

# Package ‘rcldf’

September 30, 2025

**Type** Package

**Title** Read Linguistic Data in the Cross Linguistic Data Format (CLDF)

**Version** 1.5.1

**Maintainer** Simon J. Greenhill <[simon@simon.net.nz](mailto:simon@simon.net.nz)>

**Description** Cross-Linguistic Data Format (CLDF) is a framework for storing cross-linguistic data, ensuring compatibility and ease of data exchange between different linguistic datasets see Forkel et al. (2018) <[doi:10.1038/sdata.2018.205](https://doi.org/10.1038/sdata.2018.205)>. The ‘rcldf’ package is designed to facilitate the manipulation and analysis of these datasets by simplifying the loading, querying, and visualisation of CLDF datasets making it easier to conduct comparative linguistic analyses, manage language data, and apply statistical methods directly within R.

**License** Apache License (>= 2.0)

**Encoding** UTF-8

**Imports** archive, bib2df (>= 1.1.1), csvwr, digest, dplyr, jsonlite, logger, magrittr, purrr, readr, remotes, rlang, tools, urltools, utils

**Suggests** ggplot2, patchwork, testthat, mockthat, spelling, covr, knitr, rmarkdown, qpdf

**URL** <https://github.com/SimonGreenhill/rcldf>

**BugReports** <https://github.com/SimonGreenhill/rcldf/issues>

**Language** en-US

**RoxygenNote** 7.3.2

**VignetteBuilder** knitr

**NeedsCompilation** no

**Author** Simon J. Greenhill [aut, cre]

**Repository** CRAN

**Date/Publication** 2025-09-30 07:20:02 UTC

## Contents

add_dataframe . . . . .	2
as.cldf.wide . . . . .	3
cldf . . . . .	4
coalesce_truth . . . . .	5
datatype_to_type . . . . .	5
default_dialect . . . . .	6
default_schema . . . . .	6
get_cache_dir . . . . .	7
get_details . . . . .	7
get_dir_size . . . . .	8
get_filename . . . . .	8
get_from_zenodo . . . . .	9
get_separators . . . . .	9
get_tablelename . . . . .	10
get_table_from . . . . .	10
is_github . . . . .	11
is_url . . . . .	11
list_cache_files . . . . .	12
load_clts . . . . .	12
load_concepticon . . . . .	13
load_glottolog . . . . .	13
make_cache_key . . . . .	14
nullify . . . . .	14
print.cldf . . . . .	15
read_bib . . . . .	15
relabel . . . . .	16
resolve_path . . . . .	16
separate . . . . .	17
set_cache_dir . . . . .	17
summary.cldf . . . . .	18

## Index

19

---

<i>add_dataframe</i>	<i>Adds a dataframe.</i>
----------------------	--------------------------

---

## Description

Adds a dataframe.

## Usage

```
add_dataframe(table, filename, group)
```

**Arguments**

table	a metadata section from the CLDF metadata.
filename	the filename.
group	a grouping from the metadata.

**Value**

A dataframe

---

as.cldf.wide	<i>Extracts a CLDF table as a 'wide' dataframe by resolving all foreign key links</i>
--------------	---

---

**Description**

Extracts a CLDF table as a 'wide' dataframe by resolving all foreign key links

**Usage**

```
as.cldf.wide(object, table)
```

**Arguments**

object	the CLDF dataset.
table	the name of the table to extract.

**Value**

A tibble dataframe

**Examples**

```
md <- system.file("extdata/huon", "cldf-metadata.json", package = "rcldf")
cldfobj <- cldf(md)
forms <- as.cldf.wide(cldfobj, 'FormTable')
```

---

**cldf**

*Reads a Cross-Linguistic Data Format dataset into an object.*

---

**Description**

Reads a Cross-Linguistic Data Format dataset into an object.

included here to match people expecting e.g. `readr::read_csv` etc

**Usage**

```
cldf(
  mdpdath,
  load_bib = FALSE,
  cache_dir = tools::R_user_dir("rcldf", which = "cache")
)

read_cldf(
  mdpdath,
  load_bib = FALSE,
  cache_dir = tools::R_user_dir("rcldf", which = "cache")
)
```

**Arguments**

<code>mdpather</code>	the path to the directory or metadata JSON file.
<code>load_bib</code>	a boolean flag (TRUE/FALSE, default FALSE) to load the <code>sources.bib</code> BibTeX file. <code>load_bib=FALSE</code> can easily speed up loading of a CLDF dataset by an order of magnitude or two, so we do not load sources by default.
<code>cache_dir</code>	a directory to cache downloaded files to

**Value**

A `cldf` object

**Examples**

```
cldfobj <- cldf(system.file("extdata/huon", "cldf-metadata.json", package = "rcldf"))
```

---

coalesce_truth	<i>Coalesce value to truthiness</i>
----------------	-------------------------------------

---

### Description

Determine whether the input is true, with missing values being interpreted as false.

### Usage

```
coalesce_truth(x)
```

### Arguments

x	logical, NA or NULL
---	---------------------

### Value

FALSE if x is anything but TRUE

---

datatype_to_type	<i>Map csvw datatypes to R types</i>
------------------	--------------------------------------

---

### Description

Translate **csvw datatypes** to R types. This implementation currently targets `readr::cols` column specifications.

### Usage

```
datatype_to_type(datatypes)
```

### Arguments

datatypes	a list of csvw datatypes
-----------	--------------------------

### Details

rcldf adds some overrides here to add e.g. anyURI etc.

### Value

a `readr::cols` specification - a list of collectors

### Examples

```
cspec <- datatype_to_type(list("double", list(base="date", format="yyyy-MM-dd")))  
readr::read_csv(readr::readr_example("challenge.csv"), col_types=cspec)
```

<code>default_dialect</code>	<i>CSVW default dialect</i>
------------------------------	-----------------------------

### Description

The [CSVW Default Dialect specification](#) described in [CSV Dialect Description Format](#).

### Usage

```
default_dialect
```

### Format

An object of class `list` of length 13.

### Value

a list specifying a default csv dialect

<code>default_schema</code>	<i>Create a default table schema given a csv file and dialect</i>
-----------------------------	---

### Description

If neither the table nor the group have a `tableSchema` annotation, then this default schema will be used.

### Usage

```
default_schema(filename, dialect = default_dialect)
```

### Arguments

<code>filename</code>	a csv file
-----------------------	------------

<code>dialect</code>	specification of the csv's dialect (default: <code>default_dialect</code> )
----------------------	---

### Value

a table schema

---

get_cache_dir	<i>Returns the cache dir.</i>
---------------	-------------------------------

---

### Description

Returns the cache dir.

### Usage

```
get_cache_dir(cache_dir = NA)
```

### Arguments

cache\_dir      a directory to use

### Value

A string of the cache dir

---

get_details	<i>Returns a dataframe of with details on the CLDF dataset in path.</i>
-------------	---

---

### Description

Returns a dataframe of with details on the CLDF dataset in path.

### Usage

```
get_details(path, cache_dir = NA)
```

### Arguments

path            the path to resolve  
cache\_dir      a directory to cache downloaded files to

### Value

A dataframe.

---

<code>get_dir_size</code>	<i>Returns the filesize in bytes of a directory.</i>
---------------------------	--

---

### Description

Returns the filesize in bytes of a directory.

### Usage

```
get_dir_size(path)
```

### Arguments

path	a directory to size
------	---------------------

### Value

A numeric of the file size in bytes

---

<code>get_filename</code>	<i>Get a filename from url value in metadata (handles .zip files)</i>
---------------------------	---

---

### Description

Get a filename from url value in metadata (handles .zip files)

### Usage

```
get_filename(base_dir, url)
```

### Arguments

base_dir	the base_dir
url	the url statement

### Value

A string

---

get_from_zenodo	<i>Downloads and installs a CLDF dataset from a Zenodo endpoint</i>
-----------------	---

---

**Description**

Downloads and installs a CLDF dataset from a Zenodo endpoint

**Usage**

```
get_from_zenodo(zid, load_bib = FALSE, cache_dir = NULL)
```

**Arguments**

zid	Zenodo endpoint conceptid
load_bib	load sources (TRUE/FALSE, default FALSE)
cache_dir	A cache_dir to use. If NULL it will use get_cache_dir

**Value**

A cldf object

---

get_separators	<i>Identifies the separator characters specified by the CLDF metadata.</i>
----------------	--

---

**Description**

Identifies the separator characters specified by the CLDF metadata.

**Usage**

```
get_separators(metadata)
```

**Arguments**

metadata	• a CLDF metadata.
----------	--------------------

**Value**

A dataframe with three columns (name, separator, url).

---

<code>get_tablename</code>	<i>Convert a CLDF URL tablename to a short tablename</i>
----------------------------	--

---

### Description

Convert a CLDF URL tablename to a short tablename

### Usage

```
get_tablename(conformsto, url = NA)
```

### Arguments

<code>conformsto</code>	the dc:conforms to statement
<code>url</code>	the url statement

### Value

A string

### Examples

```
get_tablename("http://cldf.clld.org/v1.0/terms.rdf#ValueTable")
```

---

<code>get_table_from</code>	<i>Extracts a single table from a CLDF dataset.</i>
-----------------------------	---

---

### Description

Extracts a single table from a CLDF dataset.

### Usage

```
get_table_from(
  table,
  mdpPath,
  cache_dir = tools::R_user_dir("rcldf", which = "cache")
)
```

### Arguments

<code>table</code>	a CLDF table type
<code>mdpPath</code>	a path to a CLDF file
<code>cache_dir</code>	a directory to cache downloaded files to

**Value**

a dataframe

**Examples**

```
md_json <- system.file("extdata/huon", "cldf-metadata.json", package = "rcldf")
df <- get_table_from("LanguageTable", md_json)
```

---

is\_github

*Returns TRUE if url looks like a github URL*

---

**Description**

Returns TRUE if url looks like a github URL

**Usage**

```
is_github(url)
```

**Arguments**

url                  A string

**Value**

A boolean TRUE/FALSE

**Examples**

```
is_github('https://github.com/SimonGreenhill/rcldf/')
```

---

is\_url

*Returns TRUE if url looks like a URL*

---

**Description**

Returns TRUE if url looks like a URL

**Usage**

```
is_url(url)
```

**Arguments**

url                  A string

**Value**

A boolean TRUE/FALSE

**Examples**

```
is_url('http://simon.net.nz')
```

<code>list_cache_files</code>	<i>Returns a dataframe of directories in the cache dir</i>
-------------------------------	--

**Description**

Returns a dataframe of directories in the cache dir

**Usage**

```
list_cache_files(cache_dir = NULL)
```

**Arguments**

`cache_dir` the cache directory to use. If NULL then R\_user\_dir will be used.

**Value**

A dataframe of the directories

<code>load_clts</code>	<i>Returns a CLDF dataset object of the latest CLTS version.</i>
------------------------	--

**Description**

Returns a CLDF dataset object of the latest CLTS version.

**Usage**

```
load_clts(load_bib = FALSE, cache_dir = NULL)
```

**Arguments**

`load_bib` load sources (TRUE/FALSE, default FALSE)

`cache_dir` A cache\_dir to use. If NULL it will use get\_cache\_dir

**Value**

A cldf object

---

load_concepticon	<i>Returns a CLDF dataset object of the latest Concepticon version.</i>
------------------	---

---

**Description**

Returns a CLDF dataset object of the latest Concepticon version.

**Usage**

```
load_concepticon(load_bib = FALSE, cache_dir = NULL)
```

**Arguments**

load_bib	load sources (TRUE/FALSE, default FALSE)
cache_dir	A cache_dir to use. If NULL it will use get_cache_dir

**Value**

A cldf object

---

load_glottolog	<i>Returns a CLDF dataset object of the latest glottolog version.</i>
----------------	---

---

**Description**

Returns a CLDF dataset object of the latest glottolog version.

**Usage**

```
load_glottolog(load_bib = FALSE, cache_dir = NULL)
```

**Arguments**

load_bib	load sources (TRUE/FALSE, default FALSE)
cache_dir	A cache_dir to use. If NULL it will use get_cache_dir

**Value**

A cldf object

`make_cache_key`      *Returns the cachekey for the given path.*

### Description

Returns the cachekey for the given path.

### Usage

`make_cache_key(path)`

### Arguments

`path`      a path to generate the cachekey for.

### Value

A string.

`nullify`      *Converts all values specified in the CLDF metadata as null to R's NA.*

### Description

Note that this is run by default on loading a dataset with `cldf()`

### Usage

`nullify(cldfobj, nulls = NULL)`

### Arguments

`cldfobj`      a CLDF Object  
`nulls`      a dataframe of null values to replace (default=NULL).

### Value

A `cldf` object

### Examples

```
cldfobj <- cldf(system.file("extdata/huon", "cldf-metadata.json", package = "rcldf"))
cldfobj <- nullify(cldfobj)
```

---

print.cldf	<i>Summarises the CLDF file</i>
------------	---------------------------------

---

### Description

Summarises the CLDF file

### Usage

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

### Arguments

x	the CLDF dataset
...	Arguments to be passed to or from other methods. Currently not used.

### Value

No return value, called for side effects.

### Examples

```
cldfobj <- cldf(system.file("extdata/huon", "cldf-metadata.json", package = "rcldf"))  
print(cldfobj)
```

---

---

read_bib	<i>Adds BibTeX source information into a CLDF dataset</i>
----------	---

---

### Description

Adds BibTeX source information into a CLDF dataset

### Usage

```
read_bib(object)
```

### Arguments

object	A CLDF object
--------	---------------

### Value

A tibble dataframe

<code>relabel</code>	<i>Relabels a column in a dataset for merging.</i>
----------------------	--

### Description

Relabels a column in a dataset for merging.

### Usage

```
relabel(column, table)
```

### Arguments

<code>column</code>	the tablename.
<code>table</code>	the tablename.

### Value

A string of "column.table"

<code>resolve_path</code>	<i>Helper function to resolve the path (e.g. directory or md.json file)</i>
---------------------------	---

### Description

Helper function to resolve the path (e.g. directory or md.json file)

### Usage

```
resolve_path(path, cache_dir = NA)
```

### Arguments

<code>path</code>	the path to resolve
<code>cache_dir</code>	a directory to cache downloaded files to

### Value

A list of two items: `path` - string containing the path to the `metadata.json` file `metadata` - a `csvwr` `metadata` object

---

separate	<i>Expands all values with separators.</i>
----------	--

---

## Description

Note that this is run by default on loading a dataset with `cldf()`

## Usage

```
separate(cldfobj, separators = NULL)
```

## Arguments

<code>cldfobj</code>	a CLDF Object
<code>separators</code>	a dataframe of separator values to replace (default=NULL).

## Value

A `cldf` object

## Examples

```
cldfobj <- cldf(system.file("extdata/huon", "cldf-metadata.json", package = "rcldf"))
cldfobj <- separate(cldfobj)
```

---

set_cache_dir	<i>Sets the cache dir for the current session.</i>
---------------	--

---

## Description

Sets the cache dir for the current session.

## Usage

```
set_cache_dir(cache_dir = NA)
```

## Arguments

<code>cache_dir</code>	a directory to use
------------------------	--------------------

## Value

NULL. Sets an environment value.

---

`summary.cldf`                    *Summarises the CLDF file*

---

## Description

Summarises the CLDF file

## Usage

```
## S3 method for class 'cldf'  
summary(object, ...)
```

## Arguments

<code>object</code>	the CLDF dataset
<code>...</code>	Arguments to be passed to or from other methods. Currently not used.

## Value

None

## Examples

```
cldfobj <- cldf(system.file("extdata/huon", "cldf-metadata.json", package = "rcldf"))  
summary(cldfobj)
```

# Index

\* **datasets**  
    default\_dialect, 6  
    add\_dataframe, 2  
    as.cldf.wide, 3  
  
    cldf, 4  
    coalesce\_truth, 5  
  
    datatype\_to\_type, 5  
    default\_dialect, 6  
    default\_schema, 6  
  
    get\_cache\_dir, 7  
    get\_details, 7  
    get\_dir\_size, 8  
    get\_filename, 8  
    get\_from\_zendodo, 9  
    get\_separators, 9  
    get\_table\_from, 10  
    get\_tablename, 10  
  
    is\_github, 11  
    is\_url, 11  
  
    list\_cache\_files, 12  
    load\_clts, 12  
    load\_concepticon, 13  
    load\_glottolog, 13  
  
    make\_cache\_key, 14  
  
    nullify, 14  
  
    print.cldf, 15  
  
    read\_bib, 15  
    read\_cldf (cldf), 4  
    readr::cols, 5  
    relabel, 16  
    resolve\_path, 16

    separate, 17  
    set\_cache\_dir, 17  
    summary.cldf, 18