entrez(tidybulk::counts_mini, .transcript = transcript, .sample = sample)
df_entrez = aggregate_duplicates(df_entrez, aggregation_function = sum, .sample = sample, .transcript = entrez)
df_entrez = mutate(df_entrez, do_test = transcript %in% c("TNFRSF4", "PLCH2", "PADI4", "PAX7"))

test_gene_overrepresentation(
  df_entrez,
  .sample = sample,
  .entrez = entrez,
  .do_test = do_test,
  species="Homo sapiens",
  gene_set=c("C2")
)
```

---

tidybulk	<i>Creates a 'tt' object from a 'tbl' or 'SummarizedExperiment' object</i>
----------	--

---

## Description

tidybulk() creates a 'tt' object from a 'tbl' formatted as |<SAMPLE>|<TRANSCRIPT>|<COUNT>|<...>|

## Usage

```
tidybulk(.data, .sample, .transcript, .abundance, .abundance_scaled = NULL)

## S4 method for signature 'spec_tbl_df'
tidybulk(.data, .sample, .transcript, .abundance, .abundance_scaled = NULL)

## S4 method for signature 'tbl_df'
tidybulk(.data, .sample, .transcript, .abundance, .abundance_scaled = NULL)

## S4 method for signature 'SummarizedExperiment'
tidybulk(.data, .sample, .transcript, .abundance, .abundance_scaled = NULL)

## S4 method for signature 'RangedSummarizedExperiment'
tidybulk(.data, .sample, .transcript, .abundance, .abundance_scaled = NULL)
```

## Arguments

.data	A 'tbl' formatted as  <SAMPLE> <TRANSCRIPT> <COUNT> <...>
.sample	The name of the sample column
.transcript	The name of the transcript/gene column
.abundance	The name of the transcript/gene abundance column
.abundance_scaled	The name of the transcript/gene scaled abundance column

## Details

### [Maturing]

This function creates a tidybulk object and is useful if you want to avoid to specify .sample, .transcript and .abundance arguments all the times. The tidybulk object have an attribute called internals where these three arguments are stored as metadata. They can be extracted as attr(<object>, "internals").

## Value

A 'tidybulk' object  
 A 'tidybulk' object  
 A 'tidybulk' object  
 A 'tidybulk' object  
 A 'tidybulk' object

## Examples

```
my_tt = tidybulk(tidybulk::counts_mini, sample, transcript, count)
```

---

tidybulk_SAM_BAM	<i>Creates a 'tt' object from a list of file names of BAM/SAM</i>
------------------	---

---

## Description

tidybulk\_SAM\_BAM() creates a 'tt' object from a 'tbl' formatted as | <SAMPLE> | <TRANSCRIPT> | <COUNT> | <...> |

## Usage

```
tidybulk_SAM_BAM(file_names, genome = "hg38", ...)
```

```
## S4 method for signature 'character,character'
tidybulk_SAM_BAM(file_names, genome = "hg38", ...)
```

## Arguments

file_names	A character vector
genome	A character string
...	Further parameters passed to the function Rsubread::featureCounts

## Details

### [Maturing]

This function is based on FeatureCounts package (DOI: 10.1093/bioinformatics/btt656). This function creates a tidybulk object and is useful if you want to avoid to specify .sample, .transcript and .abundance arguments all the times. The tidybulk object have an attribute called internals where these three arguments are stored as metadata. They can be extracted as attr(<object>, "internals").

Underlying core function Rsubread::featureCounts(annot.inbuilt = genome,nthreads = n\_cores, ...)

## Value

A 'tidybulk' object

A 'tidybulk' object



---

`unnest`*unnest*

---

## Description

`unnest``nest`

## Usage

```
unnest(  
  .data,  
  cols,  
  ...,  
  keep_empty = FALSE,  
  ptype = NULL,  
  names_sep = NULL,  
  names_repair = "check_unique"  
)  
  
## Default S3 method:  
unnest(  
  .data,  
  cols,  
  ...,  
  keep_empty = FALSE,  
  ptype = NULL,  
  names_sep = NULL,  
  names_repair = "check_unique"  
)  
  
## S3 method for class 'nested_tidybulk'  
unnest(  
  .data,  
  cols,  
  ...,  
  keep_empty = FALSE,  
  ptype = NULL,  
  names_sep = NULL,  
  names_repair = "check_unique"  
)  
  
nest(.data, ...)  
  
## Default S3 method:  
nest(.data, ...)  
  
## S3 method for class 'tidybulk'  
nest(.data, ...)
```

**Arguments**

<code>.data</code>	A tbl. (See <code>tidyr</code> )
<code>cols</code>	<[‘tidy-select’][ <code>tidyr::tidy_select</code> ]> Columns to unnest. If you ‘ <code>unnest()</code> ’ multiple columns, parallel entries must be of compatible sizes, i.e. they’re either equal or length 1 (following the standard tidyverse recycling rules).
<code>...</code>	Name-variable pairs of the form <code>new_col = c(col1, col2, col3)</code> (See <code>tidyr</code> )
<code>keep_empty</code>	See <code>tidyr::unnest</code>
<code>ptype</code>	See <code>tidyr::unnest</code>
<code>names_sep</code>	If ‘NULL’, the default, the names will be left as is. In ‘ <code>nest()</code> ’, inner names will come from the former outer names; in ‘ <code>unnest()</code> ’, the new outer names will come from the inner names.  If a string, the inner and outer names will be used together. In ‘ <code>nest()</code> ’, the names of the new outer columns will be formed by pasting together the outer and the inner column names, separated by ‘ <code>names_sep</code> ’. In ‘ <code>unnest()</code> ’, the new inner names will have the outer names (+ ‘ <code>names_sep</code> ’) automatically stripped. This makes ‘ <code>names_sep</code> ’ roughly symmetric between nesting and unnesting.
<code>names_repair</code>	See <code>tidyr::unnest</code>

**Value**

A tt object

A tt object

**Examples**

```
library(dplyr)

nest(tidybulk(tidybulk::counts_mini, sample, transcript, count), data = -transcript) %>%
unnest(data)

nest(tidybulk(tidybulk::counts_mini, sample, transcript, count), data = -transcript)
```

---

vignette\_manuscript\_signature\_boxplot

*Needed for vignette vignette\_manuscript\_signature\_boxplot*


---

**Description**

Needed for vignette vignette\_manuscript\_signature\_boxplot

**Usage**

vignette\_manuscript\_signature\_boxplot

**Format**An object of class `tbl_df` (inherits from `tbl`, `data.frame`) with 899 rows and 12 columns.

---

`vignette_manuscript_signature_tsne`*Needed for vignette vignette\_manuscript\_signature\_tsne*

---

**Description**

Needed for vignette vignette\_manuscript\_signature\_tsne

**Usage**

```
vignette_manuscript_signature_tsne
```

**Format**

An object of class `spec_tbl_df` (inherits from `tbl_df`, `tbl`, `data.frame`) with 283 rows and 10 columns.

---

`vignette_manuscript_signature_tsne2`*Needed for vignette vignette\_manuscript\_signature\_tsne2*

---

**Description**

Needed for vignette vignette\_manuscript\_signature\_tsne2

**Usage**

```
vignette_manuscript_signature_tsne2
```

**Format**

An object of class `tbl_df` (inherits from `tbl`, `data.frame`) with 283 rows and 9 columns.

---

`X_cibersort`*Cibersort reference*

---

**Description**

Cibersort reference

**Usage**

```
X_cibersort
```

**Format**

An object of class `data.frame` with 547 rows and 22 columns.

# Index

## \* datasets

breast\_tcga\_mini, 11  
 counts, 14  
 counts\_ensembl, 15  
 counts\_mini, 15  
 ensembl\_symbol\_mapping, 19  
 flybaseIDs, 23  
 se, 56  
 se\_mini, 57  
 test\_deseq2\_df, 59  
 vignette\_manuscript\_signature\_boxplot, 74  
 vignette\_manuscript\_signature\_tsne, 75  
 vignette\_manuscript\_signature\_tsne2, 75  
 X\_cibersort, 75

## \* grouping functions

group\_by, 25

## \* single table verbs

arrange, 8  
 filter, 22  
 mutate, 38  
 rename, 49  
 summarise, 57

adjust\_abundance, 4

adjust\_abundance, RangedSummarizedExperiment-method  
 (adjust\_abundance), 4

adjust\_abundance, spec\_tbl\_df-method  
 (adjust\_abundance), 4

adjust\_abundance, SummarizedExperiment-method  
 (adjust\_abundance), 4

adjust\_abundance, tbl\_df-method  
 (adjust\_abundance), 4

adjust\_abundance, tidybulk-method  
 (adjust\_abundance), 4

aggregate\_duplicates, 6

aggregate\_duplicates, RangedSummarizedExperiment-method  
 (aggregate\_duplicates), 6

aggregate\_duplicates, spec\_tbl\_df-method  
 (aggregate\_duplicates), 6

aggregate\_duplicates, SummarizedExperiment-method  
 (aggregate\_duplicates), 6

aggregate\_duplicates, tbl\_df-method  
 (aggregate\_duplicates), 6

aggregate\_duplicates, tidybulk-method  
 (aggregate\_duplicates), 6

arrange, 8, 23, 40, 49, 58

as\_matrix, 10

bind, 10

bind\_cols (arrange), 8

bind\_rows (arrange), 8

breast\_tcga\_mini, 11

cluster\_elements, 12

cluster\_elements, RangedSummarizedExperiment-method  
 (cluster\_elements), 12

cluster\_elements, spec\_tbl\_df-method  
 (cluster\_elements), 12

cluster\_elements, SummarizedExperiment-method  
 (cluster\_elements), 12

cluster\_elements, tbl\_df-method  
 (cluster\_elements), 12

cluster\_elements, tidybulk-method  
 (cluster\_elements), 12

counts, 14

counts\_ensembl, 15

counts\_mini, 15

deconvolve\_cellularity, 15

deconvolve\_cellularity, RangedSummarizedExperiment-method  
 (deconvolve\_cellularity), 15

deconvolve\_cellularity, spec\_tbl\_df-method  
 (deconvolve\_cellularity), 15

deconvolve\_cellularity, SummarizedExperiment-method  
 (deconvolve\_cellularity), 15

deconvolve\_cellularity, tbl\_df-method  
 (deconvolve\_cellularity), 15

deconvolve\_cellularity, tidybulk-method  
 (deconvolve\_cellularity), 15

describe\_transcript, 18

distinct, 18

ensembl\_symbol\_mapping, 19

ensembl\_to\_symbol, 19

ensembl\_to\_symbol, spec\_tbl\_df-method  
 (ensembl\_to\_symbol), 19

- ensembl\_to\_symbol, tbl\_df-method  
(ensembl\_to\_symbol), 19
- ensembl\_to\_symbol, tidybulk-method  
(ensembl\_to\_symbol), 19
- fill\_missing\_abundance, 20
- fill\_missing\_abundance, RangedSummarizedExperiment-method  
(fill\_missing\_abundance), 20
- fill\_missing\_abundance, spec\_tbl\_df-method  
(fill\_missing\_abundance), 20
- fill\_missing\_abundance, SummarizedExperiment-method  
(fill\_missing\_abundance), 20
- fill\_missing\_abundance, tbl\_df-method  
(fill\_missing\_abundance), 20
- fill\_missing\_abundance, tidybulk-method  
(fill\_missing\_abundance), 20
- filter, 10, 22, 40, 49, 58
- flybaseIDs, 23
- full\_join, 24
- get\_bibliography, 24
- get\_bibliography, tidybulk-method  
(get\_bibliography), 24
- group\_by, 25
- identify\_abundant, 26
- identify\_abundant, RangedSummarizedExperiment-method  
(identify\_abundant), 26
- identify\_abundant, spec\_tbl\_df-method  
(identify\_abundant), 26
- identify\_abundant, SummarizedExperiment-method  
(identify\_abundant), 26
- identify\_abundant, tbl\_df-method  
(identify\_abundant), 26
- identify\_abundant, tidybulk-method  
(identify\_abundant), 26
- impute\_missing\_abundance, 28
- impute\_missing\_abundance, RangedSummarizedExperiment-method  
(impute\_missing\_abundance), 28
- impute\_missing\_abundance, spec\_tbl\_df-method  
(impute\_missing\_abundance), 28
- impute\_missing\_abundance, SummarizedExperiment-method  
(impute\_missing\_abundance), 28
- impute\_missing\_abundance, tbl\_df-method  
(impute\_missing\_abundance), 28
- impute\_missing\_abundance, tidybulk-method  
(impute\_missing\_abundance), 28
- inner\_join, 31
- keep\_abundant, 31
- keep\_abundant, RangedSummarizedExperiment-method  
(keep\_abundant), 31
- keep\_abundant, spec\_tbl\_df-method  
(keep\_abundant), 31
- keep\_abundant, SummarizedExperiment-method  
(keep\_abundant), 31
- keep\_abundant, tbl\_df-method  
(keep\_abundant), 31
- keep\_abundant, tidybulk-method  
(keep\_abundant), 31
- keep\_variable, 34
- keep\_variable, RangedSummarizedExperiment-method  
(keep\_variable), 34
- keep\_variable, spec\_tbl\_df-method  
(keep\_variable), 34
- keep\_variable, SummarizedExperiment-method  
(keep\_variable), 34
- keep\_variable, tbl\_df-method  
(keep\_variable), 34
- keep\_variable, tidybulk-method  
(keep\_variable), 34
- left\_join, 36
- log10\_reverse\_trans, 37
- logit\_trans, 38
- mutate, 10, 23, 38, 49, 58
- nest (unnest), 73
- parse\_formula\_survival, 40
- pivot\_sample, 40
- pivot\_sample, spec\_tbl\_df-method  
(pivot\_sample), 40
- pivot\_sample, tbl\_df-method  
(pivot\_sample), 40
- pivot\_sample, tidybulk-method  
(pivot\_sample), 40
- pivot\_transcript, 41
- pivot\_transcript, spec\_tbl\_df-method  
(pivot\_transcript), 41
- pivot\_transcript, tbl\_df-method  
(pivot\_transcript), 41
- pivot\_transcript, tidybulk-method  
(pivot\_transcript), 41
- reduce\_dimensions, 42
- reduce\_dimensions, RangedSummarizedExperiment-method  
(reduce\_dimensions), 42
- reduce\_dimensions, spec\_tbl\_df-method  
(reduce\_dimensions), 42
- reduce\_dimensions, SummarizedExperiment-method  
(reduce\_dimensions), 42
- reduce\_dimensions, tbl\_df-method  
(reduce\_dimensions), 42
- reduce\_dimensions, tidybulk-method  
(reduce\_dimensions), 42

- remove\_redundancy, 46
- remove\_redundancy, RangedSummarizedExperiment-method (remove\_redundancy), 46
- remove\_redundancy, spec\_tbl\_df-method (remove\_redundancy), 46
- remove\_redundancy, SummarizedExperiment-method (remove\_redundancy), 46
- remove\_redundancy, tbl\_df-method (remove\_redundancy), 46
- remove\_redundancy, tidybulk-method (remove\_redundancy), 46
- rename, 10, 23, 40, 49, 58
- right\_join, 50
- rotate\_dimensions, 50
- rotate\_dimensions, RangedSummarizedExperiment-method (rotate\_dimensions), 50
- rotate\_dimensions, spec\_tbl\_df-method (rotate\_dimensions), 50
- rotate\_dimensions, SummarizedExperiment-method (rotate\_dimensions), 50
- rotate\_dimensions, tbl\_df-method (rotate\_dimensions), 50
- rotate\_dimensions, tidybulk-method (rotate\_dimensions), 50
- rowwise, 53
- scale\_abundance, 54
- scale\_abundance, RangedSummarizedExperiment-method (scale\_abundance), 54
- scale\_abundance, spec\_tbl\_df-method (scale\_abundance), 54
- scale\_abundance, SummarizedExperiment-method (scale\_abundance), 54
- scale\_abundance, tbl\_df-method (scale\_abundance), 54
- scale\_abundance, tidybulk-method (scale\_abundance), 54
- se, 56
- se\_mini, 57
- summarise, 10, 23, 40, 49, 57
- symbol\_to\_entrez, 58
- test\_deseq2\_df, 59
- test\_differential\_abundance, 59
- test\_differential\_abundance, RangedSummarizedExperiment-method (test\_differential\_abundance), 59
- test\_differential\_abundance, spec\_tbl\_df-method (test\_differential\_abundance), 59
- test\_differential\_abundance, SummarizedExperiment-method (test\_differential\_abundance), 59
- test\_differential\_abundance, tbl\_df-method (test\_differential\_abundance), 59
- test\_differential\_abundance, tidybulk-method (test\_differential\_abundance), 59
- test\_differential\_cellularity, 63
- test\_differential\_cellularity, RangedSummarizedExperiment-method (test\_differential\_cellularity), 63
- test\_differential\_cellularity, spec\_tbl\_df-method (test\_differential\_cellularity), 63
- test\_differential\_cellularity, SummarizedExperiment-method (test\_differential\_cellularity), 63
- test\_differential\_cellularity, tbl\_df-method (test\_differential\_cellularity), 63
- test\_differential\_cellularity, tidybulk-method (test\_differential\_cellularity), 63
- test\_gene\_enrichment, 66
- test\_gene\_enrichment, spec\_tbl\_df-method (test\_gene\_enrichment), 66
- test\_gene\_enrichment, tbl\_df-method (test\_gene\_enrichment), 66
- test\_gene\_enrichment, tidybulk-method (test\_gene\_enrichment), 66
- test\_gene\_overrepresentation, 69
- test\_gene\_overrepresentation, spec\_tbl\_df-method (test\_gene\_overrepresentation), 69
- test\_gene\_overrepresentation, tbl\_df-method (test\_gene\_overrepresentation), 69
- test\_gene\_overrepresentation, tidybulk-method (test\_gene\_overrepresentation), 69
- tidybulk, 71
- tidybulk, RangedSummarizedExperiment-method (tidybulk), 71
- tidybulk, spec\_tbl\_df-method (tidybulk), 71
- tidybulk, SummarizedExperiment-method (tidybulk), 71
- tidybulk, tbl\_df-method (tidybulk), 71
- tidybulk\_SAM\_BAM, 72
- tidybulk\_SAM\_BAM, character, character-method (tidybulk\_SAM\_BAM), 72
- ungroup (arrange), 8
- unnest, 73

vignette\_manuscript\_signature\_boxplot,  
74  
vignette\_manuscript\_signature\_tsne, 75  
vignette\_manuscript\_signature\_tsne2,  
75  
X\_cibersort, 75