

Quick Guide to AssessORFData

Deepank Korandla

11 February 2025

Package

AssessORFData 1.25.1

Contents

1	Introduction	2
2	Strains and Gene Sources	2
3	Getting the Objects	2
4	Saving the Genome	3
5	Session Info	4

1 Introduction

AssessORFData is an accompaniment to the AssessORF package, providing access to mapping and results objects generated by AssessORF as well as the genome sequences for the strains corresponding to those objects. Briefly, a mapping object stores the mapping of proteomics evidence and evolutionary conservation evidence to a particular strain's genome, and a results object stores how much evidence there is supporting or against each gene in a set of predicted genes for a particular strain's genome. Detailed descriptions of the structure and content of those two types of objects can be found in the documentation for the AssessORF package.

2 Strains and Gene Sources

AssessORFData has data for 20 strains and their IDs (within the package) are listed below:

```
AssessORFData::GetStrainIDs()
## [1] "ATCC11842" "ATCC13032" "ATCC17978" "ATCC700084" "CCMP1375"
## [6] "CECT5344" "CNRZ327" "COH1" "D_UW_3_CX" "EGD_e"
## [11] "H37Rv" "Houston_1" "I11403" "K_12_MG1655" "LAL14_1"
## [16] "MGAS5005" "PA01" "SL1344" "Strain168" "TCH1516"
## [21] "AP1" "BW25113" "HG001" "MG1363" "NCIB_3610"
## [26] "Strain10403S"
```

For each strain, there are 5 objects: 1 mapping object and 4 results objects. The 4 results objects per strain differ in that for each one, the set of predicted genes came from a different program or database. The same 4 gene sources were used for all 20 strains, and their IDs (within the package) are listed below:

```
AssessORFData::GetGeneSources()
## [1] "Prodigal" "GeneMarkS2" "GenBank" "Glimmer"
```

3 Getting the Objects

While the `data` function can be used to pull the desired object from the package into the user's workspace, using the `data` function may be inconvenient for some users because there are 100 mapping and results objects, each of which has a long name. For this reason, AssessORFData has alternative functions, `GetDataMapObj` and `GetResultsObj`, to accomplish a similar task. These functions allow the user to get the object of interest and then save it in their workspace under a different name. Examples on how to use the two functions are described below:

```
library(AssessORFData)
## Loading required package: RSQLite

## Can replace the character string specifying the strain ID (first
## parameter) with any of the other 19 strain IDs listed above
mapObj <- GetDataMapObj("MGAS5005")
resObj1 <- GetResultsObj("MGAS5005", "Prodigal")
resObj2 <- GetResultsObj("MGAS5005", "GenBank")
resObj3 <- GetResultsObj("MGAS5005", "GeneMarkS2")
resObj4 <- GetResultsObj("MGAS5005", "Glimmer")
```

4 Saving the Genome

The `SaveGenomeToPath` function allows the user to save the genome for a strain of their choosing to a file path on their local machine in situations where the user wants to run their own analyses on their strain's genome, e.g. predict genes for the genome using a different gene finding program. An example of how to use the function is provided below:

```
library(AssessORFData)

## A path to a temporary file is used in this example.
tmpFile <- paste0(tempfile(), ".fasta")

## Replace the second parameter below with a character string specifying
## the desired file path, making sure the file is of type FASTA.
SaveGenomeToPath("MGAS5005", tmpFile)

unlink(tmpFile)
```

5 Session Info

All of the output in this vignette was produced under the following conditions:

```
## R Under development (unstable) (2025-01-20 r87609)
## Platform: x86_64-pc-linux-gnu
## Running under: Ubuntu 24.04.1 LTS
##
## Matrix products: default
## BLAS: /home/biocbuild/bbs-3.21-bioc/R/lib/libRblas.so
## LAPACK: /usr/lib/x86_64-linux-gnu/lapack/liblapack.so.3.12.0 LAPACK version 3.12.0
##
## attached base packages:
## [1] stats      graphics  grDevices  utils      datasets  methods   base
##
## other attached packages:
## [1] AssessORFData_1.25.1 RSQLite_2.3.9      BiocStyle_2.35.0
##
## loaded via a namespace (and not attached):
## [1] vctrs_0.6.5      crayon_1.5.3      http_1.4.7
## [4] cli_3.6.3        knitr_1.49        rlang_1.1.5
## [7] xfun_0.50        DBI_1.2.3         UCSC.utils_1.3.1
## [10] generics_0.1.3   jsonlite_1.8.9    bit_4.5.0.1
## [13] S4Vectors_0.45.2 Biostrings_2.75.3  htmltools_0.5.8.1
## [16] stats4_4.5.0     rmarkdown_2.29    evaluate_1.0.3
## [19] DECIPHER_3.3.2   fastmap_1.2.0     yaml_2.3.10
## [22] IRanges_2.41.2   GenomeInfoDb_1.43.4 memoise_2.0.1
## [25] bookdown_0.42    BiocManager_1.30.25 compiler_4.5.0
## [28] blob_1.2.4       XVector_0.47.2    digest_0.6.37
## [31] R6_2.5.1         GenomeInfoDbData_1.2.13 bit64_4.6.0-1
## [34] tools_4.5.0      cachem_1.1.0      BiocGenerics_0.53.6
```