

AnnotationDbi: Introduction To Bioconductor Annotation Packages

Marc Carlson

December 23, 2015

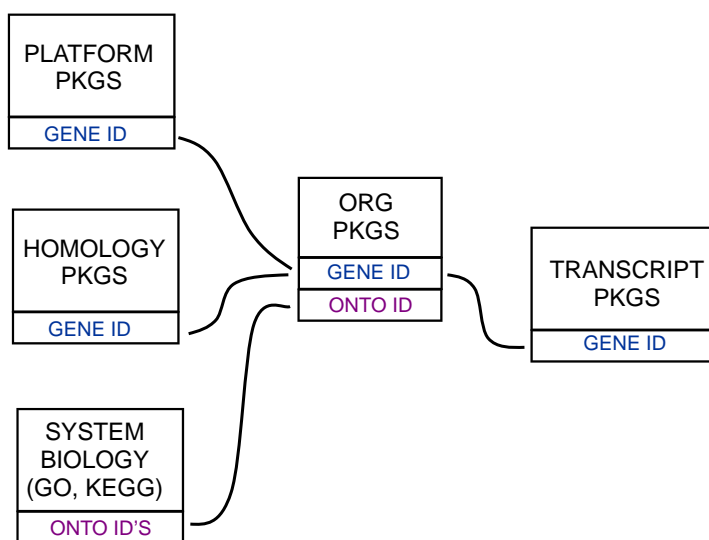


Figure 1: Annotation Packages: the big picture

Bioconductor provides extensive annotation resources. These can be *gene centric*, or *genome centric*. Annotations can be provided in packages curated by *Bioconductor*, or obtained from web-based resources. This vignette is primarily concerned with describing the annotation resources that are available as packages. More advanced users who wish to learn about how to make new annotation packages should see the vignette titled "Creating select Interfaces for custom Annotation resources" from the *AnnotationForge* package.

Gene centric *AnnotationDbi* packages include:

- Organism level: e.g. *org.Mm.eg.db*.
- Platform level: e.g. *hgu133plus2.db*, *hgu133plus2.probes*, *hgu133plus2.cdf*.
- Homology level: e.g. *hom.Dm.inp.db*.
- System-biology level: *GO.db*

Genome centric *GenomicFeatures* packages include

- Transcriptome level: e.g. *TxDb.Hsapiens.UCSC.hg19.knownGene*
- Generic genome features: Can generate via *GenomicFeatures*

One web-based resource accesses [biomart](#), via the *biomaRt* package:

- Query web-based 'biomart' resource for genes, sequence, SNPs, and etc.

The most popular annotation packages have been modified so that they can make use of a new set of methods to more easily access their contents. These four methods are named: `columns`, `keytypes`, `keys` and `select`. And they are described in this vignette. They can currently be used with all `chip`, `organism`, and `TxDb` packages along with the popular `GO.db` package.

For the older less popular packages, there are still convenient ways to retrieve the data. The *How to use bimap from the ".db" annotation packages* vignette in the `AnnotationDbi` package is a key reference for learnign about how to use `bimap` objects.

Finally, all of the '.db' (and most other *Bioconductor* annotation packages) are updated every 6 months corresponding to each release of *Bioconductor*. Exceptions are made for packages where the actual resources that the packages are based on have not themselves been updated.

0.1 AnnotationDb objects and the select method

As previously mentioned, a new set of methods have been added that allow a simpler way of extracting identifier based annotations. All the annotation packages that support these new methods expose an object named exactly the same as the package itself. These objects are collectively called *AnntoationDb* objects for the class that they all inherit from. The more specific classes (the ones that you will actually see in the wild) have names like `OrgDb`, `ChipDb` or `TxDb` objects. These names correspond to the kind of package (and underlying schema) being represented. The methods that can be applied to all of these objects are `columns`, `keys`, `keytypes` and `select`.

In addition, another accessor has recently been added which allows extraction of one column at at time. the `mapIds` method allows users to extract data into either a named character vector, a list or even a `SimpleCharacterList`. This method should work with all the different kinds of *AnntoationDb* objects described below.

0.2 ChipDb objects and the select method

An extremely common kind of Annotation package is the so called platform based or chip based package type. This package is intended to make the manufacturer labels for a series of probes or probesets to a wide range of gene-based features. A package of this kind will load an `ChipDb` object. Below is a set of examples to show how you might use the standard 4 methods to interact with an object of this type.

First we need to load the package:

```
library(hgu95av2.db)
```

If we list the contents of this package, we can see that one of the many things loaded is an object named after the package "hgu95av2.db":

```
ls("package:hgu95av2.db")
```

```
## [1] "hgu95av2"                "hgu95av2.db"            "hgu95av2ACCNUM"
## [4] "hgu95av2ALIAS2PROBE"    "hgu95av2CHR"           "hgu95av2CHRENGTHS"
## [7] "hgu95av2CHRLOC"        "hgu95av2CHRLOCEND"    "hgu95av2ENSEMBL"
## [10] "hgu95av2ENSEMBL2PROBE" "hgu95av2ENTREZID"     "hgu95av2ENZYME"
```

```
## [13] "hgu95av2ENZYME2PROBE" "hgu95av2GENENAME" "hgu95av2G0"
## [16] "hgu95av2G02ALLPROBES" "hgu95av2G02PROBE" "hgu95av2MAP"
## [19] "hgu95av2MAPCOUNTS" "hgu95av2OMIM" "hgu95av2ORGANISM"
## [22] "hgu95av2ORGPKG" "hgu95av2PATH" "hgu95av2PATH2PROBE"
## [25] "hgu95av2PFAM" "hgu95av2PMID" "hgu95av2PMID2PROBE"
## [28] "hgu95av2PROSITE" "hgu95av2REFSEQ" "hgu95av2SYMBOL"
## [31] "hgu95av2UNIGENE" "hgu95av2UNIPROT" "hgu95av2_dbInfo"
## [34] "hgu95av2_dbconn" "hgu95av2_dbfile" "hgu95av2_dbschema"
```

We can look at this object to learn more about it:

```
hgu95av2.db

## ChipDb object:
## | DBSCHEMAVERSION: 2.1
## | Db type: ChipDb
## | Supporting package: AnnotationDbi
## | DBSCHEMA: HUMANCHIP_DB
## | ORGANISM: Homo sapiens
## | SPECIES: Human
## | MANUFACTURER: Affymetrix
## | CHIPNAME: Human Genome U95 Set
## | MANUFACTURERURL: http://www.affymetrix.com/support/technical/byproduct.affx?product=hgu95
## | EGSOURCEDATE: 2015-Sep27
## | EGSOURCENAME: Entrez Gene
## | EGSOURCEURL: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA
## | CENTRALID: ENTREZID
## | TAXID: 9606
## | GOSOURCENAME: Gene Ontology
## | GOSOURCEURL: ftp://ftp.geneontology.org/pub/go/godatabase/archive/latest-lite/
## | GOSOURCEDATE: 20150919
## | GOEGSOURCEDATE: 2015-Sep27
## | GOEGSOURCEURL: ftp://ftp.ncbi.nlm.nih.gov/gene/DATA
## | KEGGSOURCEURL: ftp://ftp.genome.jp/pub/kegg/genomes
## | KEGGSOURCEDATE: 2011-Mar15
## | GPSOURCEURL: ftp://hgdownload.cse.ucsc.edu/goldenPath/hg19
## | GPSOURCEDATE: 2010-Mar22
## | ENSOURCEDATE: 2015-Jul16
## | ENSOURCEURL: ftp://ftp.ensembl.org/pub/current_fasta
## | UPSOURCEURL: http://www.uniprot.org/
## | UPSOURCEDATE: Thu Oct 1 23:31:58 2015

##
## Please see: help('select') for usage information
```


And as you can see, when you call the code above, `select` will try to return a data.frame with all the things you asked for matched up to each other.

Finally if you wanted to extract only one column of data you could instead use the `mapIds` method like this:

```
#1st get some example keys
k <- head(keys(hgu95av2.db, keytype="PROBEID"))
# then call mapIds
mapIds(hgu95av2.db, keys=k, column=c("GENENAME"), keytype="PROBEID")

##                               1000_at
##                               "mitogen-activated protein kinase 3"
##                               1001_at
## "tyrosine kinase with immunoglobulin-like and EGF-like domains 1"
##                               1002_f_at
## "cytochrome P450, family 2, subfamily C, polypeptide 19"
##                               1003_s_at
## "chemokine (C-X-C motif) receptor 5"
##                               1004_at
## "chemokine (C-X-C motif) receptor 5"
##                               1005_at
##                               "dual specificity phosphatase 1"
```

0.3 OrgDb objects and the select method

An organism level package (an 'org' package) uses a central gene identifier (e.g. Entrez Gene id) and contains mappings between this identifier and other kinds of identifiers (e.g. GenBank or Uniprot accession number, RefSeq id, etc.). The name of an org package is always of the form *org.iAbj.jidj.db* (e.g. *org.Sc.sgd.db*) where *iAbj* is a 2-letter abbreviation of the organism (e.g. *Sc* for *Saccharomyces cerevisiae*) and *jidj* is an abbreviation (in lower-case) describing the type of central identifier (e.g. *sgd* for gene identifiers assigned by the Saccharomyces Genome Database, or *eg* for Entrez Gene ids).

Just as the chip packages load a *ChipDb* object, the org packages will load a *OrgDb* object. The following exercise should acquaint you with the use of these methods in the context of an organism package.

Exercise 1

Display the *OrgDb* object for the *org.Hs.eg.db* package.

Use the `columns` method to discover which sorts of annotations can be extracted from it. Is this the same as the result from the `keytypes` method? Use the `keytypes` method to find out.

Finally, use the `keys` method to extract UNIPROT identifiers and then pass those keys in to the `select` method in such a way that you extract the gene symbol and KEGG pathway information for each. Use the help system as needed to learn which values to pass in to `columns` in order to achieve this.

Solution:

```
library(org.Hs.eg.db)
columns(org.Hs.eg.db)

## [1] "ACCNUM" "ALIAS" "ENSEMBL" "ENSEMBLPROT" "ENSEMBLTRANS"
```

```
## [6] "ENTREZID"      "ENZYME"      "EVIDENCE"    "EVIDENCEALL" "GENENAME"
## [11] "GO"            "GOALL"       "IPI"         "MAP"          "OMIM"
## [16] "ONTOLOGY"      "ONTOLOGYALL" "PATH"        "PFAM"         "PMID"
## [21] "PROSITE"       "REFSEQ"      "SYMBOL"      "UCSCKG"      "UNIGENE"
## [26] "UNIPROT"
```

```
help("SYMBOL") ## for explanation of these columns and keytypes values
```

```
keytypes(org.Hs.eg.db)
```

```
## [1] "ACCNUM"      "ALIAS"       "ENSEMBL"     "ENSEMBLPROT" "ENSEMBLTRANS"
## [6] "ENTREZID"    "ENZYME"      "EVIDENCE"    "EVIDENCEALL" "GENENAME"
## [11] "GO"          "GOALL"       "IPI"         "MAP"          "OMIM"
## [16] "ONTOLOGY"    "ONTOLOGYALL" "PATH"        "PFAM"         "PMID"
## [21] "PROSITE"     "REFSEQ"      "SYMBOL"      "UCSCKG"      "UNIGENE"
## [26] "UNIPROT"
```

```
uniKeys <- head(keys(org.Hs.eg.db, keytype="UNIPROT"))
```

```
cols <- c("SYMBOL", "PATH")
```

```
select(org.Hs.eg.db, keys=uniKeys, columns=cols, keytype="UNIPROT")
```

```
## 'select()' returned 1:many mapping between keys and columns
```

```
##   UNIPROT SYMBOL  PATH
## 1  P04217  A1BG  <NA>
## 2  V9HWD8  A1BG  <NA>
## 3  P01023   A2M 04610
## 4  P18440  NAT1 00232
## 5  P18440  NAT1 00983
## 6  P18440  NAT1 01100
## 7  Q400J6  NAT1 00232
## 8  Q400J6  NAT1 00983
## 9  Q400J6  NAT1 01100
## 10 F5H5R8  NAT1 00232
## 11 F5H5R8  NAT1 00983
## 12 F5H5R8  NAT1 01100
```

So how could you use `select` to annotate your results? This next exercise should help you to understand how that should generally work.

Exercise 2

Please run the following code snippet (which will load a fake data result that I have provided for the purposes of illustration):

```
load(system.file("extdata", "resultTable.Rda", package="AnnotationDbi"))
```

```
head(resultTable)
```

```
##           logConc      logFC LR.statistic      PValue      FDR
## 100418920 -9.639471 -4.679498    378.0732 3.269307e-84 2.613484e-80
## 100419779 -10.638865 -4.264830    291.1028 2.859424e-65 1.142912e-61
## 100271867 -11.448981 -4.009603    222.3653 2.757135e-50 7.346846e-47
```

```
## 100287169 -11.026699 -3.486593      206.7771 6.934967e-47 1.385953e-43
## 100287735 -11.036862  3.064980      204.1235 2.630432e-46 4.205535e-43
## 100421986 -12.276297 -4.695736      190.5368 2.427556e-43 3.234314e-40
```

The rownames of this table happen to provide entrez gene identifiers for each row (for human). Find the gene symbol and gene name for each of the rows in resultTable and then use the merge method to attach those annotations to it.

Solution:

```
annots <- select(org.Hs.eg.db, keys=rownames(resultTable),
                 columns=c("SYMBOL","GENENAME"), keytype="ENTREZID")
## 'select()' returned 1:1 mapping between keys and columns

resultTable <- merge(resultTable, annots, by.x=0, by.y="ENTREZID")
head(resultTable)

##   Row.names  logConc    logFC LR.statistic      PValue      FDR      SYMBOL
## 1 100127888 -10.57050  2.758937    182.8937 1.131473e-41 1.130624e-38 SLC04A1-AS1
## 2 100131223 -12.37808 -4.654318    179.2331 7.126423e-41 6.329847e-38 LOC100131223
## 3 100271381 -12.06340  3.511937    188.4824 6.817155e-43 7.785191e-40 RPS28P8
## 4 100271867 -11.44898 -4.009603    222.3653 2.757135e-50 7.346846e-47 MPVQTL1
## 5 100287169 -11.02670 -3.486593    206.7771 6.934967e-47 1.385953e-43 <NA>
## 6 100287735 -11.03686  3.064980    204.1235 2.630432e-46 4.205535e-43 TTY13B
##
##                               GENENAME
## 1                               SLC04A1 antisense RNA 1
## 2 ADP-ribosylation factor-like 8B pseudogene
## 3           ribosomal protein S28 pseudogene 8
## 4                               Mean platelet volume QTL1
## 5                               <NA>
## 6 testis-specific transcript, Y-linked 13B
```

0.4 Using select with GO.db

When you load the GO.db package, a *GODb* object is also loaded. This allows you to use the columns, keys, keytypes and select methods on the contents of the GO ontology. So if for example, you had a few GO IDs and wanted to know more about it, you could do it like this:

```
library(GO.db)
GOIDs <- c("GO:0042254", "GO:0044183")
select(GO.db, keys=GOIDs, columns="DEFINITION", keytype="GOID")
## 'select()' returned 1:1 mapping between keys and columns

##           GOID
## 1 GO:0042254
## 2 GO:0044183
##
```

```
## 1      A cellular process that results in the biosynthesis of constituent macromolecules, asse
## 2 Interacting selectively and non-covalently with any protein or protein complex (a complex of
```

0.5 Using select with TxDb packages

A *TxDb* package (a 'TxDb' package) connects a set of genomic coordinates to various transcript oriented features. The package can also contain Identifiers to features such as genes and transcripts, and the internal schema describes the relationships between these different elements. All TxDb containing packages follow a specific naming scheme that tells where the data came from as well as which build of the genome it comes from.

Exercise 3

Display the TxDb object for the *TxDb.Hsapiens.UCSC.hg19.knownGene* package.

As before, use the *columns* and *keytypes* methods to discover which sorts of annotations can be extracted from it.

Use the *keys* method to extract just a few gene identifiers and then pass those keys in to the *select* method in such a way that you extract the transcript ids and transcript starts for each.

Solution:

```
library(TxDb.Hsapiens.UCSC.hg19.knownGene)

## Loading required package: GenomicFeatures
## Loading required package: GenomeInfoDb
## Loading required package: GenomicRanges

txdb <- TxDb.Hsapiens.UCSC.hg19.knownGene
txdb

## TxDb object:
## # Db type: TxDb
## # Supporting package: GenomicFeatures
## # Data source: UCSC
## # Genome: hg19
## # Organism: Homo sapiens
## # Taxonomy ID: 9606
## # UCSC Table: knownGene
## # Resource URL: http://genome.ucsc.edu/
## # Type of Gene ID: Entrez Gene ID
## # Full dataset: yes
## # miRBase build ID: GRCh37
## # transcript_nrow: 82960
## # exon_nrow: 289969
## # cds_nrow: 237533
## # Db created by: GenomicFeatures package from Bioconductor
## # Creation time: 2015-10-07 18:11:28 +0000 (Wed, 07 Oct 2015)
## # GenomicFeatures version at creation time: 1.21.30
## # RSQLite version at creation time: 1.0.0
```



```
## # DBSCHEMAVERSION: 1.1

columns(txdb)

## [1] "CDSCHROM" "CSEND" "CDSID" "CDSNAME" "CDSSTART" "CDSSTRAND"
## [7] "EXONCHROM" "EXONEND" "EXONID" "EXONNAME" "EXONRANK" "EXONSTART"
## [13] "EXONSTRAND" "GENEID" "TXCHROM" "TXEND" "TXID" "TXNAME"
## [19] "TXSTART" "TXSTRAND" "TXTYPE"

keytypes(txdb)

## [1] "CDSID" "CDSNAME" "EXONID" "EXONNAME" "GENEID" "TXID" "TXNAME"

keys <- head(keys(txdb, keytype="GENEID"))
cols <- c("TXID", "TXSTART")
select(txdb, keys=keys, columns=cols, keytype="GENEID")

## 'select()' returned 1:many mapping between keys and columns

##      GENEID  TXID  TXSTART
## 1         1 70455 58858172
## 2         1 70456 58859832
## 3        10 31944 18248755
## 4        100 72132 43248163
## 5       1000 65378 25530930
## 6       1000 65379 25530930
## 7      10000 7895 243651535
## 8      10000 7896 243663021
## 9      10000 7897 243663021
## 10 100008586 75890 49217763
```

As is widely known, in addition to providing access via the `select` method, `TxDb` objects also provide access via the more familiar `transcripts`, `exons`, `cds`, `transcriptsBy`, `exonsBy` and `cdsBy` methods. For those who do not yet know about these other methods, more can be learned by seeing the vignette called: *Making and Utilizing TxDb Objects* in the *GenomicFeatures* package.

The version number of R and packages loaded for generating the vignette were:

```
## R version 3.2.3 (2015-12-10)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 14.04.3 LTS
##
## locale:
## [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C              LC_TIME=en_US.UTF-8
## [4] LC_COLLATE=C             LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
## [7] LC_PAPER=en_US.UTF-8     LC_NAME=C                 LC_ADDRESS=C
## [10] LC_TELEPHONE=C          LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
## [1] parallel stats4 stats graphics grDevices utils datasets methods
## [9] base
```

```
##
## other attached packages:
## [1] TxDb.Hsapiens.UCSC.hg19.knownGene_3.2.2 GenomicFeatures_1.22.7
## [3] GenomicRanges_1.22.2 GenomeInfoDb_1.6.1
## [5] GO.db_3.2.2 hgu95av2.db_3.2.2
## [7] AnnotationForge_1.12.1 org.Hs.eg.db_3.2.3
## [9] RSQLite_1.0.0 DBI_0.3.1
## [11] AnnotationDbi_1.32.3 IRanges_2.4.6
## [13] S4Vectors_0.8.5 Biobase_2.30.0
## [15] BiocGenerics_0.16.1 knitr_1.11
##
## loaded via a namespace (and not attached):
## [1] XVector_0.10.0 magrittr_1.5 GenomicAlignments_1.6.1
## [4] zlibbioc_1.16.0 BiocParallel_1.4.3 stringr_1.0.0
## [7] highr_0.5.1 tools_3.2.3 SummarizedExperiment_1.0.1
## [10] lambda.r_1.1.7 futile.logger_1.4.1 rtracklayer_1.30.1
## [13] formatR_1.2.1 futile.options_1.0.0 bitops_1.0-6
## [16] RCurl_1.95-4.7 biomaRt_2.26.1 evaluate_0.8
## [19] stringi_1.0-1 Rsamtools_1.22.0 Biostrings_2.38.2
## [22] XML_3.98-1.3 BiocStyle_1.8.0
```