

Package: **homologene** (via r-universe)

October 27, 2024

Type Package

Title Quick Access to Homologene and Gene Annotation Updates

Version 1.7.68.23.10.31

Depends R (>= 3.1.2)

Imports dplyr (>= 0.7.4), magrittr (>= 1.5), purrr (>= 0.2.5), readr (>= 1.3.1), R.utils (>= 2.8.0), assertthat (>= 0.2.1), rvest (>= 1.0.0), xml2 (>= 1.3.2)

Suggests testthat (>= 1.0.2)

Date 2023-10-31

BugReports <https://github.com/oganm/homologene/issues>

URL <https://github.com/oganm/homologene>

Description A wrapper for the homologene database by the National Center for Biotechnology Information ('NCBI'). It allows searching for gene homologs across species. Data in this package can be found at <ftp://ftp.ncbi.nih.gov/pub/HomoloGene/build68/>. The package also includes an updated version of the homologene database where gene identifiers and symbols are replaced with their latest (at the time of submission) version and functions to fetch latest annotation data to keep updated.

License MIT + file LICENSE

LazyData true

RoxygenNote 7.2.3

Repository <https://oganm.r-universe.dev>

RemoteUrl <https://github.com/oganm/homologene>

RemoteRef HEAD

RemoteSha 9a9f99c4b596ccdd05a1ea1d7f62323bffb3b721

Contents

| | |
|-----------------------------|-----------|
| autoTranslate | 2 |
| diopt | 3 |
| getGeneHistory | 4 |
| getGeneInfo | 5 |
| getHomologene | 5 |
| homologene | 6 |
| homologeneData | 6 |
| homologeneData2 | 7 |
| homologeneVersion | 7 |
| human2mouse | 8 |
| mouse2human | 8 |
| taxData | 9 |
| updateHomologene | 9 |
| updateIDs | 10 |
| Index | 11 |

| | |
|---------------|---|
| autoTranslate | <i>Attempt to automatically translate a gene list</i> |
|---------------|---|

Description

Given a list of query gene list and a target gene list, the function tries find the homology pairing that matches the query list to the target list. The query list is a short list of genes while the target list is supposed to represent a large number of genes from the target species. The default output will be the largest possible list. If `returnAllPossible = TRUE` then all possible pairings with any matches are returned. It is possible to limit the search by setting `possibleOrigins` and `possibleTargets`. Note that gene symbols of some species are more similar to each other than others. Using this with small gene lists and without providing any `possibleOrigins` or `possibleTargets` might return multiple hits, or if `returnAllPossible = TRUE` a wrong match can be returned.

Usage

```
autoTranslate(
  genes,
  targetGenes,
  possibleOrigins = NULL,
  possibleTargets = NULL,
  returnAllPossible = FALSE,
  db = homologene::homologeneData
)
```

Arguments

| | |
|-------------------|---|
| genes | A list of genes to match the target. Symbols or NCBI ids |
| targetGenes | The target list. This list is supposed to represent a large number of genes from the target species. |
| possibleOrigins | Taxonomic identifiers of possible origin species |
| possibleTargets | Taxonomic identifiers of possible target species |
| returnAllPossible | if TRUE returns all possible pairings with non zero gene matches. If FALSE (default) returns the best match |
| db | Homologene database to use. |

Value

A data frame if returnAllPossible = FALSE and a list of data frames if TRUE

| | |
|-------|-----------------------------|
| diopt | <i>Query DIOPT database</i> |
|-------|-----------------------------|

Description

Query DIOPT database (https://www.flyrnai.org/cgi-bin/DRSC_orthologs.pl) for orthologues. DIOPT database uses multiple tools to find gene orthologues. Sadly they don't have an API so this function queries by visiting the site and filling up the form. By default each query will take a minimum of 10 seconds due to delay parameter. This is taken from their robots.txt at the time this function is written. Note that DIOPT is not necessarily in sync with homologene database as provided in this package.

Usage

```
diopt(genes, inTax, outTax, delay = 10)
```

Arguments

| | |
|--------|--|
| genes | A vector of gene identifiers. Anything that DIOPT accepts |
| inTax | taxid of the species that the input genes are coming from |
| outTax | taxid of the species that you are seeking homology. 0 to query all species. |
| delay | How many seconds of delay should be between queries. Default is 10 based on the robots.txt at the time this function is written. |

Details

DIOPT does not support all species available in the homologene database. The supported species are:

4896 Schizosaccharomyces pombe

4932 Saccharomyces cerevisiae

6239 Caenorhabditis elegans

7227 Drosophila melanogaster

7955 Danio rerio

8364 Xenopus (Silurana) tropicalis

9606 Homo sapiens

10090 Mus musculus

10116 Rattus norvegicus

3702 Arabidopsis thaliana

Value

A data frame

| | |
|----------------|-----------------------------------|
| getGeneHistory | <i>Download gene history file</i> |
|----------------|-----------------------------------|

Description

Downloads and reads the gene history file from NCBI website. This file is needed for other functions

Usage

```
getGeneHistory(destfile = NULL, justRead = FALSE)
```

Arguments

| | |
|----------|--|
| destfile | Path of the output file. If NULL a temp file will be used |
| justRead | If TRUE and destfile exists, it reads the file instead of downloading the latest one from NCBI |

Value

A data frame with latest gene history information

| | |
|-------------|---|
| getGeneInfo | <i>Download gene symbol information</i> |
|-------------|---|

Description

This function downloads the gene_info file from NCBI website and returns the gene symbols for current IDs.

Usage

```
getGeneInfo(destfile = NULL, justRead = FALSE, chunk_size = 1e+06)
```

Arguments

| | |
|------------|---|
| destfile | Path of the output file. If NULL a temp file will be used |
| justRead | If TRUE and destfile exists, it reads the file instead of downloading the latest one from NCBI |
| chunk_size | Chunk size to be used with <code>link[readr]{read_tsv_chunked}</code> . The gene_info file is big enough to make its intake difficult. If you don't have large amounts of free memory you may have to reduce this number to read the file in smaller chunks |

Value

A data frame with gene symbols for each current gene id

| | |
|---------------|---------------------------------------|
| getHomologene | <i>Get the latest homologene file</i> |
|---------------|---------------------------------------|

Description

This function downloads the latest homologene file from NCBI. Note that Homologene has not been updated since 2014 so the output will be identical to [homologeneData](#) included in this package. This function is here for futureproofing purposes.

Usage

```
getHomologene(destfile = NULL, justRead = FALSE)
```

Arguments

| | |
|----------|--|
| destfile | Path of the output file. If NULL a temp file will be used |
| justRead | If TRUE and destfile exists, it reads the file instead of downloading the latest one from NCBI |

Value

A data frame with homology groups, gene ids and gene symbols

| | |
|------------|--------------------------------------|
| homologene | <i>Get homologues of given genes</i> |
|------------|--------------------------------------|

Description

Given a list of genes and a taxid, returns a data frame including the genes and their corresponding homologues

Usage

```
homologene(genes, inTax, outTax, db = homologene::homologeneData)
```

Arguments

| | |
|--------|---|
| genes | A vector of gene symbols or NCBI ids |
| inTax | taxid of the species that the input genes are coming from |
| outTax | taxid of the species that you are seeking homology |
| db | Homologene database to use. |

Examples

```
homologene(c('Eno2', '17441'), inTax = 10090, outTax = 9606)
```

| | |
|----------------|-----------------------|
| homologeneData | <i>homologeneData</i> |
|----------------|-----------------------|

Description

List of gene homologues used by homologene functions

Usage

```
homologeneData
```

Format

An object of class `data.frame` with 275237 rows and 4 columns.

| | |
|-----------------|------------------------|
| homologeneData2 | <i>homologeneData2</i> |
|-----------------|------------------------|

Description

A modified copy of the homologene database. Homologene was updated at 2014 and many of its gene IDs and symbols are out of date. Here the IDs and symbols are replaced with their most current version Last update: Tue Oct 31 18:41:52 2023

Usage

homologeneData2

Format

An object of class data.frame with 266573 rows and 4 columns.

| | |
|-------------------|-----------------------------------|
| homologeneVersion | <i>Version of homologene used</i> |
|-------------------|-----------------------------------|

Description

Version of homologene used

Usage

homologeneVersion

Format

An object of class integer of length 1.

`human2mouse`*Human/mouse wrapper for homologue*

Description

Human/mouse wrapper for homologue

Usage

```
human2mouse(genes, db = homologue::homologueData)
```

Arguments

| | |
|--------------------|--------------------------------------|
| <code>genes</code> | A vector of gene symbols or NCBI ids |
| <code>db</code> | Homologue database to use. |

Examples

```
human2mouse(c('EN02', '4340'))
```

`mouse2human`*Mouse/human wrapper for homologue*

Description

Mouse/human wrapper for homologue

Usage

```
mouse2human(genes, db = homologue::homologueData)
```

Arguments

| | |
|--------------------|--------------------------------------|
| <code>genes</code> | A vector of gene symbols or NCBI ids |
| <code>db</code> | Homologue database to use. |

Examples

```
mouse2human(c('Eno2', '17441'))
```

| | |
|---------|--|
| taxData | <i>Names and ids of included species</i> |
|---------|--|

Description

Names and ids of included species

Usage

```
taxData
```

Format

An object of class `data.frame` with 21 rows and 2 columns.

| | |
|------------------|-----------------------------------|
| updateHomologene | <i>Update homologene database</i> |
|------------------|-----------------------------------|

Description

Creates an updated version of the homologene database. This is done by downloading the latest gene annotation information and tracing changes in gene symbols and identifiers over history. [homologeneData2](#) was created using this function over the original [homologeneData](#). This function requires downloading large amounts of data from the NCBI ftp servers.

Usage

```
updateHomologene(
  destfile = NULL,
  baseline = homologene::homologeneData2,
  gene_history = NULL,
  gene_info = NULL
)
```

Arguments

| | |
|--------------|---|
| destfile | Optional. Path of the output file. |
| baseline | The baseline homologene file to be used. By default uses the homologeneData2 that is included in this package. The more ids to update, the more time is needed for the update which is why the default option uses an already updated version of the original database. |
| gene_history | A gene history data frame, possibly returned by getGeneHistory function. Use this if you want to have a static gene_history file to update up to a specific date. An up to date gene_history object can be set to update to a specific date by trimming rows that have recent dates. Note that the same is not possible for the gene_info. If not provided, the latest file will be downloaded. |

gene_info A gene info data frame that contains ID-symbol matches, possibly returned by [getGeneInfo](#). Use this if you want a static version. Should be in sync with the gene_history file. Note that there is no easy way to track changes in gene symbols back in time so if you want to update it up to a specific date, make sure you don't lose that file.

Value

Homologene database in a data frame with updated gene IDs and symbols

| | |
|-----------|------------------------|
| updateIDs | <i>Update gene IDs</i> |
|-----------|------------------------|

Description

Given a list of gene ids and gene history information, traces changes in the gene's name to get the latest valid ID

Usage

```
updateIDs(ids, gene_history)
```

Arguments

ids Gene ids
gene_history Gene history information, probably returned by [getGeneHistory](#)

Value

A character vector. New ids for genes that changed ids, or "-" for discontinued genes. the input itself.

Examples

```
## Not run:  
gene_history = getGeneHistory()  
updateIDs(c("4340964", "4349034", "4332470", "4334151", "4323831"),gene_history)  
  
## End(Not run)
```

Index

* datasets

- homologeneData, [6](#)
- homologeneData2, [7](#)
- homologeneVersion, [7](#)
- taxData, [9](#)

autoTranslate, [2](#)

diopt, [3](#)

getGeneHistory, [4](#), [9](#), [10](#)

getGeneInfo, [5](#), [10](#)

getHomologene, [5](#)

homologene, [6](#)

homologeneData, [5](#), [6](#), [9](#)

homologeneData2, [7](#), [9](#)

homologeneVersion, [7](#)

human2mouse, [8](#)

mouse2human, [8](#)

taxData, [9](#)

updateHomologene, [9](#)

updateIDs, [10](#)