

# Retrieve clinical trial information

Ralf Herold

2024-05-12

## Get started

### Attach package ctrdata

```
library(ctrdata)
citation("ctrdata")
```

Remember to respect the registers' terms and conditions (see `ctrOpenSearchPagesInBrowser(copyright = TRUE)`). Please cite this package in any publication as follows: Ralf Herold (2024). `ctrdata`: Retrieve and Analyze Clinical Trials in Public Registers. R package version 1.18.0. <https://cran.r-project.org/package=ctrdata>

### Open register's advanced search page in browser

These functions open the browser, where the user can start searching for trials of interest.

```
# Please review and respect register copyrights:
ctrOpenSearchPagesInBrowser(
  copyright = TRUE
)
# Open browser with example search:
ctrOpenSearchPagesInBrowser(
  url = "cancer&age=under-18",
  register = "EUCTR"
)
```

### Adjust search parameters and execute search in browser

Refine the search until the trials of interest are listed in the browser. The total number of trials that can be retrieved with package `ctrdata` is intentionally limited to queries with at most 10,000 result records.

### Copy address from browser address bar to clipboard

Use functions or keyboard shortcuts according to the operating system. See here for our automation to copy the URLs of a user's queries in any of the supported clinical trial registers.

## Get address from clipboard

The next steps are executed in the R environment:

```
q <- ctrGetQueryUrl()
# * Using clipboard content as register query URL: https://www.clinicaltrialsregister.eu/ctr-search/search
# * Found search query from EUCTR: query=cancer&age=under-18&status=completed&phase=phase-one

q
#
# 1 query=cancer&age=under-18&status=completed&phase=phase-one      query-term  query-register      EUCTR

# To check, this opens a browser with the query
ctrOpenSearchPagesInBrowser(url = q)
```

## Retrieve protocol-related information, transform, save to database, check

```
# Count number of trial records
ctrLoadQueryIntoDb(
  queryterm = q,
  only.count = TRUE
)$n
# * Found search query from EUCTR: query=cancer&age=under-18&status=completed&phase=phase-one
# * Checking trials in EUCTR...
# Retrieved overview, multiple records of 103 trial(s) from 6 page(s) to be downloaded (estimate: 10 MB)
# [1] 103

# Connect to a database and chose a collection (table)
db <- nodbi::src_sqlite(
  dbname = "sqlite_file.sql",
  collection = "test"
)

# Retrieve records, download into database
ctrLoadQueryIntoDb(
  queryterm = q,
  con = db
)
# * Checking trials in EUCTR...
# Retrieved overview, multiple records of 103 trial(s) from 6 page(s) to be downloaded (estimate: 10 MB)
# (1/3) Downloading trials...
# Note: register server cannot compress data, transfer takes longer (estimate: 100 s)
# Download status: 6 done; 0 in progress. Total size: 8.73 Mb (100%)... done!
# (2/3) Converting to NDJSON (estimate: 2 s)...
# (3/3) Importing records into database...
# = Imported or updated 414 records on 103 trial(s)
# No history found in expected format.
# Updated history ("meta-info" in "test")
# $n
# [1] 414

# Show which queries have been downloaded into database
```

```

dbQueryHistory(con = db)
#           query-timestamp query-register query-records
# 1 2024-05-12 17:56:11          EUCTR          414 query=cancer&age=under-18&status=completed&phase=p

```

With a file-base SQLite database, this takes about 20 seconds for about 300 records, with most of the time needed for internet-retrieval with is slow from this register. Speed is higher with other registers, with using MongoDB and with memory-based SQLite.

## Repeat and update a previous query

Instead of “last”, an integer number can be specified for `querytoupdate` that corresponds to the number when using `dbQueryHistory()`.

```

ctrLoadQueryIntoDb(
  querytoupdate = "last",
  con = db
)

```

Depending on the register, an update (differential update) is possible or the original query is executed fully again.

## Retrieve results

For EUCTR, result-related trial information has to be requested to be retrieved, because it will take longer to download and store. For CTGOV, any results are always included in the retrieval. Note that trial documents, including any results reports, can be downloaded by specifying parameter `documents.path`, see `help(ctrLoadQueryIntoDb)`.

```

ctrLoadQueryIntoDb(
  queryterm = q,
  euctrresults = TRUE,
  con = db
)
# * Checking trials in EUCTR...
# Retrieved overview, multiple records of 103 trial(s) from 6 page(s) to be downloaded (estimate: 10 MB)
# (1/3) Downloading trials...
# Note: register server cannot compress data, transfer takes longer (estimate: 100 s)
# (2/3) Converting to NDJSON (estimate: 2 s)...
# (3/3) Importing records into database...
# = Imported or updated 414 records on 103 trial(s)
# * Checking results if available from EUCTR for 103 trials:
# (1/4) Downloading results...
# Download status: 103 done; 0 in progress. Total size: 59.71 Mb (100%)... done!
# Download status: 27 done; 0 in progress. Total size: 108.69 Kb (100%)... done!
# - extracting results (. = data, F = file[s] and data, x = none):
# . . . . . F . . . . . F . . . . . F . . . . .
# . . . F . F F . F . . F . . . F F . . . . . F . . . . . F . . . F .
# (2/4) Converting to NDJSON (estimate: 8 s)...
# (3/4) Importing results into database (may take some time)...
# (4/4) Results history: not retrieved (euctrresultshistory = FALSE)
# = Imported or updated results for 76 trials

```

```
# Updated history ("meta-info" in "test")
# $n
# [1] 414
```

With a file-base SQLite database, this takes about 4 minutes for about 300 records, with most of the time needed for merging result- and protocol-related information in SQLite; this is much faster with MongoDB and PostgreSQL.

The download or presence of results is not recorded in `dbQueryHistory()` because the availability of results increases over time.

## Add information from another register

The same collection can be used to store (and analyse) trial information from different registers. Example:

```
ctrLoadQueryIntoDb(
  queryterm = "https://classic.clinicaltrials.gov/ct2/results?cond=neuroblastoma&recrs=e&age=0&intr=Drug
  con = db
)
# * Appears specific for CTGOV Classic website
# Since 2024-06-25, the classic CTGOV servers are no longer available. Package
# ctrdata has translated the classic CTGOV query URL from this call of function
# ctrLoadQueryIntoDb(queryterm = ...) into a query URL that works with the current
# CTGOV2. This is printed below and is also part of the return value of this function,
# ctrLoadQueryIntoDb(...)$url. This URL can be used with ctrdata functions.
# Note that the fields and data schema of trials differ between CTGOV and CTGOV2.
#
# Replace this URL:
#
# https://classic.clinicaltrials.gov/ct2/results?cond=neuroblastoma&recrs=e&age=0&intr=Drug
#
# with this URL:
#
# https://clinicaltrials.gov/search?cond=neuroblastoma&intr=Drug&aggFilters=ages:child,status:com
#
# * Found search query from CTGOV2: cond=neuroblastoma&intr=Drug&aggFilters=ages:child,status:com
# * Checking trials using CTGOV REST API 2.0, found 222 trials
# (1/3) Downloading in 1 batch(es) (max. 1000 trials each; estimate: 22 MB total)
# Download status: 1 done; 0 in progress. Total size: 9.79 Mb (615%)... done!
# (2/3) Converting to NDJSON...
# (3/3) Importing records into database...
# JSON file #: 1 / 1
# = Imported or updated 222 trial(s)
# No history found in expected format.
# Updated history ("meta-info" in "test")
# $n
# [1] 222
```

With a file-base SQLite database, this takes about 10 seconds for about 200 records.

Note that in this example, a warning message may be issued from importing an NDJSON file with trial records. The warning arises from the high level of complexity of some of the XML content of some of the trial records. The issue can be resolved by increasing in the operating system the stack size available to R, see: <https://github.com/rfhhb/ctrdata/issues/22>

## Add records from CTIS register into the same collection

As of 29 June 2024, more than 4300 trials are publicly accessible in CTIS, which has to be used since January 2023 for new clinical trials in the European Union (EU). Queries in the CTIS search interface can be automatically copied to the clipboard so that a user can paste them into `queryterm`, see here.

```
# Retrieve trials from another register:
ctrLoadQueryIntoDb(
  queryterm = 'searchCriteria={"containAll":"","containAny":"neonates","containNot":""}',
  register = "CTIS",
  con = db
)
# * Found search query from CTIS: searchCriteria={"containAll":"","containAny":"neonates","containNot":
# * Checking trials in CTIS...
# (1/4) Downloading trial list(s), found 7 trials
# (2/4) Downloading and processing trial data... (estimate: 0.6 Mb)
# Download status: 7 done; 0 in progress. Total size: 262.58 Kb (100%)... done!
# (3/4) Importing records into database...
# (4/4) Updating with additional data: .
# = Imported 7, updated 7 record(s) on 7 trial(s)
# Updated history ("meta-info" in "test")
# $n
# [1] 7
```

## Add personal annotations

When downloading trial information, the user can specify an annotation to all records that are downloaded. By default, annotations are accumulated if trial records are loaded again or updated; alternatively, annotations can be replaced.

Annotations are useful for analyses, for example to specially identify subsets of records in the database.

```
ctrLoadQueryIntoDb(
  queryterm = paste0(
    "https://classic.clinicaltrials.gov/ct2/results?",
    "cond=neuroblastoma&recrs=e&age=0&type=Intr&cntry=DE"),
  annotation.text = "site_DE ",
  annotation.mode = "append",
  con = db
)
# * Appears specific for CTGOV Classic website
# Since 2024-06-25, the classic CTGOV servers are no longer available. Package
# ctrdata has translated the classic CTGOV query URL from this call of function
# ctrLoadQueryIntoDb(queryterm = ...) into a query URL that works with the current
# CTGOV2. This is printed below and is also part of the return value of this function,
# ctrLoadQueryIntoDb(...) $url. This URL can be used with ctrdata functions.
# Note that the fields and data schema of trials differ between CTGOV and CTGOV2.
#
# Replace this URL:
#
# https://classic.clinicaltrials.gov/ct2/results?cond=neuroblastoma&recrs=e&age=0&type=Intr&cntry=DE
#
# with this URL:
```

```
#
# https://clinicaltrials.gov/search?cond=neuroblastoma&country=Germany&aggFilters=ages:child,studyType:
#
# * Found search query from CTGOV2:
#   cond=neuroblastoma&country=Germany&aggFilters=ages:child,studyType:int,status:com
# * Checking trials using CTGOV REST API 2.0, found 16 trials
# (1/3) Downloading in 1 batch(es) (max. 1000 trials each; estimate: 1.6 MB total)
# Download status: 1 done; 0 in progress. Total size: 487.83 Kb (492%)... done!
# (2/3) Converting to NDJSON...
# (3/3) Importing records into database...
# JSON file #: 1 / 1
# = Imported or updated 16 trial(s)
# = Annotated retrieved records (16 records)
# Updated history ("meta-info" in "test")
# $n
# [1] 16
```

## Add information using trial identifiers

When identifiers of clinical trials of interest are already known, this example shows how they can be processed to import the trial information into a database collection. This involves constructing a query that combines the identifiers and then iterating over the sets of identifiers. Note to combine identifiers using “+OR+” into the `queryterm`, and that `register` has to be specified.

```
# ids of trials of interest
ctIds <- c(
  "NCT00001209", "NCT00001436", "NCT00187109", "NCT01516567", "NCT01471782", "NCT03042429",
  "NCT00357084", "NCT00357500", "NCT00365755", "NCT00407433", "NCT00410657", "NCT00436657",
  "NCT00436852", "NCT00445965", "NCT00450307", "NCT00450827", "NCT00471679", "NCT00486564",
  "NCT00492167", "NCT00499616", "NCT00503724", "NCT00509353", "NCT00520936", "NCT00536601",
  "NCT00567567", "NCT00578864", "NCT00601003", "NCT00644696", "NCT00646230", "NCT00659984",
  "NCT00716976", "NCT00743496", "NCT00793845", "NCT00806182", "NCT00831844", "NCT00867568",
  "NCT00877110", "NCT00885326", "NCT00918320", "NCT00923351", "NCT00939770", "NCT00960739"
)

# split into sets of each 25 trial ids
# (larger sets e.g. 50 may still work)
idSets <- split(ctIds, ceiling(seq_along(ctIds) / 25))

# variable to collect import results
result <- NULL

# iterate over sets of trial ids
for (idSet in idSets) {

  # import
  setResult <- ctrLoadQueryIntoDb(

    # for CTGOV2 use:
    queryterm = paste0("term=", paste0(idSet, collapse = " ")),

    # specify register that holds the information
    register = "CTGOV2",
```

```

    con = db
  )

  # append results
  result <- c(result, list(setResult))
}

# inspect results
as.data.frame(do.call(rbind, result))[, c("n", "failed")]
#      n failed
# 1 25    NULL
# 2 17    NULL

```

Note that from CTIS, trial information *identified by trial identifiers* can be retrieved only one-by-one, repeating queries for each trial such as <https://euclinicaltrials.eu/app/#/search?number=2023-503994-39-00>.

## Find synonyms of active substance names

Not all registers automatically expand search terms to include alternative terms, such as codes and other names of active substances. To obtain a character vector of synonyms for any active substance name, use:

```

ctrFindActiveSubstanceSynonyms(
  activesubstance = "imatinib"
)
# [1] "imatinib"           "Carcemia"           "Cemivil"
# [4] "CGP 57148"          "CGP-57148B"         "CGP57148B"
# [7] "Gleevac"            "Gleevec"            "Gleevec (Imatinib Mesylate)"
# [10] "Gleevec"            "Glivec"             "Imatinib"
# [13] "Imatinib Mesylate"  "Imatinib-AFT"       "NSC #716051"
# [16] "NSC-716051"         "QTI571"             "STI571"
# [19] "STI 571"           "STI-571"            "STI571"
# [22] "tyrosine kinase inhibitors" "Ziatar"             "Zoleta"

```

These names can then be used in queries in any register.

## Using a MongoDB database

This example works with a free service here. Note that the user name and password need to be encoded. The format of the connection string is documented at <https://docs.mongodb.com/manual/reference/connection-string/>.

For recommended databases, see vignette `Install R package ctrdata`.

```

# Specify base uri for remote MongoDB server,
# as part of the encoded connection string
db <- nodbi::src_mongo(
  # Note: this provides read-only access
  url = "mongodb+srv://DWbJ7Wh:bdTHh5cS@cluster0-b9wpw.mongodb.net",
  db = "dbperm",
  collection = "dbperm")

```

```

# Since the above access is read-only,
# just obtain fields of interest:
dbGetFieldsIntoDf(
  fields = c("a2_eudract_number",
             "e71_human_pharmacology_phase_i"),
  con = db)

```

#	<i>_id</i>	<i>a2_eudract_number</i>	<i>e71_human_pharmacology_phase_i</i>
# 1	2010-024264-18-3RD	2010-024264-18	TRUE
# 2	2010-024264-18-AT	2010-024264-18	TRUE
# 3	2010-024264-18-DE	2010-024264-18	TRUE
# 4	2010-024264-18-GB	2010-024264-18	TRUE
# 5	2010-024264-18-IT	2010-024264-18	TRUE
# 6	2010-024264-18-NL	2010-024264-18	TRUE