

Package ‘igblastr’

May 14, 2026

Title User-friendly R Wrapper to IgBLAST

Description The igblastr package provides functions to conveniently install and use a local IgBLAST installation from within R. The package also includes a set of built-in IgBLAST-compatible germline databases from OGRDB, the AIRR Community’s Open Germline Receptor Database, for various organisms. It provides functions to create additional IgBLAST-compatible germline databases using reference sequences retrieved from IMGT/V-QUEST or local FASTA files supplied by the user. When possible, the FWR/CDR boundaries on the V alleles (a.k.a “internal data”) are computed and stored in the germline database, so can be used as a replacement for the internal data shipped with IgBLAST. IgBLAST is described at <<https://pubmed.ncbi.nlm.nih.gov/23671333/>>. IgBLAST web interface: <<https://www.ncbi.nlm.nih.gov/igblast/>>. OGRDB: <<https://ogrdb.airr-community.org/>>. IMGT/V-QUEST download site: <<https://www.imgt.org/download/V-QUEST/>>.

biocViews Immunology, Immunogenetics, ImmunoOncology, CellBiology

URL <https://bioconductor.org/packages/igblastr>

BugReports <https://github.com/HyrienLab/igblastr/issues>

Version 1.3.1

License Artistic-2.0

Encoding UTF-8

Depends R (>= 4.2.0), tibble, Biostrings

Imports methods, utils, stats, tools, R.utils, curl, httr, xml2, rvest, xtable, jsonlite, S4Vectors, IRanges, GenomeInfoDb

Suggests GenomicAlignments, parallel, testthat, knitr, rmarkdown, BiocStyle, ggplot2, dplyr, scales, ggseqlogo

VignetteBuilder knitr

Collate utils.R internet-utils.R long_to_wide_airr.R igdata-IO.R translate_codons.R allele2gene.R parse_imgt_fasta_headers.R file-utils.R loci-utils.R db-utils.R ndm_data-IO.R compute_V_gene_delineations.R clean_allele_set.R auxdata-IO.R compute_auxdata.R V_alleles-inspect.R J_alleles-inspect.R LATIN_NAMES.R IMGT-utils.R IMGT-c_region-utils.R precompiled-igblast-utils.R cache-utils.R get_igblast_root.R edit_imgt_file.R igblast_info.R update_live_igdata.R intndata-utils.R auxdata-utils.R install_igblast.R

make_blastdbs.R create_region_db.R create_germline_db.R
 create_c_region_db.R reset_germline_dbs.R list_germline_dbs.R
 use_germline_db.R reset_c_region_dbs.R list_c_region_dbs.R
 use_c_region_db.R install_custom_germline_db.R OGRDB-utils.R
 OGRDB-API.R download_OGRDB_germline_sequences.R
 download_OGRDB_germline_json.R install_OGRDB_germline_db.R
 download_IMGT_germline_sequences.R install_IMGT_germline_db.R
 augment_germline_db.R prepare_igblastn_cmdline_args.R
 read_igblastn_fmt7_output.R read_igblastn_AIRR_output.R
 igblastn.R igbrowser.R summarizeMismatches.R OAS-utils.R zzz.R

git_url <https://git.bioconductor.org/packages/igblastr>

git_branch devel

git_last_commit ab6ea27

git_last_commit_date 2026-05-12

Repository Bioconductor 3.24

Date/Publication 2026-05-13

Author Hervé Pagès [aut, cre] (ORCID: <<https://orcid.org/0009-0002-8272-4522>>),
 Ollivier Hyrien [aut, fnd] (ORCID:
 <<https://orcid.org/0000-0003-1909-2542>>),
 Kellie MacPhee [ctb] (ORCID: <<https://orcid.org/0009-0008-0993-4009>>),
 Michael Duff [ctb] (ORCID: <<https://orcid.org/0009-0008-4279-0756>>),
 Jason Taylor [ctb]

Maintainer Hervé Pagès <hpages.on.github@gmail.com>

Contents

| | |
|---|----|
| allele2gene | 3 |
| augment_germline_db | 4 |
| auxdata-IO | 7 |
| auxdata-utils | 8 |
| compute_auxdata | 11 |
| compute_V_gene_delineations | 13 |
| download_IMGT_germline_sequences | 16 |
| download_OGRDB_germline_json | 18 |
| download_OGRDB_germline_sequences | 22 |
| get_igblast_root | 25 |
| igblastn | 26 |
| igblastr_usage_report | 32 |
| igblast_info | 33 |
| IGBLAST_ROOT | 35 |
| igbrowser | 36 |
| install_custom_germline_db | 37 |
| install_igblast | 42 |
| install_IMGT_germline_db | 43 |
| intdata-utils | 46 |
| J_alleles-inspect | 49 |
| list_c_region_dbs | 51 |
| list_germline_dbs | 53 |
| ndm_data-IO | 55 |

| | |
|-------------------------------------|----|
| OAS-utils | 57 |
| parse_imgt_fasta_headers | 60 |
| read_igblastn_AIRR_output | 63 |
| read_igblastn_fmt7_output | 64 |
| reset_c_region_dbs | 66 |
| reset_germline_dbs | 67 |
| summarizeMismatches | 68 |
| translate_codons | 68 |
| update_live_igdata | 70 |
| use_c_region_db | 72 |
| use_germline_dbs | 74 |
| V_alleles-inspect | 75 |

| | |
|--------------|-----------|
| Index | 77 |
|--------------|-----------|

| | |
|-------------|--|
| allele2gene | <i>Go from germline gene allele names to germline gene names</i> |
|-------------|--|

Description

A simple convenience function to remove the allele suffix from a vector of germline gene allele names.

Usage

```
allele2gene(allele_names)
```

Arguments

`allele_names` A character vector of germline gene allele names.

Details

Germline gene allele names use specific nomenclature, often including a gene name followed by an asterisk and allele number, like in IGHD3-16*03.

The `allele2gene()` function simply removes the asterisk and anything that follows it to keep the gene name only. Note that the function is *vectorized*, that is, its input can be a *vector* of germline gene allele names, in which case the corresponding germline gene names are returned in a vector of the same length as the input vector. The names on the input vector are propagated, and so are the NAs in it.

Value

A character vector *parallel* to the input vector.

See Also

- [load_germline_db](#) to load the nucleotide sequences of the gene regions stored in a cached germline db.
- [load_c_region_db](#) to load the nucleotide sequences of the gene regions stored in a cached C-region db.
- [intdata_utils](#) to access IgBLAST *internal data*.

Examples

```
allele2gene(c("IGHV1-2*04", "IGHV1-2*06", "IGHV5-51*01", "IGHD3-16*03"))

J_alleles <- load_germline_db("_OGRDB.human.IGH+IGK+IGL.202410",
                             region_types="J")

names(J_alleles)
allele2gene(names(J_alleles))

C_alleles <- load_c_region_db("_IMGT.human.IGH+IGK+IGL.202412")
names(C_alleles)
allele2gene(names(C_alleles))
```

augment_germline_db *Add novel gene alleles to a germline db*

Description

WARNING: Some shortcomings were identified in the design of the `augment_germline_db_[VDJ]()` functions so they are now deprecated. Please do not use them. They will be replaced with a better alternative in future versions of the package.

Three functions to add novel V, D, or J gene alleles to a germline db.

Note that these functions can also be used to combine germline databases from two different organisms. See "COMBINE GERMLINE DATABASES FROM TWO ORGANISMS" in the Examples section below for how to do this.

Usage

```
augment_germline_db_V(db_name, novel_alleles,
                      destdir=".", overwrite=FALSE, verbose=FALSE)

augment_germline_db_D(db_name, novel_alleles,
                      destdir=".", overwrite=FALSE, verbose=FALSE)

augment_germline_db_J(db_name, novel_alleles,
                      destdir=".", overwrite=FALSE, verbose=FALSE)
```

Arguments

| | |
|---------------|---|
| db_name | A single string that is the name of the cached germline db that contains the set of gene alleles to augment. Use <code>list_germline_dbs()</code> to list the cached germline dbs. The exact function used (i.e. <code>augment_germline_db_V()</code> , <code>augment_germline_db_D()</code> , or <code>augment_germline_db_J()</code>) determines the set of alleles to augment (i.e. alleles from the V, D, or J region). |
| novel_alleles | A single string that is the path to a FASTA file (possibly gz-compressed) where the novel alleles are stored. Alternatively, the novel alleles can be supplied as a <i>named</i> <code>DNAStrngSet</code> object. |
| destdir | A single string that is the path to the "destination directory", that is, the directory where the augmented V-, D-, or J-region db is to be created. This directory will be created if it doesn't exist already. Note that, by default, the augmented region db will be created in the current directory. |

| | |
|-----------|---|
| overwrite | If the "destination directory" already contains a V-, D-, or J-region db, should it be overwritten? |
| verbose | Set to TRUE to have the function display some details about its internal operations. |

Value

These functions don't return anything (invisible NULL).

See Also

- The `igblastn` function to run the `igblastn` *standalone executable* included in IgBLAST from R. This is the main function in the **igblastr** package.
- `list_germline_dbs` to list the cached germline dbs.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## Not run:
if (!has_igblast()) install_igblast()

query <- system.file(package="igblastr", "extdata",
                     "BCR", "heavy_sequences.fasta")

use_c_region_db("_IMGT.human.IGH+IGK+IGL.202412")

## -----
## USE HUMAN GERMLINE DATABASE FROM AIRR
## -----

use_germline_db("_OGRDB.human.IGH+IGK+IGL.202410")

AIRR_df <- igblastn(query)

## -----
## ADD NOVEL V ALLELES
## -----

## 'fake_human_V_alleles.fasta' contains made-up novel V alleles:
## - 2 novel alleles for gene IGHV1-8: IGHV1-8*fake1, IGHV1-8*fake2
## - 1 novel allele for gene IGHV4-61: IGHV4-61*fake
my_novel_V_alleles <- system.file(package="igblastr", "extdata",
                                  "novel_germline_alleles",
                                  "fake_human_V_alleles.fasta")

## Take a quick look at these novel V alleles:
readDNAStrngSet(my_novel_V_alleles)

## Create a new V germline database that combines the V alleles
## from _OGRDB.human.IGH+IGK+IGL.202410 with our novel V alleles:
my_Vdb_path <- file.path(tempdir(), "myVdb")
augment_germline_db_V("_OGRDB.human.IGH+IGK+IGL.202410",
                     my_novel_V_alleles,
                     destdir=my_Vdb_path)
```

```

## To use this new augmented V germline database with igblastn(),
## supply its path via the 'germline_db_V' argument:
AIRR_df2 <- igblastn(query, germline_db_V=my_Vdb_path)

## -----
## A QUICK COMPARISON BETWEEN 'AIRR_df' AND 'AIRR_df2'
## -----

## Index of rows where "v_call" has changed between 'AIRR_df'
## and 'AIRR_df2':
idx <- which(AIRR_df$v_call != AIRR_df2$v_call)
idx # 2 rows

AIRR_df[idx, c("v_call", "v_cigar", "v_identity")]

AIRR_df2[idx, c("v_call", "v_cigar", "v_identity")]

## Besides these 2 rows, all the other rows are the same:
stopifnot(all.equal(AIRR_df[-idx, ], AIRR_df2[-idx, ]))

## -----
## COMBINE GERMLINE DATABASES FROM TWO ORGANISMS
## -----

## The augment_germline_db_[VDJ]() functions can be used to combine
## germline databases from two different organisms. This can be useful
## for example when working with BCR sequences from mice that have been
## engineered to have both mouse and some human immunoglobulin genes.
##
## To create a hybrid human/mouse V germline database, we can either:
##
## (1) Add all (or a subset of) mouse V alleles to all human V alleles.
##     This is done by extracting mouse V germline allele sequences from
##     a cached germline database and using them to augment a cached
##     germline database for human.
##
## (2) Add all (or a subset of) human V alleles to all mouse V alleles.
##     This is done by extracting human V germline allele sequences from
##     a cached germline database and using them to augment a cached
##     germline database for mouse.
##
## Note that:
## - We can choose to subset or not the V germline allele sequences
##   extracted from one V germline database before adding them to the
##   other V germline database.
## - The two approaches above are equivalent if we don't subset, that
##   is, if we combine **all** human V alleles with **all** mouse V
##   alleles.
## - However if our engineered mice only have a small known subset of
##   human immunoglobulin genes (e.g. IGHV1-2), then we might want to
##   create a hybrid human/mouse germline database that only adds the
##   human alleles for genes IGHV1-2 to the mouse V alleles. In this
##   case we need to use (2).

## Let's do (2):

db_name1 <- "_OGRDB.mouse.PWD_PhJ.IGH+IGK+IGL.202410"

```

```

db_name2 <- "_OGRDB.human.IGH+IGK+IGL.202410"

## Extract human V germline alleles:
human_V_alleles <- load_germline_db(db_name2, "V")

## Subset to keep only alleles for genes IGHV1-2:
idx <- grep("^IGHV[12]", names(human_V_alleles))
human_V12_alleles <- human_V_alleles[idx]

## Create a new V germline database that combines the mouse V
## alleles from 'db_name1' with the alleles in 'human_V12_alleles':
engmouseVdb_path <- file.path(tempdir(), "engmouseVdb")
augment_germline_db_V(db_name1, human_V12_alleles,
                      destdir=engmouseVdb_path)

## End(Not run)

## Then, assuming that 'query' contains BCR sequences from the
## engineered mice:
## Not run:
use_germline_db(db_name1)
use_c_region_db("_IMGT.mouse.IGH.202509")
igblastn(query, germline_db_V=engmouseVdb_path, ...)

## End(Not run)

## Note that, by default, the mouse-only D and J databases that we
## selected above with 'use_germline_db(db_name1)' are being used.
## If we also want to create hybrid D and J databases, we need
## to repeat the above steps for each of them. Then we need to
## specify the paths to the 3 hybrid databases when we call igblastn():
## Not run:
igblastn(query, germline_db_V=engmouseVdb_path,
          germline_db_D=engmouseDdb_path,
          germline_db_J=engmouseJdb_path,
          ...)

## End(Not run)

```

auxdata-IO

Read/write IgBLAST auxiliary data files

Description

IgBLAST auxiliary data files (a.k.a. .aux files) are tab-delimited text files used by IgBLAST to annotate germline J gene allele sequences.

The **igblastR** package provides low-level functions to read/write this type of file.

Usage

```

read_auxdata(filepath)
write_auxdata(auxdata, file="")

```

Arguments

| | |
|----------|---|
| filepath | The path to the IgBLAST auxiliary data file to read. |
| auxdata | A data.frame with 1 row per germline J gene allele sequence and the same columns as the data.frame returned by read_auxdata(). See Value section below for the details. |
| file | The path to the IgBLAST auxiliary data file to write. |

Value

read_auxdata() returns a data.frame with 1 row per germline J gene allele sequence and the following columns:

1. allele_name: allele name;
2. coding_frame_start: first coding frame start position (0-based);
3. chain_type: chain type as a 2-letter code e.g. JK for "J allele from the Kappa locus" (BCR locus) or JG for "J allele from the Gamma locus" (TCR locus);
4. cdr3_end: CDR3 end position (0-based);
5. extra_bps: extra base pairs beyond J coding end.

write_auxdata() doesn't return anything (i.e. invisible NULL).

See Also

- [compute_auxdata](#) to annotate a set of germline J gene allele sequences.
- [extract_auxdata_from_ogrdb_json](#) to extract the coding frame and CDR3 end information for the J alleles annotated in an AIRR-C JSON file.
- [auxdata_utils](#) to access IgBLAST *auxiliary data*.
- The [igblastn](#) function to run the *igblastn standalone executable* included in IgBLAST from R. This is the main function in the **igblastR** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## COMING SOON...
```

| | |
|---------------|--------------------------------------|
| auxdata-utils | <i>Access IgBLAST auxiliary data</i> |
|---------------|--------------------------------------|

Description

IgBLAST *auxiliary data* is expected to annotate all the known germline J gene alleles for a given organism. It is provided by NCBI and is typically included in a standard IgBLAST installation.

The *auxiliary data* informs about the coding frame and CDR3/FWR4 boundary on the J allele sequences. Although this information is not strictly needed by IgBLAST, the latter needs it in order to compute some of the AIRR fields like `vj_in_frame`, `productive`, `cdr3`, `fwr4`, etc...

`get_auxdata_path()` and `load_auxdata()` can be used to access the *auxiliary data* included in IgBLAST or in one of the *cached germline dbs* managed by **igblastR**.

We also provide `compute_germline_db_auxdata()` to compute the *auxiliary data* associated with the J alleles in a cached germline dbs managed by **igblastR**.

Usage

```
## Access auxiliary data:
get_auxdata_path(organism, which=c("live", "original"))
load_auxdata(organism, which=c("live", "original"))

## Compute auxiliary data associated with the J alleles in germline db:
compute_germline_db_auxdata(db_name, ...)
```

Arguments

| | |
|----------|---|
| organism | A single string containing the name of an organism as returned by <code>list_igblast_organisms()</code> . Alternatively, this can be the name of a cached germline db. Note that this works only for germline dbs that include their own <i>auxiliary data</i> . You can use: <code>list_germline_dbs(with.auxdata.only=TRUE)</code> to list them. See <code>?list_germline_dbs</code> for more information. |
| which | By default, <code>get_auxdata_path()</code> and <code>load_auxdata()</code> access the "live IgBLAST data", that is, the IgBLAST data that the user has possibly updated with <code>update_live_igdata()</code> . Depending on whether updates were applied or not, the "live IgBLAST data" might differ from the original IgBLAST data. Set which to "original" if you want to access the original IgBLAST data instead. See <code>?update_live_igdata</code> for more information about "live" and "original" IgBLAST data. |
| db_name | A single string specifying the name of a cached germline db. Use <code>list_germline_dbs()</code> to list all the cached germline dbs. |
| ... | Extra arguments to be passed to the internal call to <code>compute_auxdata()</code> . See <code>?compute_auxdata</code> for what arguments are supported. |

Details

IgBLAST *auxiliary data* is typically included in a standard IgBLAST installation. It's located in the `optional_file/` directory which is itself a subdirectory of IgBLAST *root directory*.

The data consists of one tabulated file per organism. Each file indicates the germline J gene coding frame start position, the J gene type, and the CDR3 end position for all known germline J allele sequences. See <https://ncbi.github.io/igblast/cook/How-to-set-up.html> for additional details.

IMPORTANT NOTE: Using IgBLAST with *auxiliary data* that is not compatible with the germline J gene sequences in the germline db can cause it to return improper frame status or CDR3 information (other returned information will still be correct).

Value

`get_auxdata_path()` returns a single string containing the path to the *auxiliary data* included in the IgBLAST installation used by **igblastr**, for the specified organism. Not necessarily suitable to use with `igblastn()` (see WARNING below).

`load_auxdata()` returns the *auxiliary data* in a data.frame with 1 row per supplied J allele sequence and the same columns as the data.frame returned by `read_auxdata()`. See `?read_auxdata` for more information.

`compute_germline_db_auxdata()` returns the computed *auxiliary data* in a data.frame with 1 row per J allele in the germline db and the same columns as the data.frame returned by `compute_auxdata()` or `read_auxdata()`. See `?compute_auxdata` for more information.

See Also

- `compute_auxdata` to annotate a set of germline J gene allele sequences.
- `intdata_utils` to access IgBLAST *internal data*.
- `update_live_igdata` for more information about "live" and "original" IgBLAST data.
- `list_igblast_organisms` to list the organisms supported by IgBLAST.
- `list_germline_dbs` to list the cached germline dbs.
- <https://ncbi.github.io/igblast/cook/How-to-set-up.html> for important information about IgBLAST *auxiliary data*.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastR** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

## -----
## Access IgBLAST auxiliary data for a given organism
## -----

## IgBLAST only includes auxiliary data for the following organisms:
list_igblast_organisms()

get_auxdata_path("human")

human_auxdata <- load_auxdata("human")
head(human_auxdata)

## -----
## Access the auxiliary data included in a germline db
## -----

## List germline dbs with auxiliary data:
list_germline_dbs(with_auxdata.only=TRUE)

## Access the auxiliary data included in germline db
## _OGRDB.rhesus_monkey.IGH+IGK+IGL.202602 (note that this data was
## extracted from the AIRR-C JSON files provided by OGRDB for rhesus
## monkey, see '?extract_auxdata_from_ogrdb_json' for more information):
db_name <- "_OGRDB.rhesus_monkey.IGH+IGK+IGL.202602"

get_auxdata_path(db_name)

rhesus_monkey_auxdata <- load_auxdata(db_name)
head(rhesus_monkey_auxdata)

## -----
## Compute the auxiliary data associated with the J alleles in a
## germline db
```

```
## -----

## The auxiliary data included in germline db
## _OGRDB.rhesus_monkey.IGH+IGK+IGL.202602 can be computed
## with compute_germline_db_auxdata():
auxdata2 <- compute_germline_db_auxdata(db_name)

## Replace the negative "cdr3_end" value with NA:
replace_idx <- which(auxdata2[, "cdr3_end"] < 0L)
auxdata2[replace_idx, "cdr3_end"] <- NA

## Sanity check:
stopifnot(identical(auxdata2, rhesus_monkey_auxdata))
```

compute_auxdata

Compute IgBLAST auxiliary data

Description

A utility function to annotate a set of germline J gene allele sequences in a way similar to how they are annotated in *IgBLAST auxiliary data*.

Note that the annotation produced by the function can be used by *IgBLAST* as a substitute to the *auxiliary data* shipped with the *IgBLAST* software (and typically included in a standard *IgBLAST* installation).

See [?load_auxdata](#) for more information about *IgBLAST auxiliary data*.

Usage

```
compute_auxdata(J_alleles, codon_starts=NULL, no.warnings=FALSE)
```

Arguments

- | | |
|--------------|--|
| J_alleles | A DNAStringSet object containing germline J gene allele sequences. |
| codon_starts | When supplied, must be a named integer vector <i>parallel</i> to J_alleles, that is, it must have the same length and names as J_alleles. For each allele in J_alleles, codon_starts should indicate the start position (1-based) of the first codon in the allele. This should be a number ≥ 1 and ≤ 3 . NAs are allowed when the coding frame is not known. Note that the "coding frame start" is 0-based so can simply be inferred from the "codon start" by subtracting 1. compute_auxdata() will use the supplied "codon starts" to either validate the "coding frame starts" that it could determine, or to set the "coding frame starts" that it could not determine. |
| no.warnings | Set to TRUE to suppress the warnings that compute_auxdata() otherwise issues when the returned data.frame has NAs in its coding_frame_start column or negative values in its cdr3_end column. |

Details

The FWR4 region is expected to start with the following amino acid motifs (X represents any amino acid):

- "WGXG" on the heavy chain;
- "FGXG" on the light chain.

compute_auxdata() searches for the "WGXG" and "FGXG" motifs in the supplied allele sequences to determine the start of their FWR4 region. From there it can easily infer the cdr3_end, coding_frame_start, and extra_bps columns.

Note that the function will emit a warning if the start of the FWR4 region (and therefore the CDR3 end) could not be found for some alleles, or if a stop codon was found in some alleles.

Value

Returns the computed *auxiliary data* in a data.frame with 1 row per supplied J allele sequence and the same columns as the data.frame returned by read_auxdata(). See ?read_auxdata for more information.

See Also

- [load_auxdata](#) to access IgBLAST *auxiliary data*.
- [compute_V_gene_delineations](#) to annotate a set of germline V gene allele sequences.
- <https://ncbi.github.io/igblast/cook/How-to-set-up.html> for important information about IgBLAST *auxiliary data*.
- [DNAStringSet](#) objects in the **Biostrings** package.
- The [igblastn](#) function to run the *igblastn standalone executable* included in IgBLAST from R. This is the main function in the **igblastR** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## -----
## BASIC EXAMPLE
## -----

## Let's load a set of human J allele sequences:
db_name <- "_OGRDB.human.IGH+IGK+IGL.202410"
J_alleles <- load_germline_db(db_name, region_types="J")
J_alleles # DNAStringSet object

computed_auxdata <- compute_auxdata(J_alleles)
head(computed_auxdata)

## -----
## SANITY CHECK
## -----

## Note that 'computed_auxdata' is in agreement with the auxiliary
## data included in IgBLAST for human, except for alleles IGHJ6*02
## and IGHJ6*03 (the 'extra_bps' column contains incorrect values
## for these two alleles, 0's instead of 1's):
human_auxdata0 <- load_auxdata("human", which="original")
```

```

bad_alleles <- c("IGHJ6*02", "IGHJ6*03")
subset(human_auxdata0, allele_name %in% bad_alleles)

## We manually correct this:
fixme <- human_auxdata0[ , "allele_name"] %in% bad_alleles
human_auxdata0[fixme, "extra_bps"] <- 1L # replace 0L with 1L

## Now the data in 'computed_auxdata' matches exactly the corresponding
## data in 'human_auxdata0':
m <- match(computed_auxdata[ , "allele_name"],
           human_auxdata0[ , "allele_name"])
human_auxdata <- human_auxdata0[m, ]
rownames(human_auxdata) <- NULL
stopifnot(identical(computed_auxdata, human_auxdata))

```

```

compute_V_gene_delineations
      Compute V gene delineations

```

Description

A utility function to compute the FWR/CDR boundaries (a.k.a. V gene delineations) for a set of germline V gene allele sequences. Note that the input sequences must be *gapped*, that is, they must contain the so-called "IMGT gaps" whose purpose is to convey the delineation information.

Usage

```
compute_V_gene_delineations(gapped_V_alleles, as.IRangesList=FALSE)
```

Arguments

gapped_V_alleles
A [DNAStrngSet](#) object containing germline V gene allele *gapped* sequences.

as.IRangesList Experts only! Set to TRUE to return the delineations in an [IRangesList](#) object instead of a data.frame. In this case, the returned IRangesList object will have 1 list element per supplied V allele sequence and will carry the following metadata columns (accessible with [mcols\(\)](#)): `seq_len`, `coding_frame_start`, `starting_gap`, `all_gaps_in_frame`, `all_gaps_contained`. See Value section below for what these columns contain.

Details

See <https://www.imgt.org/IMGTScientificChart/Numbering/IMGTIGVLSuperfamily.html> and https://www.imgt.org/IMGTScientificChart/Numbering/IMGT-Kabat_part1.html for how gaps are used in the context of "IMGT unique numbering for V-DOMAIN and V-LIKE-DOMAIN".

Note that the annotation returned by `compute_V_gene_delineations()` can be used by IgBLAST as a substitute to the *internal data* shipped with the IgBLAST software (and typically included in a standard IgBLAST installation).

See [?load_intdata](#) for more information about IgBLAST *internal data*.

Value

Returns the computed *internal data* in a data.frame with 1 row per supplied V allele sequence and a set of columns similar to the data.frame returned by `read_ndm_data()`. With the following differences:

- The data.frame returned by `compute_V_gene_delineations()` has no `chain_type` column.
- The data.frame returned by `compute_V_gene_delineations()` has the following additional columns:
 1. `seq_len`: length of the sequence after removal of the gaps;
 2. `starting_gap`: size (in number of nucleotides) of gap block located at the very beginning of the sequence if any (0 if no such block);
 3. `all_gaps_in_frame`: indicates whether the gap blocks align with the underlying coding frame or not;
 4. `all_gaps_contained`: TRUE if the gap blocks don't cross the FWR/CDR boundaries, and FALSE otherwise.

See Also

- <https://www.imgt.org/IMGTScientificChart/Numbering/IMGTIGVLsuperfamily.html> for how the "IMGT unique numbering for V-DOMAIN and V-LIKE-DOMAIN" works.
- `load_intdata` to access IgBLAST *internal data*.
- `compute_auxdata` to annotate a set of germline J gene allele sequences.
- `DNAStrngSet` objects in the **Biostrings** package.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastR** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## -----
## 1. USING HUMAN GERMLINE V GENE ALLELE SEQUENCES FROM OGRDB
## -----
## Visit https://ogrdb.airr-community.org/germline_sets/Homo%20sapiens
## for the list of human germline sets available at OGRDB.

## Download OGRDB germline gene allele sequences for human locus IGK
## to a temporary directory:
germline_set <- c(`IGKappa_VJ`=4)

dir.create(tmpdir <- tempfile())
download_OGRDB_germline_sequences("Homo sapiens", germline_set,
                                  destdir=tmpdir)

## List the downloaded files:
list.files(tmpdir) # 2 FASTA files (one per IMGT group)

## The V allele sequences are **gapped**:
IGKV_alleles <- readDNAStrngSet(file.path(tmpdir, "IGKV.fasta"))
as.character(IGKV_alleles[1:3])

## Compute the V gene delineations for 'IGKV_alleles':
IGKV_gene_delinea <- compute_V_gene_delineations(IGKV_alleles)
```

```

head(IGKV_gene_delinea)

## As a sanity check, we're going to compare 'IGKV_gene_delinea' with
## the V gene delineations provided by the corresponding AIRR-C JSON
## file from OGRDB.

download_OGRDB_germline_json("Homo sapiens", germline_set,
                             destdir=tmpdir)

list.files(tmpdir) # 1 JSON file

json_path <- file.path(tmpdir, "IGKappa_VJ.json")
IGKV_gene_delinea2 <- extract_intdata_from_ogrdb_json(json_path)
head(IGKV_gene_delinea2)

## Keep only the shared columns:
shared_colnames <- intersect(colnames(IGKV_gene_delinea),
                             colnames(IGKV_gene_delinea2))
IGKV_gene_delinea <- IGKV_gene_delinea[ , shared_colnames]
IGKV_gene_delinea2 <- IGKV_gene_delinea2[ , shared_colnames]

## Let's use same_alleles_annot(), a small undocumented helper function,
## to verify that 'IGKV_gene_delinea' and 'IGKV_gene_delinea2' contain
## the same delineation data (possibly with their rows in different
## order):
stopifnot(same_alleles_annot(IGKV_gene_delinea, IGKV_gene_delinea2))

## Remove temporary directory:
unlink(tmpdir, recursive=TRUE)

## -----
## 2. USING MOUSE GERMLINE V GENE ALLELE SEQUENCES FROM OGRDB
## -----
## Visit https://ogrdb.airr-community.org/germline\_sets/Mus%20musculus
## for the list of mouse germline sets available at OGRDB.

## Download OGRDB germline V gene allele sequences for locus IGL of
## mouse strain MSM/MsJ to a temporary directory:
germline_set <- c(`MSM/MsJ IGLV`=1)

dir.create(tmpdir <- tempfile())
download_OGRDB_germline_sequences("Mus musculus", germline_set,
                                 destdir=tmpdir)

## List the downloaded files:
list.files(tmpdir) # 1 FASTA file

## Compute the V gene delineations for 'IGKV_alleles':
IGLV_alleles <- file.path(tmpdir, "IGLV.fasta")
IGLV_gene_delinea <- compute_V_gene_delineations(IGLV_alleles)
IGLV_gene_delinea

## Note that for the first two alleles, the length of the ungapped
## sequence (seq_len) is smaller than the FWR3 end position (fwr3_end):
IGLV_gene_delinea[ , c("fwr3_end", "seq_len")]

## In other words, the sequences of alleles IGLV0-37NY*00 and

```

```

## IGLV0-COOR*00 end before the end of their FWR3.

## This is reflected in the V gene delineations provided by the
## corresponding AIRR-C JSON file from OGRDB:
download_OGRDB_germline_json("Mus musculus", germline_set,
                             destdir=tmpdir)

list.files(tmpdir)
json_path <- file.path(tmpdir, "IGLV.json")
IGLV_gene_delinea2 <- extract_intdata_from_ogrdb_json(json_path)
IGLV_gene_delinea2

## As you can see, when a V allele sequence ends before the end of
## its FWR3, the FWR3 end position reported in the AIRR-C JSON file
## is the position of the last nucleotide in the (ungapped) sequence.

## Remove temporary directory:
unlink(tmpdir, recursive=TRUE)

```

download_IMGT_germline_sequences

Download germline sequences from IMGT

Description

The `download_IMGT_germline_sequences()` function downloads germline V/D/J gene allele sequences (nucleotides) from the IMGT/V-QUEST download site for a given IMGT/V-QUEST release and organism. The sequences are downloaded in several FASTA files, one per *IMGT group*.

CONDITIONS OF USE AND LICENSE: The IMGT data is provided to the academic users and NPO's (Not for Profit Organization(s)) under the CC BY-NC-ND 4.0 license. See <https://creativecommons.org/licenses/by-nc-nd/4.0/>. Any other use of IMGT material, from the private sector, needs a financial arrangement with CNRS.

Usage

```

download_IMGT_germline_sequences(release, organism="Homo sapiens",
                                tcr.db=FALSE,
                                destdir=".", overwrite=FALSE,
                                recache=FALSE, ...)

```

```

## Related utilities:
list_IMGT_releases(recache=FALSE)
list_IMGT_organisms(release)
IMGT_is_up()

```

Arguments

| | |
|----------|--|
| release | A single string specifying the IMGT/V-QUEST release to download the germline sequences from (or to list the organisms from for <code>list_IMGT_organisms()</code>). Use <code>list_IMGT_releases()</code> to list all releases. |
| organism | A single string specifying the latin name of the organism for which to download the germline sequences. |

| | |
|-----------|---|
| tcr.db | <p>Should the germline sequences to download be B-cell Receptor (a.k.a. BCR) sequences or T-cell Receptor (a.k.a. TCR) sequences?</p> <p>By default (i.e. when tcr.db is omitted or set to FALSE), the V/D/J allele sequences for the three BCR loci (IGH, IGK, IGL) will be downloaded. This will result in the 7 following FASTA files being downloaded:</p> <pre style="text-align: center;">IGHV.fasta, IGHD.fasta, IGHJ.fasta, IGKV.fasta, IGKJ.fasta, IGLV.fasta, IGLJ.fasta</pre> <p>When tcr.db is set to TRUE, the V/D/J allele sequences for the four TCR loci (TRA, TRB, TRG, TRD) will be downloaded instead. This will result in the 10 following FASTA files being downloaded:</p> <pre style="text-align: center;">TRAV.fasta, TRAJ.fasta, TRBV.fasta, TRBD.fasta, TRBJ.fasta, TRGV.fasta, TRGJ.fasta, TRDV.fasta, TRDD.fasta, TRDJ.fasta</pre> <p>Note that:</p> <ul style="list-style-type: none"> • The names of the downloaded files follow the <IMGT group>.fasta convention, where <IMGT group> is one of the 7 <i>IMGT groups</i> with the "IG" prefix or one of the 10 <i>IMGT groups</i> with the "TR" prefix. • IMGT doesn't always provide germline sequences for all <i>IMGT groups</i>. This means that, depending on the IMGT/V-QUEST release and organism, some FASTA files can be missing. |
| destdir | A single string that is the path to the directory where the FASTA files containing the germline sequences should be downloaded. |
| overwrite | download_IMGT_germline_sequences() won't replace existing files, unless overwrite is set to TRUE. |
| recache | download_IMGT_germline_sequences() and list_IMGT_releases() both use a caching mechanism for the data that they download from IMGT/V-QUEST. In the case of download_IMGT_germline_sequences(), the downloaded files are cached on disk so the caching is persistent across R sessions. In the case of list_IMGT_releases(), the caching is in memory so does not persist across sessions. Set recache to TRUE to force a new download (and recaching) of the data. |
| ... | Extra arguments to be passed to the internal call to download.file(). See ?download.file in the utils package for more information. |

Value

download_IMGT_germline_sequences() returns the names of the downloaded FASTA files in an invisible character vector.

list_IMGT_releases() returns the list of IMGT/V-QUEST releases in a character vector. The releases are sorted from newest to oldest (latest release is first).

list_IMGT_organisms() returns the list of organisms included in the specified IMGT/V-QUEST release in a character vector.

IMGT_is_up() returns TRUE or FALSE, indicating whether the IMGT website at <https://www.imgt.org> is up and running or down.

See Also

- [install_IMGT_germline_db](#) to install a germline db from IMGT.
- [install_custom_germline_db](#) to install a germline db from user-supplied gene allele sequences.
- [download_OGRDB_germline_sequences](#) to download germline sequences from OGRDB.
- The IMGT website: <https://www.imgt.org/>.
- The IMGT/V-QUEST download site: <https://www.imgt.org/download/V-QUEST/>.

Examples

```

if (IMGT_is_up()) {
  ## As of April 10, 2026, the latest IMGT/V-QUEST release is 202614-2:
  list_IMGT_releases()

  list_IMGT_organisms("202614-2")

  ## Download Mouse BCR germline gene allele sequences from IMGT/V-QUEST
  ## 202614-2 to a temporary directory:
  dir.create(destdir <- tempfile("fasta_dir_"))
  download_IMGT_germline_sequences("202614-2", organism="Mus musculus",
                                   destdir=destdir)

  ## List the downloaded files:
  list.files(destdir) # 7 FASTA files

  ## Remove temporary directory:
  unlink(destdir, recursive=TRUE)

  ## Download Mouse TCR germline gene allele sequences from IMGT/V-QUEST
  ## 202614-2 to a temporary directory:
  dir.create(destdir <- tempfile("fasta_dir_"))
  download_IMGT_germline_sequences("202614-2", organism="Mus musculus",
                                   tcr.db=TRUE, destdir=destdir)

  ## List the downloaded files:
  list.files(destdir) # 10 FASTA files

  ## Remove temporary directory:
  unlink(destdir, recursive=TRUE)
}

```

download_OGRDB_germline_json

Download AIRR-C JSON files from OGRDB

Description

The `download_OGRDB_germline_json()` function downloads the AIRR-C JSON files associated with the specified germline sets from OGRDB.

The `extract_intdata_from_ogrdb_json()` function extracts the V gene delineations for the V alleles annotated in an AIRR-C JSON file. Note that it's a simple reimplementaion in R of the Python script `makeogrannotate.py` included in IgbLAST.

The `extract_auxdata_from_ogrdb_json()` function extracts the coding frame and CDR3 end information for the J alleles annotated in an AIRR-C JSON file.

If you use OGRDB in your research, please cite:

The current landscape of adaptive immune receptor genomic and repertoire data: OGRDB and VDJbase
 Lees, Peres, Klein, Amos et al., *Nucleic Acids Research*, November 2025.
<https://doi.org/10.1093/nar/gkaf1094>

Usage

```
download_OGRDB_germline_json(organism, germline_sets,
                             source_set=FALSE,
                             destdir=".", overwrite=FALSE,
                             recache=FALSE, ...)
```

```
extract_intdata_from_ogrdb_json(json_path, extra_fields=NULL)
extract_auxdata_from_ogrdb_json(json_path, extra_fields=NULL)
```

Arguments

| | |
|---------------|---|
| organism | A single string specifying the latin name of the organism for which to download the AIRR-C JSON file, e.g. "Homo sapiens". See https://ogrdb.airr-community.org/germline_sets for the organisms that are currently available at OGRDB. |
| germline_sets | A named integer vector. The names on the vector must be the names of the germline sets to download and the values in the vector must be their versions. Go to https://ogrdb.airr-community.org/germline_sets , select a Species, and click on "Show Germline Sets" to see the list of germline sets available for that species and their versions. |
| source_set | Whether to download the AIRR-C JSON files associated with the <i>Source Sets</i> instead of the <i>Reference Sets</i> . Only applies when organism is Homo sapiens. |
| destdir | A single string that is the path to the directory where the AIRR-C JSON files should be downloaded. |
| overwrite | <code>download_OGRDB_germline_json()</code> won't replace existing files, unless <code>overwrite</code> is set to TRUE. |
| recache | <code>download_OGRDB_germline_json()</code> uses a caching mechanism for the data that gets downloaded from OGRDB. The data is cached on disk so the caching is persistent across R sessions. Set <code>recache</code> to TRUE to force a new download (and recaching) of the data. |
| ... | Extra arguments to be passed to the internal call to <code>download.file()</code> . See ?download.file in the utils package for more information. |
| json_path | The path to an AIRR-C JSON file as one obtained with <code>download_OGRDB_germline_json()</code> . |
| extra_fields | NULL (the default) or a character vector of valid AIRR-C "AlleleDescription fields" to extract from the AIRR-C JSON file and to include in the returned data.frame as additional columns. See https://docs.airr-community.org/en/latest/datarep/germline.html#alleledescription-fields for the list of valid "AlleleDescription fields". |

Value

download_OGRDB_germline_json() returns the list of AIRR-C JSON files that it produced in an invisible character vector that carries the names of the corresponding germline sets.

extract_intdata_from_ogrdb_json() returns the V gene delineations in a data.frame with 1 row per V allele in the AIRR-C JSON file and with the same columns as the data.frame returned by load_intdata(). See ?load_intdata for the details. If extra_fields was specified, the returned data.frame will have 1 additional column per field in extra_fields.

extract_auxdata_from_ogrdb_json() returns the coding frame and CDR3 end information in a data.frame with 1 row per J allele in the AIRR-C JSON file and with the same columns as the data.frame returned by load_auxdata(). See ?load_auxdata for the details. If extra_fields was specified, the returned data.frame will have 1 additional column per field in extra_fields.

See Also

- [download_OGRDB_germline_sequences](#) to download germline sequences from OGRDB.
- [write_ndm_data](#) to write the data.frame returned by extract_intdata_from_ogrdb_json() in *ndm* format.
- [write_auxdata](#) to write the data.frame returned by extract_auxdata_from_ogrdb_json() in *aux* format.
- [compute_V_gene_delineations](#) to annotate a set of germline V gene allele sequences.
- [compute_auxdata](#) to annotate a set of germline J gene allele sequences.
- [intdata_utils](#) to access IgBLAST *internal data*.
- [auxdata_utils](#) to access IgBLAST *auxiliary data*.
- The OGRDB website: <https://ogrdb.airr-community.org/>.

Examples

```
## -----
## download_OGRDB_germline_json()
## -----
## We're going to use download_OGRDB_germline_json() to download the
## AIRR-C JSON files for mouse strain PWD/PhJ from OGRDB.
## Visit https://ogrdb.airr-community.org/germline_sets/Mus%20musculus
## for the list of germline sets available for mouse.

## We need the five following germline sets in order to cover all the
## BCR loci/regions for mouse strain PWD/PhJ:
germline_sets <- c(`PWD/PhJ IGH`=2, `PWD/PhJ IGKV`=1, `PWD/PhJ IGLV`=1,
                  `IGKJ (all strains)`=1, `IGLJ (all strains)`=1)

## Let's download them to a temporary directory:
dir.create(temp_json_dir <- tempfile("json_dir_"))
download_OGRDB_germline_json("Mus musculus", germline_sets,
                             destdir=temp_json_dir)

## List the downloaded files:
list.files(temp_json_dir) # 5 AIRR-C JSON files (one per germline set)

## -----
## extract_intdata_from_ogrdb_json()
## -----
## Let's use extract_intdata_from_ogrdb_json() to extract the V gene
```

```

## delineations for the V alleles annotated in the AIRR-C JSON files.
## Note that germline sets "IGKJ (all strains)" and "IGLJ (all strains)"
## don't annotate any V alleles so we exclude files IGKJ.json and
## IGLJ.json:

json_path <- file.path(temp_json_dir, "IGH.json")
IGHV_intdata <- extract_intdata_from_ogrdb_json(json_path)
head(IGHV_intdata)

json_path <- file.path(temp_json_dir, "IGKV.json")
IGKV_intdata <- extract_intdata_from_ogrdb_json(json_path)
head(IGKV_intdata)

json_path <- file.path(temp_json_dir, "IGLV.json")
IGLV_intdata <- extract_intdata_from_ogrdb_json(json_path)
IGLV_intdata

## Combine the 3 data.frames:
intdata <- rbind(IGHV_intdata, IGKV_intdata, IGLV_intdata)
dim(intdata)

## 'intdata' should be the same as the "internal data" included in
## germline db _OGRDB.mouse.PWD_PhJ.IGH+IGK+IGL.202410:
intdata0 <- load_intdata("_OGRDB.mouse.PWD_PhJ.IGH+IGK+IGL.202410")

## Let's use same_alleles_annot(), a small undocumented helper function,
## to verify that 'intdata' and 'intdata0' contain the same data (possibly
## with their rows in different order):
stopifnot(same_alleles_annot(intdata, intdata0))

## Note that calling extract_intdata_from_ogrdb_json() on an AIRR-C JSON
## file with no V alleles returns a 0-row data.frame with a warning:
json_path <- file.path(temp_json_dir, "IGKJ.json")
IGKJ_intdata <- extract_intdata_from_ogrdb_json(json_path) # warning!
stopifnot(nrow(IGKJ_intdata) == 0)
json_path <- file.path(temp_json_dir, "IGLJ.json")
IGLJ_intdata <- extract_intdata_from_ogrdb_json(json_path) # warning!
stopifnot(nrow(IGLJ_intdata) == 0)

## -----
## extract_auxdata_from_ogrdb_json()
## -----
## Let's use extract_auxdata_from_ogrdb_json() to extract the coding
## frame and CDR3 end information for the J alleles annotated in the
## AIRR-C JSON files. Note that germline sets "PWD/PhJ IGKV" and
## "PWD/PhJ IGLV" don't annotate any J alleles so we exclude files
## IGKV.json and IGLV.json:

json_path <- file.path(temp_json_dir, "IGH.json")
IGHJ_auxdata <- extract_auxdata_from_ogrdb_json(json_path)
IGHJ_auxdata

json_path <- file.path(temp_json_dir, "IGKJ.json")
IGKJ_auxdata <- extract_auxdata_from_ogrdb_json(json_path)
IGKJ_auxdata

json_path <- file.path(temp_json_dir, "IGLJ.json")

```

```

IGLJ_auxdata <- extract_auxdata_from_ogrdb_json(json_path)
IGLJ_auxdata

## Combine the 3 data.frames:
auxdata <- rbind(IGHJ_auxdata, IGKJ_auxdata, IGLJ_auxdata)
dim(auxdata)

## 'auxdata' should be the same as the "auxiliary data" included
## in germline db _OGRDB.mouse.PWD_PhJ.IGH+IGK+IGL.202410:
auxdata0 <- load_auxdata("_OGRDB.mouse.PWD_PhJ.IGH+IGK+IGL.202410")
stopifnot(same_alleles_annot(auxdata, auxdata0))

## Note that calling extract_auxdata_from_ogrdb_json() on an AIRR-C JSON
## file with no J alleles returns a 0-row data.frame with a warning:
json_path <- file.path(temp_json_dir, "IGKV.json")
IGKV_auxdata <- extract_auxdata_from_ogrdb_json(json_path) # warning!
stopifnot(nrow(IGKV_auxdata) == 0)
json_path <- file.path(temp_json_dir, "IGLV.json")
IGLV_auxdata <- extract_auxdata_from_ogrdb_json(json_path) # warning!
stopifnot(nrow(IGLV_auxdata) == 0)

## Remove temporary directory:
unlink(temp_json_dir, recursive=TRUE)

```

download_OGRDB_germline_sequences

Download germline sequences from OGRDB

Description

The `download_OGRDB_germline_sequences()` function downloads germline V/D/J gene allele sequences (nucleotides) from OGRDB for the specified germline sets. The sequences are downloaded in several FASTA files, one per *IMGT* group.

If you use OGRDB in your research, please cite:

The current landscape of adaptive immune receptor genomic and repertoire data: OGRDB and VDJbase
 Lees, Peres, Klein, Amos et al., *Nucleic Acids Research*, November 2025.
<https://doi.org/10.1093/nar/gkaf1094>

Usage

```

download_OGRDB_germline_sequences(organism, germline_sets,
                                  gapped=TRUE, source_set=FALSE,
                                  destdir=".", overwrite=FALSE,
                                  recache=FALSE, ...)

```

Arguments

organism A single string specifying the latin name of the organism for which to download the germline sequences, e.g. "Homo sapiens". See https://ogrdb.airr-community.org/germline_sets for the organisms that are currently available at OGRDB.

`germline_sets` A named integer vector. The names on the vector must be the names of the germline sets to download and the values in the vector must be their versions. Go to https://ogrdb.airr-community.org/germline_sets, select a Species, and click on "Show Germline Sets" to see the list of germline sets available for that species and their versions.

Note that each OGRDB germline set corresponds to one or more *IMGT groups*, where an *IMGT group* is a combination of locus and region type. For example:

- Germline set IGH_VDJ corresponds to *IMGT groups* IGHV, IGHD, and IGHJ.
- Germline set IGLambda_VJ corresponds to *IMGT groups* IGLV and IGLJ.
- Germline set PWD/PhJ IGH corresponds to *IMGT groups* IGHV, IGHD, and IGHJ.
- Germline set PWD/PhJ IGKV corresponds to *IMGT group* IGKV.

For reference, the 10 *IMGT groups* for B-cell Receptors are:

```
IGHV, IGHD, IGHJ, IGHC,
IGKV, IGKJ, IGKC,
IGLV, IGLJ, IGLC
```

IMPORTANT: The germline sets in your `germline_sets` vector must correspond to distinct *IMGT groups*.

Here is an example of valid `germline_sets` vector for *Mus musculus*:

```
germline_sets <- c(`PWD/PhJ IGH`=2,
  `PWD/PhJ IGKV`=1, `IGKJ (all strains)`=1,
  `PWD/PhJ IGLV`=1, `IGLJ (all strains)`=1)
```

`gapped` By default `download_OGRDB_germline_sequences()` will download the *gapped* V allele sequences. Set `gapped` to `FALSE` to download the *ungapped* V allele sequences instead.

`source_set` Whether to download the germline sequences from the *Source Sets* instead of the *Reference Sets*. Only applies when organism is *Homo sapiens*. Note that the AIRR-community/OGRDB maintainers recommend to use the *Reference Sets* for AIRR-seq analysis. See for example https://ogrdb.airr-community.org/germline_set/75. See also <https://wordpress.vdjbase.org/index.php/ogrdb/explanation-of-germline-set-formats/> for a brief explanation of the different germline set formats available from OGRDB.

`destdir` A single string that is the path to the directory where the FASTA files containing the germline sequences should be downloaded.

`overwrite` `download_OGRDB_germline_sequences()` won't replace existing files, unless `overwrite` is set to `TRUE`.

`recache` `download_OGRDB_germline_sequences()` uses a caching mechanism for the data that gets downloaded from OGRDB. The data is cached on disk so the caching is persistent across R sessions.

Set `recache` to `TRUE` to force a new download (and recaching) of the data.

... Extra arguments to be passed to the internal call to `download.file()`. See [?download.file](#) in the `utils` package for more information.

Value

The function returns the list of FASTA files that it produced in an invisible character vector that carries the names of the corresponding germline sets.

See Also

- [install_custom_germline_db](#) to install a germline db from user-supplied gene allele sequences.
- [download_IMGT_germline_sequences](#) to download germline sequences from IMGT.
- [download_OGRDB_germline_json](#) to download AIRR-C JSON files from OGRDB.
- The OGRDB website: <https://ogrdb.airr-community.org/>.
- <https://wordpress.vdjbase.org/index.php/ogrdb/explanation-of-germline-set-formats/> for a brief explanation of the different germline set formats available from OGRDB.

Examples

```
## -----
## DOWNLOAD RHESUS MONKEY GERMLINE SEQUENCES FROM OGRDB
## -----
## Visit https://ogrdb.airr-community.org/germline_sets/Macaca%20mulatta
## for the list of germline sets available for rhesus monkey.

## Download germline gene allele sequences for all BCR loci/regions
## for rhesus monkey to a temporary directory:
germline_sets <- c(IGH_VDJ=2, IGK_VJ=2, IGL_VJ=2)

dir.create(destdir <- tempfile("fasta_dir_"))
download_OGRDB_germline_sequences("Macaca mulatta", germline_sets,
                                destdir=destdir)

## List the downloaded files:
list.files(destdir) # 7 FASTA files (one per IMGT group)

## Remove temporary directory:
unlink(destdir, recursive=TRUE)

## -----
## DOWNLOAD MOUSE GERMLINE SEQUENCES FROM OGRDB
## -----
## Visit https://ogrdb.airr-community.org/germline_sets/Mus%20musculus
## for the list of germline sets available for mouse.

## Download germline gene allele sequences for all BCR loci/regions
## for mouse strain PWD/PhJ to a temporary directory:
germline_sets <- c(`PWD/PhJ IGH`=2,
                  `PWD/PhJ IGKV`=1, `IGKJ (all strains)`=1,
                  `PWD/PhJ IGLV`=1, `IGLJ (all strains)`=1)

dir.create(destdir <- tempfile("fasta_dir_"))
download_OGRDB_germline_sequences("Mus musculus", germline_sets,
                                destdir=destdir)

## List the downloaded files:
list.files(destdir) # 7 FASTA files (one per IMGT group)

## Remove temporary directory:
unlink(destdir, recursive=TRUE)
```

| | |
|------------------|--|
| get_igblast_root | <i>Control IgBLAST installation to use</i> |
|------------------|--|

Description

Get (or set) the IgBLAST installation used (or to be used) by the **igblastR** package.

Usage

```
get_igblast_root()
set_igblast_root(version_or_path)
```

Arguments

version_or_path

A single string that is either a version number (e.g. "1.22.0") or the path to an IgBLAST installation.

Details

set_igblast_root can be used to set or change the path to the IgBLAST installation to use. This can be an *internal* or *external* installation.

In the former case, version_or_path should be the version of an existing *internal* installation. The setting will be persistent.

In the latter case, it should be the full path (absolute or relative) to the *root directory* of a valid *external* installation. Note that the setting won't be persistent i.e. it won't be remembered across R sessions. See ?IGBLAST_ROOT for how to set the *external* IgBLAST installation to use in **igblastR** in a persistent manner.

Value

get_igblast_root() returns a single string containing the path to the *root directory* of the IgBLAST installation used by **igblastR**.

set_igblast_root() returns a single string containing the path to the *root directory* of the newly selected IgBLAST installation. The string is returned invisibly.

See Also

- The [igblastn](#) function to run the *igblastn standalone executable* included in IgBLAST from R. This is the main function in the **igblastR** package.
- [install_igblast](#) to perform an *internal* IgBLAST installation.
- [igblast_info](#) to collect basic information about the IgBLAST installation used by the **igblastR** package.
- [IGBLAST_ROOT](#) to set the *external* IgBLAST installation to be used by the **igblastR** package in a persistent manner.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

get_igblast_root()
```

 igblastn

BLAST for BCR/Ig and TCR sequences

Description

The `igblastn()` function is a wrapper to the `igblastn standalone executable` included in `IgBLAST`. This is the main function in the **igblastr** package.

Usage

```
igblastn(query, outfmt="AIRR",
          germline_db_V="auto", germline_db_V_seqidlist=NULL,
          germline_db_D="auto", germline_db_D_seqidlist=NULL,
          germline_db_J="auto", germline_db_J_seqidlist=NULL,
          organism="auto", c_region_db="auto",
          custom_internal_data="auto", auxiliary_data="auto",
          domain_system=c("imgt", "kabat"), ig_seqtype="auto",
          ...,
          out=NULL, parse.out=TRUE,
          show.in.browser=FALSE, show.command.only=FALSE)

igblastn_help(long.help=FALSE, show.in.browser=FALSE)
```

Arguments

| | |
|--------|--|
| query | <p>A character vector containing the paths to the input files (FASTA), or a <i>named DNAStringSet</i> object.</p> <p>If a character vector, then query must be of length ≥ 1 and each vector element must be the path to a FASTA file (possibly gz-compressed). In the context of <code>IgBLAST</code>, the DNA sequences in the FASTA files are referred to as <i>the query sequences</i>, and the sequence names found in the description lines of the FASTA records are referred to as <i>the query sequence ids</i>.</p> <p>Note that the query sequences are typically (but not always) stored in a single file, in which case query will be a single string. If more than one FASTA file is specified via query, then <code>igblastn()</code> will concatenate all the files together and pass the resulting file to the <code>igblastn standalone executable</code>.</p> <p>If query is a <i>DNAStringSet</i> object, then it must have names on it. These will be considered the query sequence ids.</p> |
| outfmt | <p>One of "AIRR", 3, 4, 7, or 19. "AIRR" is the default and is an alias for 19. outfmt can also be a string describing a customized format 7 e.g. "7 qseqid sseqid pident nident length score".</p> <p>See ?list_outfmt7_specifiers for more information about customizing format 7.</p> |

- `germline_db_V` "auto" (the default), or the path to a V-region db.
 Note that, by default (i.e. when `germline_db_V` is omitted or set to "auto"), `igblastn()` uses the V-region db that belongs to the cached germline db that is currently selected.
 See [?use_germline_db](#) for how to select the cached germline db to use with `igblastn()`.
- `germline_db_D` Same as `germline_db_V` but for the D-region db.
- `germline_db_J` Same as `germline_db_V` but for the J-region db.
- `germline_db_V_seqidlist,` `germline_db_D_seqidlist,`
`germline_db_J_seqidlist`
 Restrict search of germline database to list of gene alleles. A list of gene alleles can be specified either as a character vector of gene allele identifiers (e.g. IGHV3-23*01, IGHV3-23*04, etc...) or as the path to a file containing the identifiers (one identifier per line). In the latter case, a file object must be passed to the `germline_db_V_seqidlist`, `germline_db_D_seqidlist`, or `germline_db_J_seqidlist` argument. The file object will typically be constructed with something like `file("path/to/some/file")`.
- `organism` "auto" (the default), or the organism associated with the query sequences. Supported organisms include human, mouse, rat, rabbit and rhesus_monkey. Use [list_igblast_organisms\(\)](#) to obtain this list programmatically.
 Note that, by default (i.e. when `organism` is omitted or set to "auto"), `igblastn()` infers the organism from the name of the cached germline db that is currently selected.
 See [?use_germline_db](#) for how to select the cached germline db to use with `igblastn()`.
- `c_region_db` "auto" (the default), NULL, or the path to a C-region db.
 Note that, by default (i.e. when `c_region_db` is omitted or set to "auto"), `igblastn()` uses the cached C-region db that is currently selected.
 See [?use_c_region_db](#) for how to select the cached C-region db to use with `igblastn()`.
- `custom_internal_data`
 "auto" (the default), or a data.frame containing custom FWR/CDR annotation, or the path to a file containing custom FWR/CDR annotation, or NULL.
 By default (i.e. when `custom_internal_data` and `germline_db_V` are both omitted or set to "auto"), `igblastn()` will use:
- the *internal data* included in the cached germline db that is currently selected, if it has any;
 - the organism-specific IgBLAST *internal data* from NCBI if the cached germline db that is currently selected does not include its own *internal data*.
- Use [list_germline_dbs](#)(with `.intdata.only=TRUE`) to list the cached germline dbs that include their own *internal data*.
 See [?use_germline_db](#) for how to select the cached germline db to use with `igblastn()`.
 Note that when `custom_internal_data` is set to NULL, `igblastn()` always uses the organism-specific IgBLAST *internal data* from NCBI.
- `auxiliary_data` "auto" (the default), or a data.frame containing auxiliary data, or the path to a file containing auxiliary data (.aux file), or NULL.
 By default (i.e. when `auxiliary_data` and `germline_db_J` are both omitted or set to "auto"), `igblastn()` will use:

- the *auxiliary data* included in the cached germline db that is currently selected, if it has any;
- the organism-specific IgBLAST *auxiliary data* from NCBI if the cached germline db that is currently selected does not include its own *auxiliary data*.

Use `list_germline_dbs` (with `.auxdata.only=TRUE`) to list the cached germline dbs that include their own *auxiliary data*.

See `?use_germline_db` for how to select the cached germline db to use with `igblastn()`.

IMPORTANT NOTE: When `auxiliary_data` is set to `NULL`, then no *auxiliary data* is used. In this case, `igblastn()` can emit a significant number of the following warning:

```
Warning: Auxiliary data file could not be found
```

and various columns of the returned AIRR-formatted `tibble` (e.g. columns `vj_in_frame`, `productive`, `cdr3`, `fwr4`, and others) will be filled with NAs.

| | |
|----------------------------|---|
| <code>domain_system</code> | Set to "imgt" or "kabat". |
| <code>ig_seqtype</code> | Set to "Ig" or "TCR" depending on whether the query sequences are BCR/Ig or TCR sequences. Note that, by default (i.e. when <code>ig_seqtype</code> is omitted or set to "auto"), the value of <code>ig_seqtype</code> is inferred from the germline loci that appear in the name of the cached germline db that is currently selected (this name can be obtained with <code>use_germline_db()</code>). If these are BCR/Ig germline loci (i.e. IGH, IGK, IGL), then the inferred value will be "Ig". If they are TCR germline loci (TRA, TRB, TRG, TRD), then it will be "TCR". See <code>?use_germline_db</code> for how to select the cached germline db to use with <code>igblastn()</code> . |
| <code>...</code> | Extra arguments to be passed to the <code>igblastn standalone executable</code> . The list of valid arguments can be displayed with <code>igblastn_help()</code> . Note that the argument/value pairs must be passed to the <code>igblastn()</code> function in the usual R fashion. For example, what would be passed as <code>-num_alignments 1 -num_threads 8</code> when invoking the <code>igblastn standalone executable</code> in a terminal should be passed as <code>num_alignments_V=1, num_threads=8</code> when calling the <code>igblastn()</code> function: <pre>igblastn(query, num_alignments_V=1, num_threads=8)</pre> For options that don't require a value (e.g. <code>-extend_align5end</code> , <code>-extend_align3end</code> , <code>-ungapped</code> , etc...), pass the empty string (or a white string) to the argument. For example: <pre>igblastn(query, extend_align5end="", extend_align3end="")</pre> |
| <code>out</code> | <code>NULL</code> (the default), or the path to the file where the <code>igblastn standalone executable</code> should write its output. Note that, by default (i.e. when <code>out</code> is omitted or set to <code>NULL</code>), <code>igblastn()</code> instructs the <code>igblastn standalone executable</code> to write its output to a temporary file. |
| <code>parse.out</code> | Whether <code>igblastn()</code> should parse the plain-text output produced by the <code>igblastn standalone executable</code> or not, before returning it to the user. <code>TRUE</code> by default. If set to <code>FALSE</code> , then <code>igblastn()</code> returns the output as-is in a character vector, with one line per element in the vector. Note that <code>igblastn()</code> sets the |

"igblastn_raw_output" class attribute on this character vector, which allows compact display of the vector (this is achieved via a dedicated `print()` method defined in the **igblastr** package). The class attribute can be dropped with `unclass()`.

`show.in.browser`

For `igblastn()`: Whether the output of the `igblastn standalone executable` should also be displayed in a browser or not (in addition to being returned by the `igblastn()` function call). FALSE by default.

For `igblastn_help()`: Whether the help printed by the `igblastn standalone executable` (when invoked with the `-h` or `-help` argument) should be displayed in a browser or not. FALSE by default.

`show.command.only`

TRUE or FALSE. If set to TRUE, `igblastn()` won't invoke the `igblastn standalone executable` and instead will display the full command that shows how it would have invoked it. Note that the command is also returned in an invisible character vector. FALSE by default.

`long.help`

TRUE or FALSE. If set to FALSE (the default), the `igblastn standalone executable` is invoked with the `-h` argument. Otherwise, it's invoked with the `-help` argument.

Value

`igblastn()` captures the output produced by the `igblastn standalone executable` and returns it as:

- A [tibble](#) with 1 row per query sequence if `outfmt` is "AIRR" or 19 and `parse.out` is TRUE.
- A nested list with two top-level components (records and footer) if `outfmt` is 7 (or a customized format 7) and `parse.out` is TRUE. See [?read_igblastn_fmt7_output](#) for more information.
- A character vector with class attribute "igblastn_raw_output" on it in all other cases, that is, if `parse.out` is FALSE or `outfmt` is 3 or 4. See the `parse.out` argument above for more information.

Note

By default, the NCBI BLAST+ and IgBLAST programs will "call home" to report usage when they run on a computer with internet access. See <https://www.ncbi.nlm.nih.gov/books/NBK569851/> for the details. This can induce a significant slowdown in some situations e.g. when the `igblastn standalone executable` is called in a loop on a small set of query sequences at each iteration.

For this reason, the "call home" feature is disabled in **igblastr** by default, unless environment variable `BLAST_USAGE_REPORT` is set to true. See [?igblastr_usage_report](#) for more information.

See Also

- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.
- [install_igblast](#) to perform an *internal* IgBLAST installation.
- [igblast_info](#) to collect basic information about the IgBLAST installation used by the **igblastr** package.
- [install_IMGT_germline_db](#) to install a germline db from IMGT.
- [install_custom_germline_db](#) to install a germline db from user-supplied gene allele sequences.

- `use_germline_db` to select the cached germline db to use with `igblastn()`.
- `use_c_region_db` to select the cached C-region db to use with `igblastn()`.
- `igbrowser` to display the annotated sequences returned by `igblastn()` in a browser.
- `list_outfmt7_specifiers` for how to customize output format 7.
- `list_igblast_organisms` to list the organisms supported by IgBLAST.
- `igblastr_usage_report` to turn "Usage Reporting" on or off.
- `DNAStringSet` objects implemented in the **Biostrings** package.
- `tibble` objects implemented in the **tibble** package.

Examples

```

if (!has_igblast()) install_igblast()

igblast_info()

## -----
## Access query sequences and select germline and C-region dbs to use
## -----

## Files 'heavy_sequences.fasta' and 'light_sequences.fasta' included
## in igblastr contain 250 paired heavy- and light- chain sequences (125
## sequences in each file) downloaded from OAS (the Observed Antibody
## Space database):
filenames <- paste0(c("heavy", "light"), "_sequences.fasta")
query <- system.file(package="igblastr", "extdata", "BCR", filenames)

## Install Human germline db from IMGT:
db_name <- install_IMGT_germline_db("202614-2", "Homo_sapiens",
                                   overwrite=TRUE)

## Select germline db to use with igblastn():
use_germline_db(db_name)

## Select C-region db to use with igblastn():
use_c_region_db("_IMGT.human.IGH+IGK+IGL.202412")

## -----
## Call igblastn()
## -----

## We don't specify the 'outfmt' argument so output will be in AIRR
## format:
AIRR_df <- igblastn(query)
AIRR_df

## The result is a tibble with one row per query sequence:
class(AIRR_df)
dim(AIRR_df)

## You can call igbrowser() on 'AIRR_df' to visualize the annotated
## sequences in a browser. See '?igbrowser'.

## Note that this tibble can easily be converted to an ordinary data.frame
## with 'as.data.frame()', or to a DataFrame with 'as(., "DataFrame)":

```

```

as(AIRR_df, "DataFrame")

## To call igblastn() on a subset of the FASTA file, load the file as a
## DNASTringSet object with Biostrings::readDNASTringSet(), then subset
## the object, and finally pass the result of the subsetting operation
## to igblastn():
query_21_30 <- readDNASTringSet(query)[21:30]
query_21_30 # a DNASTringSet object with 10 sequences
igblastn(query_21_30)

## -----
## Parallel computing
## -----

## The igblastn standalone executable included in IgBLAST supports
## multithreading via command line argument -num_threads. To use it
## in igblastn(), set argument 'num_threads' to the desired value:
AIRR_df2 <- igblastn(query, num_threads=4)
stopifnot(identical(AIRR_df, AIRR_df2))

## Unfortunately, in our experience, using 'num_threads' on Linux
## doesn't achieve any significant speedup.

## Alternatively, one can parallelize execution at the R level using
## standard tools from the parallel or BiocParallel package. In this
## case it's the responsibility of the user to split the input in
## smaller batches and combine the results obtained for each batch:
if (.Platform$OS.type != "windows") {
  library(parallel)
  ## Split 'query' in 10 batches of 25 sequences each:
  batches <- split(readDNASTringSet(query), rep(1:10, each=25))
  AIRR_df3 <- mclapply(seq_along(batches),
                      function(i) igblastn(batches[[i]]),
                      mc.cores=4)
  ## Combine the results:
  AIRR_df3 <- do.call(rbind, AIRR_df3)
  stopifnot(identical(AIRR_df, AIRR_df3))
}

## Note that the actual speedup will depend on many factors like
## hardware, number and size of the batches, number of cores used,
## etc...

## -----
## TCR analysis
## -----

## NCBI IgBLAST can also be used for TCR sequence analysis, and so does
## igblastn().

## File 'SRR11341217.fasta.gz' included in igblastr contains 10,875 human
## beta chain TCR transcripts running from 5' of reverse transcription
## reaction to beginning of constant region:
filename <- "SRR11341217.fasta.gz"
query <- system.file(package="igblastr", "extdata", "TCR", filename)

## For this example, we're only keeping the first 100 sequences:

```

```

query <- head(readDNASTringSet(query), n=100)

## Install Human TCR germline db from IMGT:
db_name <- install_IMGT_germline_db("202614-2", "Homo_sapiens",
                                   tcr.db=TRUE, overwrite=TRUE)

## Select germline db to use with igblastn():
use_germline_db(db_name)

## Select C-region db to use with igblastn():
use_c_region_db("_IMGT.human.TRA+TRB+TRG+TRD.202509")

## Call igblastn(). Note that the 'ig_seqtype' argument will be
## automatically set to "TCR" (see documentation of the 'ig_seqtype'
## argument above in this man page for more information):
AIRR_df <- igblastn(query)

## -----
## More examples
## -----

## See '?read_igblastn_fmt7_output' for more examples.

```

igblast_usage_report *Turn "Usage Reporting" on or off*

Description

By default, the NCBI BLAST+ and IgBLAST programs will "call home" to report usage when they run on a computer connected to the internet. See <https://www.ncbi.nlm.nih.gov/books/NBK569851/> for the details. This can induce a significant slowdown in some situations e.g. when the *igblastn standalone executable* is called in a loop on a small set of query sequences at each iteration.

For this reason, the "call home" feature is disabled in **igblast** by default, unless environment variable `BLAST_USAGE_REPORT` is set to true.

More precisely, the "call home" feature is controlled by global option `igblast_usage_report` in **igblast**. On package startup, this option is set to TRUE if environment variable `BLAST_USAGE_REPORT` is set to true. Otherwise (i.e. if `BLAST_USAGE_REPORT` is not set, or is set to false or gibberish) it is set to FALSE.

Details

The user can change the value of global option `igblast_usage_report` any time with:

```
options(igblast_usage_report=TRUE)
```

or with:

```
options(igblast_usage_report=FALSE)
```

To get the value of this option, use:

```
getOption("igblast_usage_report")
```

Note that changing the value of a global option interactively with `options(...)` won't be remembered across R sessions. For a persistent change, you can either:

- Put the `options(...)` command in your `.Rprofile` file. See [?Rprofile](#) for more information. Note that this is the standard way of setting global options persistently.
- In the particular case of global option `igblastr_usage_report` an alternative is to define environment variable `BLAST_USAGE_REPORT` outside R. The exact way to do this is OS-dependent e.g. on Linux and Mac you can define it in your user's `.profile` by adding the following line to it:

```
export BLAST_USAGE_REPORT=true
```

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## Check current status of usage reporting:
getOption("igblastr_usage_report")

## Turn on usage reporting:
options(igblastr_usage_report=TRUE)

## Turn off usage reporting:
options(igblastr_usage_report=FALSE)
```

igblast_info

Check IgBLAST used by igblastr

Description

Collect basic information about the IgBLAST installation used by the **igblastr** package, or about any IgBLAST installation on the user machine.

Usage

```
igblast_info(igblast_root=get_igblast_root())

igblast_build(igblast_root=get_igblast_root())
igblastn_version(igblast_root=get_igblast_root(), raw.version=FALSE)
makeblastdb_version(igblast_root=get_igblast_root(), raw.version=FALSE)
list_igblast_organisms(igblast_root=get_igblast_root())

has_igblast()
```

Arguments

- `igblast_root` A single string that is the path to an IgBLAST installation. By default `igblast_root` is set to `get_igblast_root()`, which is the path to the IgBLAST installation used by the **igblast** package. See `?get_igblast_root` for more information.
- Note that the supplied string must contain the path to the *root directory* of an IgBLAST installation, that is, to a directory with a *bin* subdirectory in it that has the `igblastn`, `igblastp`, and `makeblastdb` *standalone executables* (on Windows these executables are files named `igblastn.exe`, `igblastp.exe`, and `makeblastdb.exe`, respectively).
- `raw.version` By default (i.e. when `raw.version` is omitted or set to `FALSE`), `igblastn_version()` and `makeblastdb_version()` return the version string of the `igblastn` and `makeblastdb` *standalone executables* included in IgBLAST. This string is extracted from the output produced by system commands:
- ```
igblastn -version
```
- and
- ```
makeblastdb -version
```
- When `raw.version` is set to `TRUE`, `igblastn_version()` and `makeblastdb_version()` return the *full output* produced by the above commands.

Value

`igblast_info()` returns a named list containing basic information about the IgBLAST installation.

`igblast_build()` returns a single string containing IgBLAST build information.

By default, `igblastn_version()` returns a single string containing the version of the `igblastn` *standalone executable* included in IgBLAST.

By default, `makeblastdb_version()` returns a single string containing the version of the `makeblastdb` *standalone executable* included in IgBLAST.

`list_igblast_organisms()` returns a character vector that lists the organisms for which IgBLAST provides *internal data*. Note that this is obtained by simply listing the content of the `internal_data` directory located in the IgBLAST installation.

`has_igblast()` returns `TRUE` or `FALSE`.

See Also

- The `igblastn` function to run the `igblastn` *standalone executable* included in IgBLAST from R. This is the main function in the **igblast** package.
- `install_igblast` to perform an *internal* IgBLAST installation.
- `get_igblast_root` to get (or set) the IgBLAST installation used (or to be used) by the **igblast** package.
- `IGBLAST_ROOT` to set the *external* IgBLAST installation to be used by the **igblast** package in a persistent manner.
- `intdata_utils` to access IgBLAST *internal data*.
- `auxdata_utils` to access IgBLAST *auxiliary data*.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

igblast_info()

list_igblast_organisms()
```

IGBLAST_ROOT

Use an external IgBLAST installation

Description

Select the *external* IgBLAST installation to use in **igblastr** in a persistent manner.

Details

The **igblastr** package can use 2 types of IgBLAST installation:

1. Internal (a.k.a. igblastr-managed): refers to an installation obtained with `install_igblast()`.
2. External: refers to an installation that is not managed by the **igblastr** package. This is usually an installation that was manually performed by you or a system administrator on your machine. It can be a system-wide installation or a per-user installation.

To use an *external* installation of IgBLAST in **igblastr**, set environment variable IGBLAST_ROOT to the path of the installation. Note that this must be the path to the *root directory* of the IgBLAST installation, that is, to a directory with a `bin` subdirectory in it that has the `igblastn`, `igblastp`, and `makeblastdb` *standalone executables* (on Windows these executables are files named `igblastn.exe`, `igblastp.exe`, and `makeblastdb.exe`, respectively).

This can be done within your current R session with `Sys.setenv(IGBLAST_ROOT="path/to/igblast_root")` for testing. However, this won't be remembered across R sessions.

To set IGBLAST_ROOT in a persistent manner, define it outside R. The exact way to do this is OS-dependent e.g. on Linux and Mac you can define it in your user's `.profile` by adding the following line to it:

```
export IGBLAST_ROOT="path/to/igblast_root"
```

See Also

- The `igblastn` function to run the `igblastn` *standalone executable* included in IgBLAST from R. This is the main function in the **igblastr** package.
- `install_igblast` to perform an *internal* IgBLAST installation.
- `igblast_info` to collect basic information about the IgBLAST installation used by the **igblastr** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

igbrowser

*Display annotated BCR sequences in a browser***Description**

Use `igbrowser()` to display the annotated BCR sequences returned by `igblastn()` in a browser. For each sequence, the V, D, and J segments are shown as well as the FWR1-4 and CDR1-3 regions. Additionally, the C segments are shown if the C-region information is available.

Usage

```
igbrowser(AIRR_df, show.full.sequence=FALSE, dna.coloring=TRUE,
          Vcolor="#FFDDD2", Dcolor="#CFC", Jcolor="#CEF", Ccolor="#EEC",
          FWRcolor="#C9D", CDRcolor="#EE4")
```

Arguments

`AIRR_df` The AIRR-formatted data.frame or `tibble` returned by `igblastn()`. Note that calling `igbrowser()` on a data.frame with thousands of rows is quite resource-intensive (it can even crash your browser!), so in this case we recommend sub-setting the data.frame before passing it to `igbrowser()` to keep the number of rows under 2000.

`show.full.sequence` By default, the part of the BCR sequences upstream of the V region is not shown. Set `show.full.sequence` to `TRUE` to show it.

`dna.coloring` Whether the nucleotides in the BCR sequences (sequence column in `AIRR_df`) should be colored or not.

`Vcolor, Dcolor, Jcolor, Ccolor` The background colors of the V, D, J, and C segments of the BCR sequences. Note that the C segments are shown only if `AIRR_df` contains C-region information.

`FWRcolor, CDRcolor` The background colors of the Framework Regions (FWR1-4) and Complementarity-Determining Regions (CDR1-3), respectively.

Value

0 or the error code returned by the internal call to `browseURL()`, invisibly.

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the `igblastR` package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.
- `tibble` objects implemented in the `tibble` package.

Examples

```

if (!has_igblast()) install_igblast()

query <- system.file(package="igblast", "extdata",
                    "BCR", "heavy_sequences.fasta")
use_germline_db("_OGRDB.human.IGH+IGK+IGL.202410")

## -----
## With C regions
## -----

use_c_region_db("_IMGT.human.IGH+IGK+IGL.202412")
AIRR_df <- igblastn(query)
igbrowser(AIRR_df)

## By default, the part of the sequences upstream of the V region is
## not shown. Use 'show.full.sequence=TRUE' to show the full sequences:
igbrowser(AIRR_df, show.full.sequence=TRUE)

## -----
## No C regions
## -----

use_c_region_db("")
AIRR_df2 <- igblastn(query)
igbrowser(AIRR_df2)

```

```
install_custom_germline_db
```

Install a germline db from user-supplied gene allele sequences

Description

The `install_custom_germline_db()` function creates a germline database from a set of user-supplied FASTA files containing germline V/D/J gene allele sequences (nucleotides), and installs it in **igblastr**'s persistent cache. The new germline database can then be used later with `igblastn()`.

Usage

```
install_custom_germline_db(db_name, fasta_dir, tcr.db=FALSE,
                           loci="auto", gapped=FALSE, with.intdata=FALSE,
                           with.auxdata=FALSE, imgt.fasta=FALSE,
                           disambiguate.allele.names=FALSE,
                           overwrite=FALSE, verbose=FALSE)
```

Arguments

| | |
|------------------------|--|
| <code>db_name</code> | A single string specifying the name of the germline db to install. The string <i>must</i> start with "cus" and cannot contain whitespace characters. |
| <code>fasta_dir</code> | The path to the directory containing the FASTA files to import in the database. Note that the files to import <i>must</i> be named: |

- IGH[VDJ].fasta, IGK[VJ].fasta, and IGL[VJ].fasta for a BCR germline loci database;
- TRA[VJ].fasta, TRB[VDJ].fasta, TRG[VJ].fasta, and TRD[VDJ].fasta for a TCR germline loci database.

See the `tcr.db` argument below for more details.

Note that the FASTA header lines are expected to contain the allele names:

- either on their own if the header lines have no |,
- or in the 2nd field if the header lines have fields separated with |.

`tcr.db`

Should the database to install be populated with allele sequences from the BCR (B-cell Receptor) or TCR (T-cell Receptor) germline loci?

The BCR germline loci are: IGH, IGK, IGL.

The TCR germline loci are: TRA, TRB, TRG, TRD.

By default, the V/D/J allele sequences from the BCR germline loci will be used. More precisely, the 7 following FASTA files will be expected to be present in `fasta_dir`:

```
IGHV.fasta, IGHD.fasta, IGHJ.fasta,
IGKV.fasta, IGKJ.fasta,
IGLV.fasta, IGLJ.fasta
```

The files that are effectively present will be imported in the database, with a warning if some of them are missing.

If `tcr.db` is set to TRUE, then the V/D/J allele sequences from the TCR germline loci will be used instead. More precisely, the 10 following FASTA files will be expected to be present in `fasta_dir`:

```
TRAV.fasta, TRAJ.fasta,
TRBV.fasta, TRBD.fasta, TRBJ.fasta,
TRGV.fasta, TRGJ.fasta,
TRDV.fasta, TRDD.fasta, TRDJ.fasta
```

The files that are effectively present will be imported in the database, with a warning if some of them are missing.

`loci`

By default, the database to install will be populated with the allele sequences from all the BCR or TCR loci. However, if you want to restrict the database to specific loci, you can use the `loci` argument to specify these loci. The subset of loci can be specified as a character vector with one element per locus (e.g. "IGH" or `c("TRA", "TRB")`), or as a +-separated list in a single string (e.g. "TRA+TRB").

`gapped`

Set to TRUE if the supplied V allele sequences are *gapped*, in which case the gaps will be removed when the sequences get imported in the database.

`with.intdata`

If the supplied V allele sequences are *gapped*, you can set `with.intdata` to TRUE to have the *internal data* (i.e. the FWR/CDR boundaries on the V alleles) computed and added to the new db.

Note that the presence of *internal data* in the new db will be indicated in the `intdata` column of the `data.frame` returned by `list_germline_dbs()`.

If a cached germline db includes its own *internal data*, then `igblastn()` will use it (as *custom internal data*) instead of the *internal data* shipped with IgBLAST. See documentation of the `custom_internal_data` argument in `?igblastn` for more information.

Also see `?intdata_utils` for more information about IgBLAST *internal data*.

`with.auxdata`

COMING SOON...

| | |
|---------------------------|--|
| imgt.fasta | COMING SOON... |
| disambiguate.allele.names | Note that <i>repeated</i> alleles (i.e. alleles with identical ungapped DNA sequences and names) are dropped. More precisely, only the first allele in each group of <i>repeated</i> alleles is kept. If, after dropping the <i>repeated</i> alleles, the names of the remaining alleles are not unique, then an error will be raised, unless <code>disambiguate.allele.names</code> is set to TRUE, in which case the <i>ambiguous</i> allele names will be disambiguated by adding a suffix to them. Finally, note that we only look at <i>repeated</i> alleles or <i>ambiguous</i> allele names within the same region type e.g. within the union of IGHV.fasta, IGKV.fasta, and IGLV.fasta (V regions), or within the union of TRBD.fasta and TRDD.fasta (D regions). But we never look at <i>repeated</i> alleles or <i>ambiguous</i> allele names across IGHV.fasta and IGHJ.fasta, or across TRBV.fasta and TRBD.fasta (these are very very unlikely to occur anyways). |
| overwrite | Set to TRUE if a germline db with the specified name is already installed in igblastr 's persistent cache, in which case the existing db will be replaced with the new one. |
| verbose | Set to TRUE to have the function display some details about its internal operations. |

Value

`install_custom_germline_db()` returns the name to the newly installed germline db as an invisible string.

Note

`install_custom_germline_db()` creates the local database by performing the instructions provided at <https://ncbi.github.io/igblast/cook/How-to-set-up.html>.

See Also

- [install_IMGT_germline_db](#) to install a germline db from IMGT.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- [list_germline_dbs](#) to list the cached germline dbs.
- [load_germline_db](#) to load the nucleotide sequences of the gene regions stored in a cached germline db.
- [use_germline_db](#) to select the cached germline db to use with `igblastn()`.
- [load_intdata](#) to access IgBLAST *internal data*.
- [download_IMGT_germline_sequences](#) to download germline sequences from IMGT.
- [download_OGRDB_germline_sequences](#) to download germline sequences from OGRDB.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

## -----
```

```

## A. WITH UNGAPPED GERMLINE V GENE ALLELE SEQUENCES
## -----

## In this admittedly artificial example, we're going to create and
## install a germline db from the gene allele sequences in built-in
## germline db _OGRDB.human.IGH+IGK+IGL.202410. More precisely, we're
## going to:
## 1. Dump built-in germline db _OGRDB.human.IGH+IGK+IGL.202410
##    in 7 FASTA files.
## 2. Use these FASTA files to create and install a new germline db.
## 3. Check that the new db contains the same sequences as the
##    original db.

## --- 1. Dump _OGRDB.human.IGH+IGK+IGL.202410 in 7 FASTA files ---

dump_germline_db <- function(db_name, destdir) {
  for (region_type in c("V", "D", "J")) {
    alleles <- load_germline_db(db_name, region_types=region_type)
    for (locus in c("IGH", "IGK", "IGL")) {
      loc_alleles <- alleles[grep(locus, names(alleles))]
      if (length(loc_alleles) == 0L)
        next
      fasta_file <- paste0(locus, region_type, ".fasta")
      writeXStringSet(loc_alleles, file.path(destdir, fasta_file))
    }
  }
}

db_name0 <- "_OGRDB.human.IGH+IGK+IGL.202410"
dir.create(exA_fasta_dir <- tempfile())

dump_germline_db(db_name0, exA_fasta_dir)

list.files(exA_fasta_dir) # 7 FASTA files

## --- 2. Use FASTA files to install new germline db ---

exA_db_name <- "cus.exA.human.IGH+IGK+IGL"
install_custom_germline_db(exA_db_name, exA_fasta_dir)
list_germline_dbs()

## --- 3. Check that new and original dbs have the same sequences ---

stopifnot(identical(
  as.character(load_germline_db(db_name0)),
  as.character(load_germline_db(exA_db_name))
))

## -----
## B. WITH GAPPED GERMLINE V GENE ALLELE SEQUENCES
## -----

## In this example, we're going to download human germline gene allele
## **gapped** sequences from AIRR-community/OGRDB and use them to create
## and install a new germline db. More precisely, we're going to:
## 1. Use download_OGRDB_germline_sequences() to download the
##    gapped sequences for the following "Reference Sets" from

```

```

##      AIRR-community/OGRDB:
##      - IGH_VDJ version 9
##      - IGKappa_VJ version 4
##      - IGLambda_VJ version 3
##      Each "Reference Set" corresponds to a BCR germline locus. For
##      each set, download_OGRDB_germline_sequences() will download the
##      FASTA file provided by AIRR-community/OGRDB that contains the
##      gapped allele sequences for that locus, and split the file into
##      one FASTA file per region type.
##      2. Use these FASTA files to create and install a new germline db.
##      3. Check that the new db contains the same sequences
##      as built-in germline db _OGRDB.human.IGH+IGK+IGL.202410.

## --- 1. Download human gapped sequences from AIRR-community/OGRDB ---

dir.create(exB_fasta_dir <- tempfile())
germline_sets <- c(IGH_VDJ=9, IGKappa_VJ=4, IGLambda_VJ=3)

download_OGRDB_germline_sequences("Homo sapiens", germline_sets,
                                  destdir=exB_fasta_dir)

list.files(exB_fasta_dir) # 7 FASTA files

## Note that the V allele sequences have gaps:
IGKV_alleles <- readDNASTringSet(file.path(exB_fasta_dir, "IGKV.fasta"))
as.character(IGKV_alleles[1:3])

## --- 2. Use FASTA files to install new germline db ---

exB_db_name <- "cus.exB.human.IGH+IGK+IGL"
install_custom_germline_db(exB_db_name, exB_fasta_dir, gapped=TRUE)
list_germline_dbs()

## --- 3. Check that the new db contains the same sequences ---
## ---          as _OGRDB.human.IGH+IGK+IGL.202410          ---

## The Reference Sets we used to populate the new db (IGH_VDJ version 9,
## IGKappa_VJ version 4, IGLambda_VJ version 3) are the same that were
## used to populate built-in germline db _OGRDB.human.IGH+IGK+IGL.202410.
## Let's check that the two dbs contain the same sequences:
stopifnot(identical(
  as.character(load_germline_db(db_name0)),
  as.character(load_germline_db(exB_db_name))
))

## -----
## C. CREATE AND INSTALL GERMLINE DB WITH COMPUTED "INTERNAL DATA"
## -----

## The gaps that germline gene allele sequence providers like IMGT and
## AIRR-community/OGRDB inject in the V allele sequences reflect their
## FWR/CDR boundaries. When the 'with.intdata' argument is set to TRUE,
## install_custom_germline_db() takes advantage of this to compute and
## add the "internal data" associated with the V alleles to the new db:
exC_db_name <- "cus.exC.human.IGH+IGK+IGL"
install_custom_germline_db(exC_db_name, exB_fasta_dir, gapped=TRUE,
                          with.intdata=TRUE)

```

```
list_germline_dbs() # 'intdata' col indicates presence of "internal data"

## Note that the "internal data" in a germline db can be loaded with
## load_intdata() (see '?load_intdata' for more information):
intdataC <- load_intdata(exC_db_name)
head(intdataC)

## Check that the "internal data" in the new db is the same as in
## _OGRDB.human.IGH+IGK+IGL.202410:
stopifnot(identical(intdataC, load_intdata(db_name0)))

## -----
## Remove the 3 custom dbs
## -----

rm_germline_db(exA_db_name)
rm_germline_db(exB_db_name)
rm_germline_db(exC_db_name)
```

install_igblast *Install IgBLAST*

Description

Download and install a pre-compiled IgBLAST from NCBI FTP site for use with **igblastr**.

Usage

```
install_igblast(release="LATEST", overwrite=FALSE, ...)
```

Arguments

| | |
|-----------|--|
| release | A single string specifying the IgBLAST release version to install. For example "LATEST" (recommended), or one of the IgBLAST release versions listed at https://ftp.ncbi.nih.gov/blast/executables/igblast/release/ (e.g. "1.21.0"). Note that old versions have not been tested and are not guaranteed to be compatible with the igblastr package. |
| overwrite | Set to TRUE to reinstall if the specified IgBLAST release version is already installed. |
| ... | Extra arguments to be passed to the internal call to <code>download.file()</code> . See ?download.file in the utils package for more information. |

Value

The path to the *root directory* of the IgBLAST installation, as an invisible string.

See Also

- The [igblastn](#) function to run the *igblastn standalone executable* included in IgBLAST from R. This is the main function in the **igblastr** package.
- [IGBLAST_ROOT](#) to use an *external* IgBLAST installation.
- [igblast_info](#) to collect basic information about the IgBLAST installation used by the **igblastr** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

igblast_info()
```

```
install_IMGT_germline_db
```

Install a germline db from IMGT

Description

The `install_IMGT_germline_db()` function downloads germline V/D/J gene allele sequences from the IMGT/V-QUEST download site for a given IMGT/V-QUEST release and organism, and stores them in a local IgBLAST-compatible germline database. This local database gets installed in **igblastr**'s persistent cache. It can then be used later with `igblastn()`.

CONDITIONS OF USE AND LICENSE: The IMGT data is provided to the academic users and NPO's (Not for Profit Organization(s)) under the CC BY-NC-ND 4.0 license. See <https://creativecommons.org/licenses/by-nc-nd/4.0/>. Any other use of IMGT material, from the private sector, needs a financial arrangement with CNRS.

Usage

```
install_IMGT_germline_db(release, organism="Homo sapiens", tcr.db=FALSE,
                        loci="auto",
                        without.intdata=FALSE, without.auxdata=FALSE,
                        overwrite=FALSE, verbose=FALSE, ...)
```

```
## Can only be used to "subset" a built-in C-region databse (see
## advanced example in Examples section below):
```

```
install_IMGT_c_region_db(organism, loci,
                        disambiguate.allele.names=FALSE,
                        overwrite=FALSE, verbose=FALSE)
```

Arguments

`release, organism, tcr.db`

See [?download_IMGT_germline_sequences](#) for the documentation of these arguments. Note that, in addition to the organism latin name (e.g. "Homo sapiens"), `install_IMGT_c_region_db()` also accepts the common name of the organism (e.g. "human").

`loci`

By default, the database to install will be populated with the allele sequences from all the BCR or TCR loci. However, if you want to restrict the database to specific loci, you can use the `loci` argument to specify these loci. The subset of loci can be specified as a character vector with one element per locus (e.g. "IGH" or c("TRA", "TRB")), or as a +-separated list in a single string (e.g. "TRA+TRB").

`without.intdata`

By default, `install_IMGT_germline_db()` will also compute and add the *internal data* (i.e. the FWR/CDR boundaries on the V alleles) to the new db. Note that `install_IMGT_germline_db()` uses the *gapped* V allele sequences provided by IMGT for that. You can skip this step by setting `without.intdata` to TRUE.

The presence of *internal data* in the new db will be indicated in the `intdata` column of the `data.frame` returned by `list_germline_dbs()`.

If a cached germline db includes its own *internal data*, then `igblastn()` will use it (as *custom internal data*) instead of the *internal data* shipped with IgBLAST. See documentation of the `custom_internal_data` argument in `?igblastn` for more information.

Also see `?intdata_utils` for more information about IgBLAST *internal data*.

`without_auxdata`

COMING SOON...

`overwrite` Set to TRUE to reinstall if the requested database is already installed.

`verbose` Set to TRUE to have the function display some details about its internal operations.

`...` Extra arguments to be passed to the internal call to `download.file()`. See `?download.file` in the `utils` package for more information.

`disambiguate_allele_names`

See documentation of the `disambiguate_allele_names` in `?install_custom_germline_db`.

Details

The following naming scheme is used to form the name of the installed database:

IMGT-<release>.<organism>.<loci>

where:

1. <release> is the IMGT/V-QUEST release e.g. 202614-2 or 202449-1. Use `list_IMGT_releases()` to get the list of releases currently available at IMGT/V-QUEST.
2. <organism> is the latin name (a.k.a. binomial name) of the organism with all spaces replaced with underscores (`_`). For example `Homo_sapiens` or `Macaca_mulatta`. Use `list_IMGT_organisms("<release>")` to get the list of organisms included in a given IMGT/V-QUEST release. Note that, starting with release 202405-2, IMGT/V-QUEST provides BCR and TCR germline gene allele sequences for mouse strain C57BL6J (`Mus_musculus_C57BL6J`).
3. <loci> is a string obtained by concatenating the germline loci together separated with the `+` sign. For example `IGH+IGK+IGL` or `TRA+TRB+TRG+TRD`. The list of loci depends on whether the germline gene allele sequences for BCR or TCR were requested. See `tcr.db` argument above for more information. Note that for some IMGT/V-QUEST releases/organisms, only a subset of the loci are available. For example, in release 202343-3, the only TCR germline loci available for `Mus_musculus_C57BL6J` are `TRA` and `TRB`. This will be automatically reflected in the name of the installed germline db.

Value

`install_IMGT_germline_db()` returns the name to the newly installed germline db as an invisible string.

`install_IMGT_c_region_db()` returns the name to the newly installed C-region db as an invisible string.

Note

`install_IMGT_germline_db()` creates the local database by performing the instructions provided at <https://ncbi.github.io/igblast/cook/How-to-set-up.html>.

See Also

- [install_custom_germline_db](#) to install a germline db from user-supplied gene allele sequences.
- [list_IMGT_releases](#) to list IMGT/V-QUEST releases.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the `igblastR` package.
- [list_germline_dbs](#) to list the cached germline dbs.
- [use_germline_db](#) to select the cached germline db to use with `igblastn()`.
- The IMGT website: <https://www.imgt.org/>.
- The IMGT/V-QUEST download site: <https://www.imgt.org/download/V-QUEST/>.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```

if (!has_igblast()) install_igblast()

if (IMGT_is_up()) {
  ## -----
  ## BASIC EXAMPLES
  ## -----

  ## As of April 10, 2026, the latest IMGT/V-QUEST release is 202614-2:
  list_IMGT_releases()

  list_IMGT_organisms("202614-2")

  ## Download human BCR germline gene allele sequences from IMGT/V-QUEST
  ## release 202614-2, and store them in a cached germline database:
  install_IMGT_germline_db("202614-2", organism="Homo sapiens",
                          overwrite=TRUE)

  ## List the cached germline databases:
  list_germline_dbs()

  ## Select newly installed germline db to use with igblastn():
  use_germline_db("IMGT-202614-2.Homo_sapiens.IGH+IGK+IGL")

  ## Download human TCR germline gene allele sequences from IMGT/V-QUEST
  ## release 202614-2, and store them in a cached germline database:
  install_IMGT_germline_db("202614-2", organism="Homo sapiens",
                          tcr.db=TRUE, overwrite=TRUE)

  list_germline_dbs()

  ## -----
  ## ADVANCED EXAMPLES
  ## -----

  ## Install an IMGT database for a subset of TCR loci:
  install_IMGT_germline_db("202614-2", organism="Homo sapiens",
                          loci="TRA+TRB", overwrite=TRUE)

  list_germline_dbs()

  ## Install the corresponding C-region database:

```

```

install_IMGT_c_region_db("human", "TRA+TRB", overwrite=TRUE)
list_c_region_dbs()

## Note that install_IMGT_c_region_db() can only be used to "subset"
## a built-in C-region database.
}

```

intdata-utils

Access IgBLAST internal data

Description

IgBLAST *internal data* is expected to annotate all the known germline V gene alleles for a given organism. It is provided by NCBI and is typically included in a standard IgBLAST installation.

The *internal data* informs about the FWR/CDR boundaries on the V allele sequences, also known as the V gene delineations. This information is needed and used internally by IgBLAST.

`get_intdata_path()` and `load_intdata()` can be used to access the *internal data* included in IgBLAST or in one of the *cached germline dbs* managed by **igblastr**.

`V_genes_with_varying_fwrcdr_boundaries()` can be used to identify germline V genes for which the FWR/CDR boundaries are not the same across all alleles.

Usage

```

## Access internal data:
get_intdata_path(organism, for.aa=FALSE, domain_system=c("imgt", "kabat"),
                 which=c("live", "original"))
load_intdata(organism, for.aa=FALSE, domain_system=c("imgt", "kabat"),
             which=c("live", "original"))

## Identify V genes with varying FWR/CDR boundaries across alleles:
V_genes_with_varying_fwrcdr_boundaries(intdata, V_segment=NULL)

```

Arguments

| | |
|----------------------------|---|
| <code>organism</code> | A single string containing the name of an organism as returned by <code>list_igblast_organisms()</code> . Alternatively, this can be the name of a cached germline db. Note that this works only for germline dbs that include their own <i>internal data</i> . You can use: <pre>list_germline_dbs(with.intdata.only=TRUE)</pre> to list them. See <code>?list_germline_dbs</code> for more information. |
| <code>for.aa</code> | By default, the data.frame returned by <code>load_intdata()</code> contains FWR/CDR boundaries reported with respect to the nucleotide sequences of the germline V alleles. Setting <code>for.aa</code> to TRUE will return a data.frame where they are reported with respect to the amino acid sequences of the germline V alleles. |
| <code>domain_system</code> | Domain system to be used for segment annotation. Must be "imgt" (the default) or "kabat". |

| | |
|-----------|---|
| which | By default, <code>get_intdata_path()</code> and <code>load_intdata()</code> access the "live IgBLAST data", that is, the IgBLAST data that the user has possibly updated with <code>update_live_igdata()</code> . Depending on whether updates were applied or not, the "live IgBLAST data" might differ from the original IgBLAST data. Set which to "original" if you want to access the original IgBLAST data instead. See <code>?update_live_igdata</code> for more information about "live" and "original" IgBLAST data. |
| intdata | A data.frame as returned by <code>load_intdata()</code> . |
| V_segment | The name of a V gene segment. This can be set to "fwr1", "cdr1", "fwr2", "cdr2", or "fwr3". By default (i.e. when <code>V_segment</code> is omitted or set to NULL), <code>V_genes_with_varying_fwrcdr_boundaries</code> will identify V genes for which any segment has varying boundaries across alleles. Otherwise, it will identify V genes for which the specified segment has varying boundaries. |

Details

IgBLAST *internal data* is typically included in a standard IgBLAST installation. It's located in the `internal_data/` directory which is itself a subdirectory of IgBLAST *root directory*.

Value

`get_intdata_path()` returns a single string containing the path to the *internal data* included in the IgBLAST installation used by **igblastr**, for the specified organism.

`load_intdata()` returns the *internal data* in a data.frame with 1 row per germline V gene allele sequence and the same columns as the data.frame returned by `read_ndm_data()`. See `?read_ndm_data` for more information.

`V_genes_with_varying_fwrcdr_boundaries()` returns a character vector containing the names of the germline V genes for which the FWR/CDR boundaries are not the same across all alleles.

See Also

- `compute_V_gene_delineations` to annotate a set of germline V gene allele sequences.
- `auxdata_utils` to access IgBLAST *auxiliary data*.
- `update_live_igdata` for more information about "live" and "original" IgBLAST data.
- `list_igblast_organisms` to list the organisms supported by IgBLAST.
- `list_germline_dbs` to list the cached germline dbs.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

igblast_info()

## -----
## Access IgBLAST internal data for a given organism
```

```

## -----
## IgBLAST only includes internal data for the following organisms:
list_igblast_organisms()

get_intdata_path("rabbit")

rabbit_intdata <- load_intdata("rabbit")
head(rabbit_intdata)

rabbit_intdata2 <- load_intdata("rabbit", for.aa=TRUE)
head(rabbit_intdata2)

## The values in the "end" cols in 'rabbit_intdata' are exactly 3 times
## those in the "end" cols in 'rabbit_intdata2':
end_colnames <- grep("_end$", colnames(rabbit_intdata), value=TRUE)
stopifnot(identical(rabbit_intdata[, end_colnames],
                    rabbit_intdata2[, end_colnames] * 3L))

## -----
## Access the internal data included in a germline db
## -----

## List germline dbs with internal data:
list_germline_dbs(with.intdata.only=TRUE)

## Access the internal data included in germline db
## _OGRDB.human.IGH+IGK+IGL.202410 (this data was inferred from
## the gaps in the germline V gene allele sequences provided by
## AIRR-community/OGRDB):
db_name <- "_OGRDB.human.IGH+IGK+IGL.202410"

get_intdata_path(db_name)

human_intdata <- load_intdata(db_name)
head(human_intdata)

## -----
## V_genes_with_varying_fwrcdr_boundaries()
## -----

## Note that the alleles of a given germline V gene don't necessarily
## share the same FWR/CDR boundaries. You can use utility function
## V_genes_with_varying_fwrcdr_boundaries() to identify them:
human_intdata0 <- load_intdata("human")
var_genes <- V_genes_with_varying_fwrcdr_boundaries(human_intdata0)
var_genes # 4 genes

## Display rows associated with gene IGHV4-31:
subset(human_intdata0, allele2gene(allele_name) == "IGHV4-31")

## Human germline V genes for which the CDR1 boundaries are not the
## same across all alleles:
var_genes <- V_genes_with_varying_fwrcdr_boundaries(human_intdata0,
                                                    V_segment="cdr1")

var_genes # 2 genes

```

```
## Display rows associated with the genes in 'var_genes':
subset(human_intdata0, allele2gene(allele_name) %in% var_genes)
```

J_alleles-inspect *Basic inspection of J allele sequences*

Description

A small set of utilities for (very) basic inspection of J allele sequences.

Note that all of them require access to the *auxiliary data* associated with the J alleles to inspect. See `auxdata` argument below for how to obtain this data.

Usage

```
translate_J_alleles(J_alleles, auxdata)
J_allele_has_stop_codon(J_alleles, auxdata)
translate_fwr4(J_alleles, auxdata, max.codons=NA)
```

Arguments

| | |
|-------------------------|---|
| <code>J_alleles</code> | A <code>DNAStrngSet</code> object containing germline J gene allele sequences. |
| <code>auxdata</code> | A <code>data.frame</code> as returned by <code>compute_auxdata()</code> , <code>compute_germline_db_auxdata()</code> , or <code>load_auxdata()</code> . |
| <code>max.codons</code> | The maximum number of FWR4 codons to translate. By default (i.e. when <code>max.codons</code> is NA) all the FWR4 codons are translated. |

Value

`translate_J_alleles()` returns a named character vector with 1 amino acid sequence per supplied allele. The vector contains an NA for any allele that is not annotated in `auxdata` or for which `auxdata$coding_frame_start` has an NA. The names on it are the names of the supplied alleles.

`J_allele_has_stop_codon()` returns a named logical vector with 1 value per supplied allele. The vector contains an NA for any allele that is not annotated in `auxdata` or for which `auxdata$coding_frame_start` has an NA. The names on it are the names of the supplied alleles.

`translate_fwr4()` returns a named character vector with 1 amino acid sequence per supplied allele. The vector contains an NA for any allele that is not annotated in `auxdata` or for which `auxdata$cdr3_end` has an NA.

See Also

- [auxdata_utils](#) to access IgBLAST *auxiliary data*.
- [intdata_utils](#) to access IgBLAST *internal data*.
- [V_alleles_inspect](#) for basic inspection of V allele sequences.
- [update_live_igdata](#) for more information about "live" and "original" IgBLAST data.
- `DNAStrngSet` objects in the **Biostrings** package.
- The [translate_codons](#) function upon which `translate_J_alleles()` and `translate_fwr4()` are based.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```

if (!has_igblast()) install_igblast()

igblast_info()

## -----
## 1. A close look at the J allele sequences for rabbit
## -----

## IgBLAST only includes auxiliary data for the following organisms:
list_igblast_organisms()

get_auxdata_path("rabbit")
rabbit_auxdata <- load_auxdata("rabbit")

## It turns out that IgBLAST auxiliary data for rabbit matches exactly
## the set of rabbit germline J alleles available at IMGT:
db_name <- install_IMGT_germline_db("202614-2", "Oryctolagus cuniculus",
                                   without_auxdata=TRUE,
                                   overwrite=TRUE)
J_alleles <- load_germline_db(db_name, region_types="J")
J_alleles # DNASTringSet object
stopifnot(setequal(names(J_alleles), rabbit_auxdata$allele_name))

## Note that this might change with future IMGT releases.

## Let's put the allele sequences in 'J_alleles' in the same order as
## in 'rabbit_auxdata':
J_alleles <- J_alleles[rabbit_auxdata$allele_name]
stopifnot(identical(names(J_alleles), rabbit_auxdata$allele_name))

## The 'coding_frame_start' column in 'rabbit_auxdata' contains integer
## values that are >= 0 and <= 2. They indicate how many nucleotides
## precede the first codon on each allele sequence. In other words,
## this is the number of nucleotides that we need to trim on the 5'
## end of the germline J allele sequence before we start translating
## it. translate_J_alleles() uses this information to translate the
## DNA sequences in 'J_alleles':
J_aa <- translate_J_alleles(J_alleles, rabbit_auxdata)
J_aa

## No sequence in 'J_aa' should contain the letter "*" which is used
## by translate_J_alleles() to represent a stop codon. However, one
## J allele in IMGT-202614-2.Oryctolagus_cuniculus.IGH+IGK+IGL seems
## to disobey:
has_stop_codon <- grepl("*", J_aa, fixed=TRUE)
rabbit_auxdata[has_stop_codon, ] # coding_frame_start = 0 for IGKJ1-2*04
J_alleles[has_stop_codon]       # first codon (TGA) is a stop codon
J_aa[has_stop_codon]           # indeed!

## -----
## 2. About the "WGXG" and "FGXG" motifs
## -----

## The FWR4 region is expected to start with the following amino acid
## motifs (X represents any amino acid):

```

```

## - "WGXG" on the heavy chain
## - "FGXG" on the light chain

## Let's use translate_fwr4() to extract and translate the first 4
## codons of the FWR4 region:
fwr4_head <- translate_fwr4(J_alleles, rabbit_auxdata, max.codons=4)

## We expect to see the "WGXG" and "FGXG" motifs here, and most of the
## time we do:
has_motif <- grepl("[FW]G.G", fwr4_head)
table(has_motif)

## However, there are a few exceptions:
fwr4_head[!has_motif]

## -----
## 3. Compute auxiliary data for a set of J allele sequences
## -----

## See '?compute_auxdata' for the details of what compute_auxdata() does.
computed_auxdata <- compute_auxdata(J_alleles)
head(computed_auxdata)

## Alleles for which the CDR3 end could not found:
cdr3_end_not_found <- is.na(computed_auxdata$cdr3_end)
stopifnot(identical(cdr3_end_not_found, !has_motif))
J_alleles[cdr3_end_not_found]
fwr4_head[cdr3_end_not_found] # déjà vu

## 'computed_auxdata' is in agreement with 'rabbit_auxdata', except for
## the 8 alleles for which the CDR3 end could not be found:
keep_idx <- which(!cdr3_end_not_found)
stopifnot(identical(computed_auxdata[keep_idx, ],
                    rabbit_auxdata[keep_idx, ]))

```

```
list_c_region_dbs      List cached C-region dbs
```

Description

List the *cached C-region dbs*, that is, the C-region databases currently installed in **igblast**'s persistent cache.

Usage

```

## List the cached C-region dbs:
list_c_region_dbs(builtin.only=FALSE, names.only=FALSE, long.listing=FALSE)

## Remove a C-region db from igblast's persistent cache:
rm_c_region_db(db_name)

```

Arguments

| | |
|---------------------------|--|
| <code>builtin.only</code> | By default <code>list_c_region_dbs()</code> returns the list of all cached C-region dbs, including built-in C-region dbs. Set <code>builtin.only</code> to <code>TRUE</code> to return only the list of built-in C-region dbs. Note that built-in dbs are prefixed with an underscore (<code>_</code>). |
| <code>names.only</code> | By default <code>list_c_region_dbs()</code> returns the list of cached C-region dbs in a <code>data.frame</code> with one db per row. Set <code>names.only</code> to <code>TRUE</code> to return the db names only, in which case they are returned in a character vector. |
| <code>long.listing</code> | <code>TRUE</code> or <code>FALSE</code> . If set to <code>TRUE</code> , then <code>list_c_region_dbs()</code> returns a named list with one list element per C-region db. Each list element is a named integer vector that indicates the number of C-region sequences per locus. Ignored if <code>names.only</code> is set to <code>TRUE</code> . |
| <code>db_name</code> | A single string specifying the name of the C-region db to remove from the cache. This cannot be a built-in db. |

Details

Cached germline dbs and C-region dbs: The **igblastr** package provides utility functions to manage the *cached germline dbs* and *cached C-region dbs* to use with `igblastn()`.

Terminology used across **igblastr** documentation:

- A *cached germline db* contains the nucleotide sequences of the germline V, D, and J gene alleles for a given organism.
- A *cached C-region db* contains the nucleotide sequences of the constant gene regions for a given organism.

This man page documents utilities that manage cached C-region dbs.

See `?list_germline_dbs` for utilities that manage cached germline dbs.

Built-in dbs: The **igblastr** package comes with preinstalled cached germline and C-region dbs that are called *built-in dbs*.

Built-in dbs have their name prefixed with an underscore (`_`).

Note that the built-in C-region dbs from IMGT were downloaded from <https://www.imgt.org/vquest/refseqh.html#constant-sets> and included in the **igblastr** package on the date indicated by the suffix of the db name.

Value

`list_c_region_dbs()` returns the list of all cached C-region dbs in a `data.frame` with one db per row (if `names.only` is `FALSE`, which is the default), or in a character vector (if `names.only` is `TRUE`). Column `C` in the `data.frame` indicates the number of C-region sequences in each db.

`rm_c_region_db()` does not return anything (invisible `NULL`).

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- `use_c_region_db` to select the cached C-region db to use with `igblastn()`.
- `reset_c_region_dbs` to reset the cached C-region dbs.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```

if (!has_igblast()) install_igblast()

## 7 built-in C-region dbs (prefixed with an underscore):
list_c_region_dbs()
list_c_region_dbs(names.only=TRUE) # db names only
list_c_region_dbs(long.listing=TRUE) # long listing

```

| | |
|-------------------|---------------------------------|
| list_germline_dbs | <i>List cached germline dbs</i> |
|-------------------|---------------------------------|

Description

List the *cached germline dbs*, that is, the germline databases currently installed in **igblastr**'s persistent cache.

Usage

```

## List the cached germline dbs:
list_germline_dbs(builtin.only=FALSE,
                  with.intdata.only=FALSE,
                  with.auxdata.only=FALSE,
                  names.only=FALSE, long.listing=FALSE)

## Remove a germline db from igblastr's persistent cache:
rm_germline_db(db_name)

```

Arguments

builtin.only By default `list_germline_dbs()` returns the list of all cached germline dbs, including built-in germline dbs. Set `builtin.only` to `TRUE` to return only the list of built-in germline dbs. Note that built-in dbs are prefixed with an underscore (`_`).

with.intdata.only

Set `with.intdata.only` to `TRUE` to return only the list of germline dbs that include their own *internal data* (i.e. the FWR/CDR boundaries on the V alleles). Note that whether a germline db includes its own *internal data* or not is reported in the `intdata` column.

If a cached germline db includes its own *internal data*, then `igblastn()` will use it (as *custom internal data*) instead of the *internal data* shipped with IgBLAST. See documentation of the `custom_internal_data` argument in `?igblastn` for more information.

with.auxdata.only

Set `with.auxdata.only` to `TRUE` to return only the list of germline dbs that include their own *auxiliary data* (i.e. annotations of the J alleles). Note that whether a germline db includes its own *auxiliary data* or not is reported in the `auxdata` column.

If a cached germline db includes its own *auxiliary data*, then `igblastn()` will use it instead of the *auxiliary data* shipped with IgBLAST. See documentation of the `auxiliary_data` argument in `?igblastn` for more information.

| | |
|--------------|--|
| names.only | By default list_germline_dbs() returns the list of cached germline dbs in a data.frame with one db per row. Set names.only to TRUE to return the db names only, in which case they are returned in a character vector. |
| long.listing | TRUE or FALSE. If set to TRUE, then list_germline_dbs() returns a named list with one list element per germline db. Each list element is an integer matrix that indicates the number of germline gene allele sequences per locus and region type. Ignored if names.only is set to TRUE. |
| db_name | A single string specifying the name of the germline db to remove from the cache. This cannot be a built-in db. |

Details

Cached germline dbs and C-region dbs: The **igblast** package provides utility functions to manage the *cached germline dbs* and *cached C-region dbs* to use with `igblastn()`.

Terminology used across **igblast** documentation:

- A *cached germline db* contains the nucleotide sequences of the germline V, D, and J gene alleles for a given organism.
- A *cached C-region db* contains the nucleotide sequences of the constant gene regions for a given organism.

This man page documents utilities that manage cached germline dbs.

See `?list_c_region_dbs` for utilities that manage cached C-region dbs.

Built-in dbs: The **igblast** package comes with preinstalled cached germline and C-region dbs that are called *built-in dbs*.

Built-in dbs have their name prefixed with an underscore (_).

Note that additional dbs can easily be installed with functions like `install_IMGT_germline_db` or `install_custom_germline_db`.

AIRR-community/OGRDB built-in cached dbs: Note that the built-in germline dbs starting with `_AIRR` are made of the AIRR-community/OGRDB germline sets available at https://ogrdb.airr-community.org/germline_sets/Homo%20sapiens for human and at https://ogrdb.airr-community.org/germline_sets/Mus%20musculus for mouse.

Each AIRR db is populated with the latest germline datasets that were available at AIRR-community/OGRDB at the time indicated by the date (in YYYYMM format) embedded in the db name.

The AIRR dbs with the `.src` suffix contain the *Source Sets*. The AIRR dbs without the `.src` suffix contain the *Reference Sets*. See https://github.com/HyrienLab/igblast/tree/devel/inst/extdata/germline_sets/AIRR/human/202410/README.md for more information.

The AIRR-community/OGRDB maintainers recommend to use the *Reference Sets* for AIRR-seq analysis. See for example https://ogrdb.airr-community.org/germline_set/75

Value

`list_germline_dbs()` returns the list of all cached germline dbs in a data.frame with one db per row (if `names.only` is FALSE, which is the default), or in a character vector (if `names.only` is TRUE). Columns V, D, J in the data.frame indicate the number of germline gene allele sequences for each region in each db.

`rm_germline_db()` does not return anything (invisible NULL).

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- `use_germline_db` to select the cached germline db to use with `igblastn()`.
- `install_IMGT_germline_db` to install a germline db from IMGT.
- `install_custom_germline_db` to install a germline db from user-supplied gene allele sequences.
- `intdata_utils` to access IgBLAST *internal data*.
- `reset_germline_dbs` to reset the cached germline dbs.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```

if (!has_igblast()) install_igblast()

## Get list of built-in germline dbs only.
list_germline_dbs(builtin.only=TRUE)
list_germline_dbs(builtin.only=TRUE, names.only=TRUE) # db names only

## Get list of germline dbs that include their own internal data:
list_germline_dbs(with.intdata.only=TRUE)
list_germline_dbs(with.intdata.only=TRUE, names.only=TRUE)

## Long listing:
list_germline_dbs(long.listing=TRUE)
list_germline_dbs(with.intdata.only=TRUE, long.listing=TRUE)

if (IMGT_is_up()) {
  ## Install Mouse germline db from IMGT:
  install_IMGT_germline_db("202614-2", "Homo sapiens", overwrite=TRUE)

  list_germline_dbs() # all germline dbs

  ## Select germline db to use with igblastn():
  db_name <- "IMGT-202614-2.Homo_sapiens.IGH+IGK+IGL"
  use_germline_db(db_name) # select germline db to use

  use_germline_db() # get current selection
}

```

ndm_data-IO

Read/write "ndm" files

Description

The *ndm* (Nucleotide Domain Mapping) file format is a specialized, tab-delimited text file used by IgBLAST to define the FWR/CDR boundaries on the germline V gene allele sequences.

The **igblastr** package provides low-level functions to read/write data in *ndm* format.

Usage

```
read_ndm_data(filepath)
write_ndm_data(ndm_data, file="", check.data=FALSE)

validate_ndm_rows(ndm_data, allow.repeated.rows=FALSE)
```

Arguments

| | |
|---------------------|--|
| filepath | The path to the "ndm" file to read. |
| ndm_data | A data.frame with 1 row per germline V gene allele sequence and the same columns as the data.frame returned by read_ndm_data(). See Value section below for the details. |
| file | The path to the "ndm" file to write. |
| check.data | COMING SOON... |
| allow.repeated.rows | COMING SOON... |

Value

read_ndm_data() returns a data.frame with 1 row per germline V gene allele sequence and the following columns:

- allele_name: allele name;
- fwr1_start, fwr1_end: FWR1 start/end positions (1-based);
- cdr1_start, cdr1_end: CDR1 start/end positions (1-based);
- fwr2_start, fwr2_end: FWR2 start/end positions (1-based);
- cdr2_start, cdr2_end: CDR2 start/end positions (1-based);
- fwr3_start, fwr3_end: FWR3 start/end positions (1-based);
- chain_type: chain type as a 2-letter code e.g. VK for "V allele from the Kappa locus" (BCR locus) or VG for "V allele from the Gamma locus" (TCR locus);
- coding_frame_start: first coding frame start position (0-based).

write_ndm_data() doesn't return anything (i.e. invisible NULL).

validate_ndm_rows() returns a logical vector with 1 value per row in the ndm_data data.frame indicating whether the data in that row is valid or not.

See Also

- [compute_V_gene_delineations](#) to annotate a set of germline V gene allele sequences.
- [extract_intdata_from_ogrdb_json](#) to extract the V gene delineations for the V alleles annotated in an AIRR-C JSON file.
- [intdata_utils](#) to access IgBLAST *internal data*.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## COMING SOON...
```

Description

Some utility functions to query the Observed Antibody Space database, a.k.a. OAS, and to download and manipulate data from OAS.

OAS's homepage: <https://opig.stats.ox.ac.uk/webapps/oas/>

Note that OAS has two databases: the "Unpaired Sequences" database and the "Paired Sequences" database. Some of the utilities documented in this man page only work for data retrieved from the latter.

Usage

```
## Read metadata/data from a single OAS unit file:
read_OAS_csv_metadata(file)
read_OAS_csv(file, skip=1, ...)
extract_sequences_from_paired_OAS_df(df, add.prefix=FALSE)

## Basic query of OAS website:
list_paired_OAS_studies(as.df=FALSE, recache=FALSE)
list_paired_OAS_units(study, as.df=FALSE, recache=FALSE)
download_paired_OAS_units(study, units=NULL, destdir=".", ...)

## Read metadata/data from a batch of downloaded OAS unit files:
extract_metadata_from_OAS_units(dir=".", pattern="\\.csv\\.gz$")
extract_sequences_from_paired_OAS_units(dir=".", pattern="\\.csv\\.gz$")

## Check OAS website:
OAS_is_up()
```

Arguments

| | |
|------------|---|
| file | A single string that is the path to an <i>OAS unit file</i> . |
| skip | The number of lines of the data file to skip before beginning to read data. The first line in an OAS unit file contains metadata in JSON format, so must always be skipped. |
| ... | For <code>read_OAS_csv()</code> : Extra arguments to be passed to the internal call to <code>read.table()</code> . See ?read.table in the utils package for more information. For <code>download_paired_OAS_units()</code> : Extra arguments to be passed to the internal call to <code>download.file()</code> . See ?download.file in the utils package for more information. |
| df | The <code>data.frame</code> or tibble returned by <code>read_OAS_csv()</code> . |
| add.prefix | TRUE or FALSE. Should the names on the returned DNAStrngSet object be the original sequence ids as-is (this is the default), or should the <code>heavy_chain_</code> and <code>light_chain_</code> prefixes be added to them? <code>extract_sequences_from_paired_OAS_df()</code> returns a DNAStrngSet object with the sequence ids as names. The sequence ids are obtained from the <code>sequence_id_heavy</code> and <code>sequence_id_light</code> columns of the supplied <code>data.frame</code> or tibble . By |

default, they are propagated as-is to the `DNAStrngSet` object, which makes it difficult to recognize which chain (heavy or light) the antibody sequences are coming from. Setting `add.prefix` to `TRUE` will add the `heavy_chain_` or `light_chain_` prefix to the names on the `DNAStrngSet` object, hence making it easy to identify which chain a given antibody sequence is coming from.

| | |
|----------------------|--|
| <code>as.df</code> | TRUE or FALSE. By default, i.e. when <code>as.df</code> is FALSE, <code>list_paired_OAS_studies()</code> and <code>list_paired_OAS_units()</code> return the list of studies or units in a character vector. Alternatively you can set <code>as.df</code> to TRUE to get the list in a 3-column data.frame that contains a directory index as displayed at https://opig.stats.ox.ac.uk/webapps/ngsdb/paired/ or at https://opig.stats.ox.ac.uk/webapps/ngsdb/paired/Jaffe_2022/csv/ . |
| <code>recache</code> | TRUE or FALSE. <code>list_paired_OAS_studies()</code> and <code>list_paired_OAS_units()</code> both cache the information retrieved from OAS website for the duration of the R session (note that this caching is done in memory so it does not persist across sessions). Set <code>recache</code> to TRUE to force a new retrieval (and recaching) of the results. |
| <code>study</code> | A single string containing the name of a study as returned by <code>list_paired_OAS_studies()</code> . |
| <code>units</code> | NULL, or a character vector that must be a subset of <code>list_paired_OAS_units(study)</code> in which case the download will be restricted to these units only. |
| <code>destdir</code> | A single string that is the path to the directory where the OAS unit files are to be downloaded. |
| <code>dir</code> | A single string that is the path to a directory containing OAS unit files. This will typically be the same as <code>destdir</code> above if the unit files were downloaded with <code>download_paired_OAS_units()</code> . |
| <code>pattern</code> | Regular expression passed to the internal call to <code>list.files()</code> to obtain the list of OAS unit files located in <code>dir</code> . No reason to change this unless you know what you are doing. |

Details

OAS delivers data in the form of *OAS unit files*. These files are typically obtained by running the `bulk_download.sh` script that OAS generates based on one's search criteria. They are compressed CSV (comma-separated values) files with the `.csv.gz` extension.

OAS unit files can vary a lot in size: from only a few KB to 25 MB or more.

The first line in an OAS unit file contains metadata in JSON format (which means that these files cannot strictly be considered CSV files).

The CSV data is MiAIRR-compliant (see The "MiAIRR format" paper in the References section below).

Value

`read_OAS_csv_metadata()` extracts the metadata from the specified OAS unit file and returns it in a named list.

`read_OAS_csv()` extracts the data from the specified OAS unit file and returns it in a [tibble](#). The tibble has 1 row per antibody sequence if the data is unpaired (i.e. comes from the "Unpaired Sequences" database), or 1 row per sequence pair if the data is paired (i.e. comes from the "Paired Sequences" database).

`extract_sequences_from_paired_OAS_df()` returns the sequence pairs in a named [DNAStrngSet](#) object where the names are the sequence ids. See `add.prefix` above for how the sequence ids are obtained.

`list_paired_OAS_studies()` returns the list of studies that populate the "Paired Sequences" database in a character vector. This list can be seen here: <https://opig.stats.ox.ac.uk/webapps/ngsdb/paired/>.

`list_paired_OAS_units()` returns the list of all the OAS unit files that belong to a given study from the "Paired Sequences" database.

`download_paired_OAS_units()` does not return anything (invisible NULL).

`extract_metadata_from_OAS_units()` returns the metadata of all the OAS unit files found in the specified directory in a data.frame with 1 row per file.

`extract_sequences_from_paired_OAS_units()` extracts the sequence pairs from all the OAS unit files found in the specified directory and returns them in a named `DNAStrngSet` object where the names are the sequence ids. The sequence ids are obtained by prefixing the original sequence ids found in the files with the name of the unit followed by `_heavy_chain_` or `_light_chain_`.

`OAS_is_up()` returns TRUE or FALSE, indicating whether the OAS website at <https://opig.stats.ox.ac.uk/webapps/oas/> is up and running or down.

References

- The OAS paper:
Tobias H. Olsen, Fergus Boyles, Charlotte M. Deane. Observed Antibody Space: A diverse database of cleaned, annotated, and translated unpaired and paired antibody sequences. Protein Science (2021). <https://doi.org/10.1002/pro.4205>
- The "MiAIRR format" paper:
Rubelt, F., Busse, C., Bukhari, S. et al. Adaptive Immune Receptor Repertoire Community recommendations for sharing immune-repertoire sequencing data. Nat Immunol 18, 1274–1278 (2017). <https://doi.org/10.1038/ni.3873>

See Also

- OAS's homepage at: <https://opig.stats.ox.ac.uk/webapps/oas/>
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the `igblastR` package.
- `tibble` objects implemented in the `tibble` package.
- `DNAStrngSet` objects implemented in the `Biostrings` package.

Examples

```
if (OAS_is_up()) {
  list_paired_OAS_studies()

  list_paired_OAS_units("Eccles_2020")

  ## Import all the pairs of antibody sequences from the Eccles_2020 study:

  download_dir <- tempdir()
  download_paired_OAS_units("Eccles_2020", destdir=download_dir)

  metadata <- extract_metadata_from_OAS_units(download_dir)
  metadata # data.frame with 1 row per unit file

  sequences <- extract_sequences_from_paired_OAS_units(download_dir)
  sequences # DNAStrngSet object
```

```

## Odd indices correspond to heavy chain sequences and even indices
## to light chain sequences:

head(names(sequences))

sequences[1:2] # 1st pair
sequences[3:4] # 2nd pair
sequences[5:6] # 3rd pair
# etc...
}

```

```
parse_imgt_fasta_headers
```

Parse IMGT FASTA headers

Description

The IMGT FASTA headers contain 15 fields separated by |. See <https://www.imgt.org/IMGIndex/Fasta.php>.

parse_imgt_fasta_headers() takes a vector of headers and parse them into a 15-column matrix with 1 row per header.

Usage

```
parse_imgt_fasta_headers(headers)
```

Arguments

headers A character vector of IMGT FASTA headers.

Value

A 15-column character matrix with 1 row per header. The column names on the matrix are:

1. IMGT_acc: IMGT/LIGM-DB accession number(s);
2. allele_name: IMGT gene and allele name;
3. organism: species/organism;
4. func: IMGT allele functionality;
5. region: exon(s), region name(s), or extracted label(s);
6. startend_in_IMGT_acc: start and end positions in the IMGT/LIGM-DB accession number(s);
7. nb_nuc: number of nucleotides in the IMGT/LIGM-DB accession number(s);
8. codon_start: codon start, or 'NR' (not relevant) for non coding labels;
9. extra_nuc_5prime: +n: number of nucleotides (nt) added in 5' compared to the corresponding label extracted from IMGT/LIGM-DB;
10. extra_nuc_3prime: +n or -n: number of nucleotides (nt) added or removed in 3' compared to the corresponding label extracted from IMGT/LIGM-DB;

11. nuc_corrected: +n, -n, and/or nS: number of added, deleted, and/or substituted nucleotides to correct sequencing errors, or not corrected if non corrected sequencing errors;
12. nb_aa: number of amino acids (AA): this field indicates that the sequence is in amino acids;
13. nb_chars: number of characters in the sequence: nt (or AA)+IMGT gaps=total;
14. partial: partial (if it is);
15. revcomp: reverse complementary (if it is).

See Also

- <https://www.imgt.org/IMGTindex/Fasta.php> for the format of IMGT FASTA headers.
- [download_IMGt_germline_sequences](#) to download germline sequences from IMGT.
- [compute_V_gene_delineations](#) to annotate a set of germline V gene allele sequences.

Examples

```

if (IMGT_is_up()) {
  ## -----
  ## Retrieve IMGT FASTA files for mouse and parse their headers
  ## -----

  ## As of April 10, 2026, the latest IMGT/V-QUEST release is 202614-2:
  list_IMGt_releases()

  list_IMGt_organisms("202614-2")

  ## Download Mouse BCR germline gene allele sequences from IMGT/V-QUEST
  ## 202614-2 to a temporary directory:
  dir.create(destdir <- tempfile("fasta_dir_"))
  download_IMGt_germline_sequences("202614-2", organism="Mus musculus",
                                   destdir=destdir)

  ## List the downloaded files:
  list.files(destdir) # 7 FASTA files

  ## Load the V allele sequences (gapped):
  IMGt_groups <- c("IGHV", "IGKV", "IGLV")
  V_fasta_files <- file.path(destdir, paste0(IMGt_groups, ".fasta"))
  gapped_V_alleles <- readDNAStrngSet(V_fasta_files)

  ## 'gapped_V_alleles' is a DNAStrngSet object that carries the IMGT
  ## FASTA headers as names:
  gapped_V_alleles          # 865 V alleles
  names(gapped_V_alleles)[1:4] # first 4 headers

  ## Let's parse the headers:
  parsed_V_headers <- parse_imgt_fasta_headers(names(gapped_V_alleles))
  dim(parsed_V_headers) # 865 x 15 matrix
  parsed_V_headers[1:4, ] # first 4 rows

  ## -----
  ## Scrutinize IMGT V alleles reported as "partial in 5'"
  ## -----

  ## The presence of "partial in 5'" in the 14th field ('partial') of

```

```

## a header indicates a sequence that is truncated at the 5' end:
partial <- parsed_V_headers[ , "partial"]
table(partial)
is_partial_in_5prime <- grepl("in 5", partial)
table(is_partial_in_5prime) # 35 sequences marked as "partial in 5'"

## The corresponding sequences are expected to start with a gap:
start_with_gap <- grepl("^\\.", as.character(gapped_V_alleles))

## However, it seems that 1 V allele (IGHV8-9*02) is marked
## as "partial in 5'" but does NOT start with a gap:
table(is_partial_in_5prime, start_with_gap)
not_ok <- is_partial_in_5prime & !start_with_gap
parsed_V_headers[not_ok, ] # partial in 5'
gapped_V_alleles[not_ok] # does NOT start with a gap!

not_ok_V_allele <- unname(parsed_V_headers[not_ok, "allele_name"])

stopifnot(identical(not_ok_V_allele, "IGHV8-9*02")) # sanity check

## -----
## Scrutinize IMGT V alleles reported as "partial in 3'"
## -----

## The presence of "partial in 3'" in the 14th field ('partial') of
## a header indicates a sequence that is truncated at the 3' end:
is_partial_in_3prime <- grepl("in 3", partial)
table(is_partial_in_3prime) # 72 sequences marked as "partial in 3'"

## Note that the FWR3 always ends at position 312 on a gapped V allele
## sequence. This means that any gapped sequence shorter than that should
## be considered truncated at the 3' end:
gapped_seq_lens <- lengths(gapped_V_alleles)
has_incomplete_fwr3 <- gapped_seq_lens < 312

## However, it seems that 12 V alleles have an incomplete FWR3 but
## are NOT marked as "partial in 3'":
table(has_incomplete_fwr3, is_partial_in_3prime)
not_ok <- has_incomplete_fwr3 & !is_partial_in_3prime
gapped_V_alleles[not_ok] # incomplete FWR3
parsed_V_headers[not_ok, ] # NOT marked as partial in 3'

parsed_V_headers[not_ok, "allele_name"] # suspect alleles

stopifnot(sum(not_ok) == 12L) # sanity check

## -----
## A quick look at the IMGT header of J alleles
## -----

## Load the J allele sequences:
IMGT_groups <- c("IGHJ", "IGKJ", "IGLJ")
J_fasta_files <- file.path(destdir, paste0(IMGT_groups, ".fasta"))
J_alleles <- readDNAStrngSet(J_fasta_files)

## 'J_alleles' is a DNAStrngSet object that carries the IMGT
## FASTA headers as names:

```

```
J_alleles          # 27 J alleles
names(J_alleles)[1:4] # first 4 headers

## Let's parse the headers:
parsed_J_headers <- parse_imgt_fasta_headers(names(J_alleles))
dim(parsed_J_headers) # 27 x 15 matrix
parsed_J_headers[1:4, ] # first 4 rows

## Codon start is the start position of the first codon:
as.integer(parsed_J_headers[ , "codon_start"])

## Remove temporary directory:
unlink(destdir, recursive=TRUE)
}
```

```
read_igblastn_AIRR_output
      igblastn output format 19 (AIRR format)
```

Description

Read `igblastn` output format 19 (AIRR format). This is the output produced by `igblastn` when `outfmt` is set to "AIRR" or 19.

This format is sometimes called "Rearrangement summary report" or simply "AIRR rearrangement tabular" format. See for example IgBLAST web interface at <https://www.ncbi.nlm.nih.gov/igblast/>.

Usage

```
read_igblastn_AIRR_output(out)
```

Arguments

| | |
|------------------|---|
| <code>out</code> | The path to a file containing the output produced by <code>igblastn</code> when <code>outfmt</code> is set to "AIRR" or 19. |
|------------------|---|

Value

A data.frame with 1 row per query sequence and many columns.

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastR** package.
- `read_igblastn_fmt7_output` to read and parse `igblastn` output format 7.
- IgBLAST web interface at <https://www.ncbi.nlm.nih.gov/igblast/>.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```

if (!has_igblast()) install_igblast()

## -----
## Access query sequences and select germline and C-region dbs to use
## -----

## Files 'heavy_sequences.fasta' and 'light_sequences.fasta' included
## in igblast contain 250 paired heavy- and light- chain sequences (125
## sequences in each file) downloaded from OAS (the Observed Antibody
## Space database):
filenames <- paste0(c("heavy", "light"), "_sequences.fasta")
query <- system.file(package="igblast", "extdata", "BCR", filenames)

## Keep only the first 10 sequences from each file:
query <- c(head(readDNAStringSet(query[[1L]]), n=10),
           head(readDNAStringSet(query[[2L]]), n=10))

## Select the germline and C-region dbs to use with igblastn():
use_germline_db("_OGRDB.human.IGH+IGK+IGL.202410")
use_c_region_db("_IMGT.human.IGH+IGK+IGL.202412")

## -----
## Call igblastn()
## -----

out <- tempfile()
igblastn(query, out=out)
AIRR_df <- read_igblastn_AIRR_output(out)

class(AIRR_df)
dim(AIRR_df) # 1 row per query sequence

tibble(AIRR_df)

```

```

read_igblastn_fmt7_output
      igblastn output format 7

```

Description

Read and parse igblastn output format 7. This is the output produced by igblastn when outfmt is set to 7.

This format is sometimes called "Tabular with comment lines" or simply "Tabular" format. See for example IgBLAST web interface at <https://www.ncbi.nlm.nih.gov/igblast/>.

Usage

```

read_igblastn_fmt7_output(out)

## Related utilities:
parse_outfmt7(out_lines)
list_outfmt7_specifiers()

```

Arguments

`out` The path to a file containing the output produced by `igblastn` when `outfmt` is set to 7.

`out_lines` The character vector returned by `igblastn(query, outfmt=7, parse.out=FALSE, ...)`.

Value

`read_igblastn_fmt7_output(out)` returns `parse_outfmt7(readLines(out))`.

`parse_outfmt7(out_lines)` returns the parsed form of `out_lines` in a list.

`list_outfmt7_specifiers()` returns the list of format specifiers supported by `igblastn` formatting option 7.

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the `igblastr` package.
- `read_igblastn_AIRR_output` to read `igblastn` output format 19 (AIRR format).
- IgBLAST web interface at <https://www.ncbi.nlm.nih.gov/igblast/>.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

## -----
## Access query sequences and select germline and C-region dbs to use
## -----

## Files 'heavy_sequences.fasta' and 'light_sequences.fasta' included
## in igblastr contain 250 paired heavy- and light- chain sequences (125
## sequences in each file) downloaded from OAS (the Observed Antibody
## Space database):
filenames <- paste0(c("heavy", "light"), "_sequences.fasta")
query <- system.file(package="igblastr", "extdata", "BCR", filenames)

## Keep only the first 10 sequences from each file:
query <- c(head(readDNASTringSet(query[[1L]]), n=10),
           head(readDNASTringSet(query[[2L]]), n=10))

## Select the germline and C-region dbs to use with igblastn():
use_germline_db("_OGRDB.human.IGH+IGK+IGL.202410")
use_c_region_db("_IMGT.human.IGH+IGK+IGL.202412")

## -----
## FIRST igblastn RUN: GET OUTPUT IN FORMAT 7
## -----

parsed_out7 <- igblastn(query, outfmt=7)

## Note that the above is equivalent to:
out <- tempfile()
igblastn(query, outfmt=7, out=out)
```

```

parsed_out7b <- read_igblastn_fmt7_output(out)
stopifnot(identical(parsed_out7b, parsed_out7))

## and to:
out_lines <- igblastn(query, outfmt=7, parse.out=FALSE)
out_lines # raw output
parsed_out7c <- parse_outfmt7(out_lines)
stopifnot(identical(parsed_out7c, parsed_out7))

## Now taking a closer look at the output...

## Output contains one record per query sequence:
length(parsed_out7$records) # 20

## Each record can have 5 or 6 sections:
## 1. query_details
## 2. VDJ_rearrangement_summary
## 3. VDJ_junction_details
## 4. subregion_sequence_details (can be missing)
## 5. alignment_summary
## 6. hit_table

## Taking a close look at the first record:
rec1 <- parsed_out7$records[[1]]
rec1

qseqid(rec1) # query sequence id associated with this record

rec1$hit_table # data.frame with the standard columns

## -----
## SECOND igblastn RUN: GET OUTPUT IN CUSTOMIZED FORMAT 7
## -----

## For this second run we request a customized format 7 by supplying
## space delimited format specifiers. Use list_outfmt7_specifiers() to
## get the list of format specifiers supported by igblastn formatting
## option 7:
list_outfmt7_specifiers()
outfmt <- "7 qseqid sseqid pident nident length score"
parsed_out7 <- igblastn(query, outfmt=outfmt)

## Taking a close look at the first record:
rec1 <- parsed_out7$records[[1]]
rec1$hit_table # data.frame with the requested columns (+ the
# automatic "chaintype" column)

```

```

reset_c_region_dbs    Reset the cached C-region dbs

```

Description

reset_c_region_dbs() will reset the cache of C-region dbs to a pristine state by:

1. Removing *all* the C-region dbs currently in **igblastr**'s persistent cache, including the built-in C-region dbs.

2. Recreating the predefined set of built-in C-region dbs.

This means that any non built-in C-region db present in the cache **WILL BE LOST!**

Also note that `reset_c_region_dbs()` will cancel the current selection if any C-region db was previously selected with `use_c_region_db()`.

Usage

```
reset_c_region_dbs(verbose=FALSE)
```

Arguments

| | |
|---------|---|
| verbose | Set to TRUE to have <code>reset_c_region_dbs()</code> display some details about its internal operations. |
|---------|---|

Value

Nothing (i.e. an invisible NULL).

See Also

- [list_c_region_dbs](#) to list the cached C-region dbs.
- [use_c_region_db](#) to select the cached C-region db to use with `igblastn()`.
- [install_IMGT_c_region_db](#) to install a C-region db from IMGT.

Examples

```
## DANGER ZONE!
## Not run:
reset_c_region_dbs(verbose=TRUE)

## End(Not run)
```

| | |
|--------------------|--------------------------------------|
| reset_germline_dbs | <i>Reset the cached germline dbs</i> |
|--------------------|--------------------------------------|

Description

`reset_germline_dbs()` will reset the cache of germline dbs to a pristine state by:

1. Removing *all* the germline dbs currently in **igblastr**'s persistent cache, including the built-in germline dbs.
2. Recreating the predefined set of built-in germline dbs.

This means that any non built-in germline db present in the cache **WILL BE LOST!**

Also note that `reset_germline_dbs()` will cancel the current selection if any germline db was previously selected with `use_germline_db()`.

Usage

```
reset_germline_dbs(verbose=FALSE)
```

Arguments

verbose Set to TRUE to have reset_germline_dbs() display some details about its internal operations.

Value

Nothing (i.e. an invisible NULL).

See Also

- [list_germline_dbs](#) to list the cached germline dbs.
- [use_germline_db](#) to select the cached germline db to use with igblastn().
- [install_IMGT_germline_db](#) to install a germline db from IMGT.
- [install_custom_germline_db](#) to install a germline db from user-supplied gene allele sequences.

Examples

```
## DANGER ZONE!
## Not run:
  reset_germline_dbs(verbose=TRUE)

## End(Not run)
```

| | |
|---------------------|---|
| summarizeMismatches | <i>Summarize mismatches and indels between query and germline sequences</i> |
|---------------------|---|

Description

TODO

| | |
|------------------|---|
| translate_codons | <i>Extract and translate codons from a set of DNA sequences</i> |
|------------------|---|

Description

The translate_codons() function extracts and translates codons from a set of DNA sequences.

Usage

```
translate_codons(dna, offset=0, with.init.codon=FALSE)
```

```
## Used internally by translate_codons():
extract_codons(dna, offset=0)
remove_gaps(dna, gap_letter=".")
```

Arguments

| | |
|-----------------|---|
| dna | A DNAStringSet (or DNAString) object containing the codons to translate. Note that if the sequences in dna contain gaps (represented by the . letter) then translate_codons() and extract_codons() will remove them (by calling remove_gaps() internally) before doing anything else. |
| offset | The number of nucleotides that precede the first codon to translate. This must be supplied as a numeric vector with one value per sequence in dna, or as a single value. If the latter, then the same offset is used for all sequences. |
| with.init.codon | Is the first codon to translate in each DNA sequence the initiation codon? By default, with.init.codon is set to FALSE, in which case translate_codons() assumes that the first codon to translate in each DNA sequence is <i>not</i> the initiation codon. See documentation of the no.init.codon argument in ?translate in the Biostrings package for more information. |
| gap_letter | The letter or symbol representing gaps. Note that IMGT and AIRR-community/OGRDB typically use dots (".") to represent gaps in germline V gene allele sequences. |

Value

translate_codons() returns:

- An [AAStringSet](#) object with one amino acid sequence per input sequence if a [DNAStringSet](#) object was supplied.
- An [AAString](#) object if a [DNAString](#) object was supplied.

extract_codons() returns:

- A [DNAStringSet](#) object with one sequence per input sequence if a [DNAStringSet](#) object was supplied.
- A [DNAString](#) object if a [DNAString](#) object was supplied.

The sequences returned by extract_codons() are obtained by (1) removing all gaps (i.e. all . letters) from the input sequences, and (2) trimming the gap-free sequences as follow:

- On their 5' end, sequences are trimmed by the amount of nucleotides specified in offset.
- On their 3' end, sequences are trimmed by the smallest amount of nucleotides that makes the length of the trimmed sequence a multiple of 3. Note that this will always be 0, 1, or 2 nucleotides.

remove_gaps() returns the gap-free versions of the input sequences in an object of the same type as the input object.

See Also

- [DNAStringSet](#) and [AAStringSet](#) objects in the **Biostrings** package.
- The [translate](#) function in the **Biostrings** package upon which translate_codons() is based.
- [list_germline_dbs](#) to list all the *cached germline dbs*, that is, all the germline databases currently installed in **igblast**'s persistent cache.

Examples

```

## -----
## translate_codons()
## -----

## Gaps will be ignored:
dna1 <- DNASTringSet(c(".TTGTCCTTTAT..A", "GAAT.CATTTATC", "CTGTCGTTTATT"))
translate_codons(dna1)

## Handling of initiation codons (see '?translate' in the Biostrings
## package for how initiation codons are handled):
translate_codons(dna1, with.init.codon=TRUE)

## Load germline V gene allele sequences for human:
list_germline_dbs()
db_name <- "_OGRDB.human.IGH+IGK+IGL.202410"
V_alleles <- load_germline_db(db_name, region_types="V")
V_alleles # DNASTringSet object

## Translate them:
V_aa <- translate_codons(V_alleles)
V_aa # AAStringSet object

## Some human germline V gene allele sequences have a stop codon:
has_stop_codon <- grepl(".*", as.character(V_aa), fixed=TRUE)
V_aa[has_stop_codon]

## -----
## extract_codons()
## -----
extract_codons(dna1)
extract_codons(dna1, offset=1)

dna2 <- DNASTringSet(c("CCCAAAGGGTTT",
                      "CCCAAAGGGTTT",
                      "CAAAGGGTTT",
                      "AAAGGGTTT"))

extract_codons(dna2)
extract_codons(dna2, offset=1)
extract_codons(dna2, offset=2)
extract_codons(dna2, offset=3)
extract_codons(dna2, offset=4)

extract_codons(dna2, offset=3:0)
extract_codons(dna2, offset=4:1)
extract_codons(dna2, offset=5:2)
extract_codons(dna2, offset=6:3)

## -----
## remove_gaps()
## -----
remove_gaps(dna1)

```

Description

A small set of low-level utility functions to update and manage IgBLAST *auxiliary data* and *internal data*.

Usage

```
update_live_igdata(check.only=FALSE)

igdata_info()

time_since_live_igdata_last_checked(units="days")

reset_live_igdata(subdirs=c("all", "internal_data", "optional_file"))
```

Arguments

| | |
|------------|---|
| check.only | By default, update_live_igdata() checks for new IgBLAST auxiliary or internal data files available at NCBI, and it installs them if any are found. Set check.only to TRUE to only do the check without installing anything. |
| units | See ?base::difftime for valid units. |
| subdirs | By default, reset_live_igdata() resets both internal_data/ and optional_file/ directories to their original states. Set subdirs to "internal_data" or "optional_file" to reset only a particular directory. |

Details

Auxiliary data and internal data: A standard IgBLAST installation – like the one used by the **igblast** package – typically includes *auxiliary data* and *internal data* that are normally found in directories internal_data/ and optional_file/, respectively. Both directories should be subdirectories of the *root directory* of the IgBLAST installation, that is, of the directory returned by `get_igblast_root()`.

We sometimes refer to this data simply as the *IgBLAST data*.

NCBI updates: NCBI occasionally updates some of the files in the internal_data/ and optional_file/ directories between IgBLAST releases, and it is recommended to use the new files. They make the new files available at <https://ftp.ncbi.nih.gov/blast/executables/igblast/release/patch/>.

To download and install these new files, simply call update_live_igdata(). This will check for new IgBLAST auxiliary or internal data files available at NCBI, and install them if any are found.

You can restore the original files at any moment with reset_live_igdata().

Value

update_live_igdata() and reset_live_igdata() don't return anything (invisible NULL).

igdata_info() returns a named list containing information about the state of the "live" and "original" IgBLAST data.

time_since_live_igdata_last_checked() returns the time passed since the last run of update_live_igdata() in the specified units (days by default).

See Also

- `intdata_utils` to access IgBLAST *internal data*.
- `auxdata_utils` to access IgBLAST *auxiliary data*.
- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- `install_igblast` to perform an *internal* IgBLAST installation.
- `get_igblast_root` to get (or set) the IgBLAST installation used (or to be used) by the **igblastr** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```

if (!has_igblast()) install_igblast()

igblast_info()

## -----
## Check for NCBI updates
## -----

igdata_info()
update_live_igdata(check.only=TRUE)
igdata_info()

## -----
## "live" vs "original" IgBLAST data
## -----

## By default, the "live IgBLAST data" gets accessed or returned:
get_auxdata_path("human")
live_human_auxdata <- load_auxdata("human")

## Access the original IgBLAST data:
get_auxdata_path("human", which="original")
orig_human_auxdata <- load_auxdata("human", which="original")

## "live" and "original" IgBLAST data can differ if the former was
## updated with update_live_igdata(). Otherwise, they'll be the same:
identical(live_human_auxdata, orig_human_auxdata)

```

```
use_c_region_db
```

```
Select cached C-region db to use with igblastn()
```

Description

`use_c_region_db()` allows the user to select the cached C-region db to use with `igblastn()`. This choice will be remembered for the duration of the current R session but can be changed anytime.

`load_c_region_db()` allows the user to load the nucleotide sequences of the germline gene alleles stored in a cached C-region db.

Usage

```
use_c_region_db(db_name=NULL, verbose=FALSE)

load_c_region_db(db_name)
```

Arguments

| | |
|---------|---|
| db_name | For use_c_region_db(): NULL or a single string specifying the name of the cached C-region db to use. Use <code>list_c_region_dbs()</code> to list all the cached C-region dbs. If set to NULL (the default), then use_c_region_db() returns the name of the cached C-region db that is currently in use, if any. Otherwise it returns the empty string (""). Note that the current selection can be cancelled with use_c_region_db(""). For load_c_region_db(): A single string specifying the name of the cached C-region db from which to load the allele sequences. Use <code>list_c_region_dbs()</code> to list all the cached C-region dbs. |
| verbose | If set to TRUE, then use_c_region_db() will display some information about its internal operations. |

Value

When called with no argument, use_c_region_db() returns a single string containing the name of the cached C-region db currently used by `igblastn()` if any, or the empty string ("") if `igblastn()` is not using any C-region db.

When called with the db_name argument, use_c_region_db(db_name) returns db_name invisibly.
load_c_region_db() returns the allele sequences from the specified C-region db in a named `DNAStrngSet` object.

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the `igblastR` package.
- `list_c_region_dbs` to list the cached C-region dbs.
- `use_germline_db` to select the cached germline db to use with `igblastn()`.
- `DNAStrngSet` objects in the `Biostrings` package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

list_c_region_dbs() # all cached C-region dbs

## Select C-region db to use with igblastn():
db_name <- "_IMGT.human.IGH+IGK+IGL.202412"
use_c_region_db(db_name)

## Get current selection:
use_c_region_db() # get current selection
```

```
## Cancel current selection:
use_c_region_db("")
use_c_region_db()

## Load C-region sequences:
load_c_region_db(db_name)
load_c_region_db("_IMGT.mouse.IGH.202509")
```

```
use_germline_dbs      Select cached germline db to use with igblastn()
```

Description

use_germline_db() allows the user to select the cached germline db to use with `igblastn()`. This choice will be remembered for the duration of the current R session but can be changed anytime.

load_germline_db() allows the user to load the nucleotide sequences of the germline gene alleles stored in a cached germline db.

Usage

```
use_germline_db(db_name=NULL, verbose=FALSE)
```

```
load_germline_db(db_name, region_types=NULL)
```

Arguments

| | |
|--------------|--|
| db_name | For use_germline_db(): NULL or a single string specifying the name of the cached germline db to use. Use <code>list_germline_dbs()</code> to list all the cached germline dbs. If set to NULL (the default), then use_germline_db() returns the name of the cached germline db that is currently in use, if any. Otherwise it raises an error. For load_germline_db(): A single string specifying the name of the cached germline db from which to load the V, D, and/or J allele sequences. Use <code>list_germline_dbs()</code> to list all the cached germline dbs. |
| verbose | If set to TRUE, then use_germline_db() will display some information about its internal operations. |
| region_types | The types of regions (V, D, and/or J) to load from the database. Specified as a single string (e.g. "DJ") or as a character vector of single-letter elements (e.g. c("D", "J")). By default (i.e. when region_types is NULL), all the regions are returned. |

Value

When called with no argument, use_germline_db() returns a single string containing the name of the cached germline db currently used by `igblastn()` if any, or it raises an error if no germline db has been selected yet.

When called with the db_name argument, use_germline_db(db_name) returns db_name invisibly.

load_germline_db() returns the allele sequences from the specified germline db in a named `DNAStrngSet` object.

See Also

- The `igblastn` function to run the `igblastn standalone executable` included in IgBLAST from R. This is the main function in the **igblastr** package.
- `list_germline_dbs` to list the cached germline dbs.
- `use_c_region_db` to select the cached C-region db to use with `igblastn()`.
- `DNAStrngSet` objects in the **Biostrings** package.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
if (!has_igblast()) install_igblast()

list_germline_dbs() # all cached germline dbs

## Select germline db to use with igblastn():
db_name <- "_OGRDB.human.IGH+IGK+IGL.202410"
use_germline_db(db_name)

## Get current selection:
use_germline_db()

## Load germline gene allele sequences:
load_germline_db(db_name)
load_germline_db(db_name, region_types="D")
load_germline_db(db_name, region_types="DJ")
```

V_alleles-inspect *Basic inspection of V allele sequences*

Description

A small set of utilities for (very) basic inspection of V allele sequences.

Note that all of them require access to the *internal data* associated with the V alleles to inspect. See `intdata` argument below for how to obtain this data.

Usage

```
translate_V_alleles(V_alleles, intdata, V_segment=NULL)
V_allele_has_stop_codon(V_alleles, intdata)
```

Arguments

| | |
|-----------|--|
| V_alleles | A <code>DNAStrngSet</code> object containing germline V gene allele sequences. |
| intdata | A <code>data.frame</code> as returned by <code>load_intdata(..., for.aa=FALSE)</code> . |
| V_segment | The name of a V gene segment. This can be set to "fwr1", "cdr1", "fwr2", "cdr2", or "fwr3". By default <code>translate_V_alleles()</code> will translate the entire coding frame in each allele. Otherwise, it will translate the specified segment only. |

Value

`translate_V_alleles()` returns a named character vector with 1 amino acid sequence per supplied allele. The vector contains an NA for any allele that is not annotated in `intdata` or for which the required information is NA. The names on it are the names of the supplied alleles.

`V_allele_has_stop_codon()` returns a named logical vector with 1 value per supplied allele. The vector contains an NA for any allele that is not annotated in `intdata` or for which `intdata$coding_frame_start` has an NA. The names on it are the names of the supplied alleles.

See Also

- [intdata_utils](#) to access IgBLAST *internal data*.
- [auxdata_utils](#) to access IgBLAST *auxiliary data*.
- [J_alleles_inspect](#) for basic inspection of J allele sequences.
- [update_live_igdata](#) for more information about "live" and "original" IgBLAST data.
- `DNAStrngSet` objects in the **Biostrings** package.
- The [translate_codons](#) function upon which `translate_V_alleles()` is based.
- [allele2gene](#) to go from germline gene allele names to germline gene names.
- IgBLAST is described at <https://pubmed.ncbi.nlm.nih.gov/23671333/>.

Examples

```
## Let's inspect the V allele sequences stored in
## germline db _OGRDB.human.IGH+IGK+IGL.202410.

## Load the V allele sequences:
db_name <- "_OGRDB.human.IGH+IGK+IGL.202410"
V_alleles <- load_germline_db(db_name, region_types="V")
V_alleles # DNAStrngSet object

## Load the corresponding internal data:
intdata <- load_intdata(db_name)

## Translate the V allele sequences:
V_aa <- translate_V_alleles(V_alleles, intdata)
head(V_aa)

## Translate the FWR2 portion of the V allele sequences:
fwr2_aa <- translate_V_alleles(V_alleles, intdata, V_segment="fwr2")
head(fwr2_aa)

## Surprisingly, 13 V alleles in _OGRDB.human.IGH+IGK+IGL.202410
## contain the stop codon:
has_stop_codon <- grepl("*", V_aa, fixed=TRUE)
table(has_stop_codon)

## Display the V allele sequences with stop codon:
V_aa[has_stop_codon] # protein
V_alleles[has_stop_codon] # nucleotides
```

Index

- * **manip**
 - igblastn, 26
 - read_igblastn_AIRR_output, 63
 - read_igblastn_fmt7_output, 64
- * **misc**
 - IGBLAST_ROOT, 35
- * **utilities**
 - allele2gene, 3
 - augment_germline_db, 4
 - auxdata-IO, 7
 - auxdata-utils, 8
 - compute_auxdata, 11
 - compute_V_gene_delineations, 13
 - download_IMGT_germline_sequences, 16
 - download_OGRDB_germline_json, 18
 - download_OGRDB_germline_sequences, 22
 - get_igblast_root, 25
 - igblast_info, 33
 - igblast_usage_report, 32
 - igbrowser, 36
 - install_custom_germline_db, 37
 - install_igblast, 42
 - install_IMGT_germline_db, 43
 - intdata-utils, 46
 - J_alleles-inspect, 49
 - list_c_region_dbs, 51
 - list_germline_dbs, 53
 - ndm_data-IO, 55
 - OAS-utils, 57
 - parse_imgt_fasta_headers, 60
 - read_igblastn_AIRR_output, 63
 - read_igblastn_fmt7_output, 64
 - reset_c_region_dbs, 66
 - reset_germline_dbs, 67
 - summarizeMismatches, 68
 - translate_codons, 68
 - update_live_igdata, 70
 - use_c_region_db, 72
 - use_germline_dbs, 74
 - V_alleles-inspect, 75
- AAStrngSet, 69
- allele2gene, 3, 76
- augment_germline_db, 4
- augment_germline_db_D
 - (augment_germline_db), 4
- augment_germline_db_J
 - (augment_germline_db), 4
- augment_germline_db_V
 - (augment_germline_db), 4
- auxdata (auxdata-utils), 8
- auxdata-IO, 7
- auxdata-utils, 8
- auxdata_IO (auxdata-IO), 7
- auxdata_utils, 8, 20, 34, 47, 49, 72, 76
- auxdata_utils (auxdata-utils), 8
- bcr_browser (igbrowser), 36
- BLAST_USAGE_REPORT
 - (igblast_usage_report), 32
- browseURL, 36
- compute_auxdata, 8–10, 11, 14, 20, 49
- compute_germline_db_auxdata, 49
- compute_germline_db_auxdata
 - (auxdata-utils), 8
- compute_imgt_intdata
 - (compute_V_gene_delineations), 13
- compute_V_gene_delineations, 12, 13, 20, 47, 56, 61
- difftime, 71
- DNAStrng, 69
- DNAStrngSet, 4, 11–14, 26, 30, 49, 57–59, 69, 73–76
- download.file, 17, 19, 23, 42, 44, 57
- download_IMGT_germline_sequences, 16, 24, 39, 43, 61
- download_OGRDB_germline_json, 18, 24
- download_OGRDB_germline_sequences, 18, 20, 22, 39
- download_paired_OAS_units (OAS-utils), 57
- extract_auxdata_from_ogrdb_json, 8

- extract_auxdata_from_ogrdb_json
(download_OGRDB_germline_json),
18
- extract_codons (translate_codons), 68
- extract_intdata_from_ogrdb_json, 56
- extract_intdata_from_ogrdb_json
(download_OGRDB_germline_json),
18
- extract_metadata_from_OAS_units
(OAS-utils), 57
- extract_sequences_from_paired_OAS_df
(OAS-utils), 57
- extract_sequences_from_paired_OAS_units
(OAS-utils), 57
- fmt7-utils (read_igblastn_fmt7_output),
64
- fmt7_utils (read_igblastn_fmt7_output),
64
- get_auxdata_path (auxdata-utils), 8
- get_igblast_auxiliary_data
(auxdata-utils), 8
- get_igblast_root, 25, 34, 71, 72
- get_intdata_path (intdata-utils), 46
- has_igblast (igblast_info), 33
- igblast_build (igblast_info), 33
- igblast_info, 25, 29, 33, 35, 42
- IGBLAST_ROOT, 25, 34, 35, 42
- igblastn, 5, 8–10, 12, 14, 25, 26, 33–39,
42–45, 47, 52–56, 59, 63, 65, 72–75
- igblastn_help (igblastn), 26
- igblastn_version (igblast_info), 33
- igblast_usage_report, 29, 30, 32
- igbrowser, 30, 36
- igdata_info (update_live_igdata), 70
- IMGT_is_up
(download_IMGT_germline_sequences),
16
- install_custom_germline_db, 18, 24, 29,
37, 44, 45, 54, 55, 68
- install_igblast, 25, 29, 34, 35, 42, 72
- install_IMGT_c_region_db, 67
- install_IMGT_c_region_db
(install_IMGT_germline_db), 43
- install_IMGT_germline_db, 18, 29, 39, 43,
54, 55, 68
- intdata (intdata-utils), 46
- intdata-utils, 46
- intdata_utils, 3, 10, 20, 34, 38, 44, 49, 55,
56, 72, 76
- intdata_utils (intdata-utils), 46
- IRangesList, 13
- J_allele_has_stop_codon
(J_alleles-inspect), 49
- J_alleles-inspect, 49
- J_alleles_inspect, 76
- J_alleles_inspect (J_alleles-inspect),
49
- list_c_region_dbs, 51, 54, 67, 73
- list_germline_dbs, 4, 5, 9, 10, 27, 28, 38,
39, 44–47, 52, 53, 68, 69, 74, 75
- list_igblast_organisms, 9, 10, 27, 30, 46,
47
- list_igblast_organisms (igblast_info),
33
- list_IMGT_organisms
(download_IMGT_germline_sequences),
16
- list_IMGT_releases, 45
- list_IMGT_releases
(download_IMGT_germline_sequences),
16
- list_outfmt7_specifiers, 26, 30
- list_outfmt7_specifiers
(read_igblastn_fmt7_output), 64
- list_paired_OAS_studies (OAS-utils), 57
- list_paired_OAS_units (OAS-utils), 57
- load_auxdata, 11, 12, 20, 49
- load_auxdata (auxdata-utils), 8
- load_c_region_db, 3
- load_c_region_db (use_c_region_db), 72
- load_germline_db, 3, 39
- load_germline_db (use_germline_dbs), 74
- load_igblast_auxiliary_data
(auxdata-utils), 8
- load_intdata, 13, 14, 20, 39, 75
- load_intdata (intdata-utils), 46
- makeblastdb_version (igblast_info), 33
- mcols, 13
- ndm_data-IO, 55
- ndm_data_IO (ndm_data-IO), 55
- OAS-utils, 57
- OAS_is_up (OAS-utils), 57
- OAS_utils (OAS-utils), 57
- outfmt7-utils
(read_igblastn_fmt7_output), 64
- outfmt7_utils
(read_igblastn_fmt7_output), 64

- parse_imgt_fasta_headers, 60
- parse_outfmt7
 - (read_igblastn_fmt7_output), 64
- print.alignment_summary
 - (read_igblastn_fmt7_output), 64
- print.auxdata_md5sum_df
 - (update_live_igdata), 70
- print.c_region_dbs_df
 - (list_c_region_dbs), 51
- print.fmt7footer
 - (read_igblastn_fmt7_output), 64
- print.fmt7record
 - (read_igblastn_fmt7_output), 64
- print.germline_dbs_df
 - (list_germline_dbs), 53
- print.hit_table
 - (read_igblastn_fmt7_output), 64
- print.igblast_info (igblast_info), 33
- print.igblastn_raw_output (igblastn), 26
- print.igdata_info (update_live_igdata), 70
- print.outfmt7_specifiers
 - (read_igblastn_fmt7_output), 64
- print.query_details
 - (read_igblastn_fmt7_output), 64
- print.subregion_sequence_details
 - (read_igblastn_fmt7_output), 64
- print.VDJ_junction_details
 - (read_igblastn_fmt7_output), 64
- print.VDJ_rearrangement_summary
 - (read_igblastn_fmt7_output), 64
- qseqid (read_igblastn_fmt7_output), 64
- read.table, 57
- read_auxdata, 9, 12
- read_auxdata (auxdata-IO), 7
- read_igblastn_AIRR_output, 63, 65
- read_igblastn_fmt7_output, 29, 63, 64
- read_ndm_data, 14, 47
- read_ndm_data (ndm_data-IO), 55
- read_OAS_csv (OAS-utils), 57
- read_OAS_csv_metadata (OAS-utils), 57
- remove_gaps (translate_codons), 68
- reset_c_region_dbs, 52, 66
- reset_germline_dbs, 55, 67
- reset_live_igdata (update_live_igdata), 70
- rm_c_region_db (list_c_region_dbs), 51
- rm_germline_db (list_germline_dbs), 53
- Rprofile, 33
- same_alleles_annot (ndm_data-IO), 55
- set_igblast_root (get_igblast_root), 25
- summarizeMismatches, 68
- summary.query_details
 - (read_igblastn_fmt7_output), 64
- tabulate_deletions
 - (summarizeMismatches), 68
- tabulate_insertions
 - (summarizeMismatches), 68
- tabulate_mismatches
 - (summarizeMismatches), 68
- tibble, 28–30, 36, 57–59
- time_since_live_igdata_last_checked
 - (update_live_igdata), 70
- translate, 69
- translate_codons, 49, 68, 76
- translate_fwr4 (J_alleles-inspect), 49
- translate_J_alleles
 - (J_alleles-inspect), 49
- translate_V_alleles
 - (V_alleles-inspect), 75
- update_live_igdata, 9, 10, 47, 49, 70, 76
- Usage_report (igblast_usage_report), 32
- usage_report (igblast_usage_report), 32
- Usage_reporting
 - (igblast_usage_report), 32
- usage_reporting
 - (igblast_usage_report), 32
- use_c_region_db, 27, 30, 52, 67, 72, 75
- use_germline_db, 27, 28, 30, 39, 45, 55, 67, 68, 73
- use_germline_db (use_germline_dbs), 74
- use_germline_dbs, 74
- V_allele_has_stop_codon
 - (V_alleles-inspect), 75
- V_alleles-inspect, 75
- V_alleles_inspect, 49
- V_alleles_inspect (V_alleles-inspect), 75
- V_genes_with_varying_fwrcdr_boundaries
 - (intdata-utils), 46
- validate_ndm_rows (ndm_data-IO), 55
- write_auxdata, 20
- write_auxdata (auxdata-IO), 7
- write_ndm_data, 20
- write_ndm_data (ndm_data-IO), 55