

# Package ‘medrxivr’

September 28, 2020

**Title** Access and Search MedRxiv and BioRxiv Preprint Data

**Version** 0.0.3

**Description** An increasingly important source of health-related bibliographic content are preprints - preliminary versions of research articles that have yet to undergo peer review. The two preprint repositories most relevant to health-related sciences are medRxiv <<https://www.medrxiv.org/>> and bioRxiv <<https://www.biorxiv.org/>>, both of which are operated by the Cold Spring Harbor Laboratory. 'medrxivr' provides programmatic access to the 'Cold Spring Harbour Laboratory (CSHL)' API <<https://api.biorxiv.org/>>, allowing users to easily download medRxiv and bioRxiv preprint metadata (e.g. title, abstract, publication date, author list, etc) into R. 'medrxivr' also provides functions to search the downloaded preprint records using regular expressions and Boolean logic, as well as helper functions that allow users to export their search results to a .BIB file for easy import to a reference manager and to download the full-text PDFs of preprints matching their search criteria.

**License** MIT + file LICENSE

**Encoding** UTF-8

**LazyData** true

**Language** en-US

**URL** <https://github.com/ropensci/medrxivr>

**BugReports** <https://github.com/ropensci/medrxivr/issues>

**Imports** methods, dplyr, curl, jsonlite, httr, stringr, rlang, vroom, bib2df, tibble, progress, lubridate

**Suggests** testthat (>= 2.1.0), knitr, rmarkdown, covr, kableExtra, spelling

**VignetteBuilder** knitr, rmarkdown

**RoxygenNote** 7.1.1

**NeedsCompilation** no

**Author** Luke McGuinness [aut, cre] (<<https://orcid.org/0000-0001-8730-9761>>),  
 Lena Schmidt [aut] (<<https://orcid.org/0000-0003-0709-8226>>),  
 Tuija Sonkkila [rev],  
 Najko Jahn [rev]

**Maintainer** Luke McGuinness <[luke.mcguinness@bristol.ac.uk](mailto:luke.mcguinness@bristol.ac.uk)>

**Repository** CRAN

**Date/Publication** 2020-09-28 09:40:02 UTC

## R topics documented:

medrxivr . . . . .	2
mx_api_content . . . . .	2
mx_api_doi . . . . .	4
mx_crosscheck . . . . .	5
mx_download . . . . .	5
mx_export . . . . .	6
mx_search . . . . .	7
mx_snapshot . . . . .	8

<b>Index</b>	<b>9</b>
--------------	----------

---

medrxivr	<i>medrxivr: Accessing medRxiv and bioRxiv preprint data from R</i>
----------	---

---

### Description

The medrxivr package enables users to access data on preprints in the medRxiv and bioRxiv preprints repositories, both of which are run by the Cold Spring Harbour Laboratory. It also provides functions to search the preprint data, export it to a .bib file, and download the PDFs associated with specified records.

---

mx_api_content	<i>Access medRxiv/bioRxiv data via the Cold Spring Harbour Laboratory API</i>
----------------	---

---

### Description

Provides programmatic access to all preprints available through the Cold Spring Harbour Laboratory API, which serves both the medRxiv and bioRxiv preprint repositories.

**Usage**

```
mx_api_content(  
  from_date = "2013-01-01",  
  to_date = as.character(Sys.Date()),  
  clean = TRUE,  
  server = "medrxiv",  
  include_info = FALSE  
)
```

**Arguments**

from_date	Earliest date of interest. Defaults to 1st June 2019 (earliest medRxiv record was posted on 25th June 2019).
to_date	Latest date of interest. Defaults to current date.
clean	Logical, defaulting to TRUE, indicating whether to clean the data returned by the API. If TRUE, variables containing absolute paths to the preprints web-page ("link_page") and PDF ("link_pdf") are generated from the "server", "DOI", and "version" variables returned by the API. The "title", "abstract" and "authors" variables are converted to title case. Finally, the "type" and "server" variables are dropped.
server	Specify the server you wish to use: "medrxiv" (default) or "biorxiv"
include_info	Logical, indicating whether to include variables containing information returned by the API (e.g. API status, cursor number, total count of papers, etc). Default is FALSE.

**Value**

Dataframe with 1 record per row

**See Also**

Other data-source: [mx\\_api\\_doi\(\)](#), [mx\\_snapshot\(\)](#)

**Examples**

```
mx_data <- mx_api_content(from_date = "2020-01-01",  
  to_date = "2020-01-07")
```

---

mx_api_doi	<i>Access data on a single medRxiv/bioRxiv record via the Cold Spring Harbour Laboratory API</i>
------------	--

---

## Description

Provides programmatic access to data on a single preprint identified by a unique Digital Object Identifier (DOI).

## Usage

```
mx_api_doi(doi, server = "medrxiv", clean = TRUE)
```

## Arguments

doi	Digital object identifier of the preprint you wish to retrieve data on.
server	Specify the server you wish to use: "medrxiv" (default) or "biorxiv"
clean	Logical, defaulting to TRUE, indicating whether to clean the data returned by the API. If TRUE, variables containing absolute paths to the preprints web-page ("link_page") and PDF ("link_pdf") are generated from the "server", "DOI", and "version" variables returned by the API. The "title", "abstract" and "authors" variables are converted to title case. Finally, the "type" and "server" variables are dropped.

## Value

Dataframe containing details on the preprint identified by the DOI.

## See Also

Other data-source: [mx\\_api\\_content\(\)](#), [mx\\_snapshot\(\)](#)

## Examples

```
mx_data <- mx_api_doi("10.1101/2020.02.25.20021568")
```

---

mx_crosscheck	<i>Check how up-to-date the maintained medRxiv snapshot is</i>
---------------	--

---

### Description

Provides information on how up-to-date the maintained medRxiv snapshot provided by ‘mx\_snapshot()’ is by checking whether there have been any records added to, or updated in, the medRxiv repository since the last snapshot was taken.

### Usage

```
mx_crosscheck()
```

### See Also

Other helper: [mx\\_download\(\)](#), [mx\\_export\(\)](#)

### Examples

```
mx_crosscheck()
```

---

mx_download	<i>Download PDF's of preprints returned by a search</i>
-------------	---

---

### Description

Download PDF's of all the papers in your search results

### Usage

```
mx_download(
  mx_results,
  directory,
  create = TRUE,
  name = c("ID", "DOI"),
  print_update = 10
)
```

### Arguments

mx_results	Vector containing the links to the medRxiv PDFs
directory	The location you want to download the PDF's to
create	TRUE or FALSE. If TRUE, creates the directory if it doesn't exist
name	How to name the downloaded PDF. By default, both the ID number of the record and the DOI are used.
print_update	How frequently to print an update

**See Also**

Other helper: [mx\\_crosscheck\(\)](#), [mx\\_export\(\)](#)

**Examples**

```
mx_results <- mx_search(mx_snapshot(), query = "10.1101/2020.02.25.20021568")
mx_download(mx_results, directory=tempdir())
```

---

mx\_export

*Export references for preprints returning by a search to a .bib file*

---

**Description**

Export references for preprints returning by a search to a .bib file

**Usage**

```
mx_export(data, file = "medrxiv_export.bib")
```

**Arguments**

data	Dataframe returned by <a href="#">mx_search()</a> or <a href="#">mx_api_*</a> () functions
file	File location to save to. Must have the .bib file extension

**Value**

Exports a formatted .BIB file, for import into a reference manager

**See Also**

Other helper: [mx\\_crosscheck\(\)](#), [mx\\_download\(\)](#)

**Examples**

```
mx_results <- mx_search(mx_snapshot(), query = "brain")
mx_export(mx_results, tempfile(fileext = ".bib"))
```

---

mx_search	<i>Search preprint data</i>
-----------	-----------------------------

---

### Description

Search preprint data

### Usage

```
mx_search(
  data = NULL,
  query = NULL,
  fields = c("title", "abstract", "authors", "category", "doi"),
  from_date = NULL,
  to_date = NULL,
  NOT = "",
  deduplicate = TRUE
)
```

### Arguments

data	The preprint dataset that is to be searched, created either using <code>mx_api_content()</code> or <code>mx_snapshot()</code>
query	Character string, vector or list
fields	Fields of the database to search - default is Title, Abstract, Authors, Category, and DOI.
from_date	Defines earliest date of interest. Written in the format "YYYY-MM-DD". Note, records published on the date specified will also be returned.
to_date	Defines latest date of interest. Written in the format "YYYY-MM-DD". Note, records published on the date specified will also be returned.
NOT	Vector of regular expressions to exclude from the search. Default is NULL.
deduplicate	Logical. Only return the most recent version of a record. Default is TRUE.

### Examples

```
# Using the daily snapshot
mx_results <- mx_search(data = mx_snapshot(), query = "dementia")
```

---

`mx_snapshot`*Access a static snapshot of the medRxiv repository*

---

**Description**

[Available for medRxiv only] Rather than downloading a copy of the medRxiv database from the API, which can become unavailable at peak usage times, this allows users to import a maintained static snapshot of the medRxiv repository.

**Usage**

```
mx_snapshot(commit = "master")
```

**Arguments**

<code>commit</code>	Commit hash for the snapshot, taken from <a href="https://github.com/mcguinlu/medrxiv-data">https://github.com/mcguinlu/medrxiv-data</a> . Allows for reproducible searching by specifying the exact snapshot used to perform the searches. Defaults to "master", which will return the most recent snapshot.
---------------------	---

**Value**

Formatted dataframe

**See Also**

Other data-source: [mx\\_api\\_content\(\)](#), [mx\\_api\\_doi\(\)](#)

**Examples**

```
mx_data <- mx_snapshot()
```



# Index

- \* **data-source**

- mx\_api\_content, 2
  - mx\_api\_doi, 4
  - mx\_snapshot, 8

- \* **helper**

- mx\_crosscheck, 5
  - mx\_download, 5
  - mx\_export, 6

- \* **main**

- mx\_search, 7

medrxivr, 2

mx\_api\_content, 2, 4, 8

mx\_api\_doi, 3, 4, 8

mx\_crosscheck, 5, 6

mx\_download, 5, 5, 6

mx\_export, 5, 6, 6

mx\_search, 7

mx\_snapshot, 3, 4, 8