

Outline

- 1 **Bioconductor Annotation Packages**
- 2 AnnotationDbi
- 3 AnnotationDbi Basics
- 4 Working with GO.db

Bioconductor annotation packages

Major types of annotation in Bioconductor.

AnnotationDbi packages:

- Organism level: [org.Mm.eg.db](#).
- Platform level: [hgu133plus2.db](#).
- System-biology level: [GO.db](#) or [KEGG.db](#).

biomaRt:

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

Other packages:

- [rtracklayer](#) – export to UCSC web browsers.
- [GenomicFeatures](#) – coming soon for transcripts.

Outline

- 1 Bioconductor Annotation Packages
- 2 AnnotationDbi**
- 3 AnnotationDbi Basics
- 4 Working with GO.db

AnnotationDbi

AnnotationDbi is a software package that enables the package annotations:

- Each supported package contains a database.
- AnnotationDbi allows access to that data via Bimap objects.
- Some databases depend on the databases in other packages.

Organism-level annotation

There are a number of organism annotation packages with names starting with `org`, e.g., `org.Hs.eg.db` – genome-wide annotation for human.

```
> library(org.Hs.eg.db)
> org.Hs.eg()
> org.Hs.eg_dbInfo()
> org.Hs.egGENENAME
> org.Hs.eg_dbschema()
```

platform based packages (chip packages)

There are a number of platform or chip specific annotation packages named after their respective platforms, e.g. `hgu95av2.db` - annotations for the hgu95av2 Affymetrix platform.

- These packages appear to contain a lot of data but it's an illusion.

```
> library(hgu95av2.db)
> hgu95av2()
> hgu95av2_dbInfo()
> hgu95av2GENENAME
> hgu95av2_dbschema()
```

Kinds of annotation

What can you hope to extract from an annotation package?

- GO IDs: [GO](#)
- KEGG pathway IDs: [KEGG](#)
- Gene Symbols: [SYMBOL](#)
- Chromosome start and stop locs: [CHRLOC](#) and [CHRLOCEND](#)
- Alternate Gene Symbols: [ALIAS](#)
- Associated Pubmed IDs: [PMID](#)
- RefSeq IDs: [REFSEQ](#)
- Unigene IDs: [UNIGENE](#)
- PFAM IDs: [PFAM](#)
- Prosite IDs: [PROSITE](#)
- ENSEMBL IDs: [ENSEMBL](#)

Outline

- 1 Bioconductor Annotation Packages
- 2 AnnotationDbi
- 3 AnnotationDbi Basics**
- 4 Working with GO.db

Basic Bimap structure and getters

Bimaps create a mapping from one set of keys to another. And they can easily be searched.

- `toTable`: converts a Bimap to a `data.frame`
- `get`: pulls data from a Bimap
- `mget`: pulls data from a Bimap for multiple things at once

```
> head(toTable(hgu95av2SYMBOL))  
> get("38187_at",hgu95av2SYMBOL)  
> mget(c("38912_at","38187_at"),hgu95av2SYMBOL,ifnotfound=NA)
```

Reversing and subsetting Bimaps

Bimaps can also be reversed and subsetted:

- `revmap`: reverses a Bimap
- `[[, [`: Bimaps are subsettable.

```
> ##revmap
> mget(c("NAT1", "NAT2"), revmap(hgu95av2SYMBOL), ifnotfound=NA)
> ##subsetting
> head(toTable(hgu95av2SYMBOL[1:3]))
> hgu95av2SYMBOL[["1000_at"]]
> revmap(hgu95av2SYMBOL)[["MAPK3"]]
> ##Or you can combine things
> toTable(hgu95av2SYMBOL[c("38912_at", "38187_at")])
```

Annotation exercise 1

Find the gene symbol, chromosome position and KEGG pathway ID for "1003_s_at".

Annotation exercise 1 solution

```
> get("1003_s_at",hgu95av2SYMBOL)
> get("1003_s_at",hgu95av2CHRLoc)
> get("1003_s_at",hgu95av2PATH)
```

Bimap keys

Bimaps create a mapping from one set of keys to another. Some important methods include:

- `keys`: centralID for the package (directional)
- `Lkeys`: centralID for the package (probe ID or gene ID)
- `Rkeys`: centralID for the package (attached data)

> `keys(hgu95av2SYMBOL[1:4])`

> `Lkeys(hgu95av2SYMBOL[1:4])`

> `Rkeys(hgu95av2SYMBOL)[1:4]`

More Bimap structure

Not all keys have a partner (or are mapped)

- `mappedkeys`: which of the key are mapped (directional)
- `mappedLkeys` `mappedRkeys`: which keys are mapped (absolute reference)
- `count.mappedkeys`: Number of mapped keys (directional)
- `count.mappedLkeys`, `count.mappedRkeys`: Number of mapped keys (absolute)

> `mappedkeys(hgu95av2SYMBOL[1:10])`

> `mappedLkeys(hgu95av2SYMBOL[1:10])`

> `mappedRkeys(hgu95av2SYMBOL[1:10])`

> `count.mappedkeys(hgu95av2SYMBOL[1:100])`

> `count.mappedLkeys(hgu95av2SYMBOL[1:100])`

> `count.mappedRkeys(hgu95av2SYMBOL[1:100])`

Annotation exercise 2

Gene symbols are often recycled by other genes making them a poor choice for identifiers. Using what you have learned, the SYMBOL Bimap, along with the `lapply`, `length` and `sort` functions, determine which gene symbols in `hgu95av2` are the worst offenders.

Annotation exercise 2 solution

```
> badRank <- lapply(as.list(revmap(hgu95av2SYMBOL)), length)
> tail(sort(unlist(badRank)))
```


Bimap Conversions

How to handle conversions from Bimaps to lists

- `as.list`: converts a Bimap to a *list*
- `unlist2`: unlists a *list* minus the name-mangling.

```
> as.list(hgu95av2SYMBOL[c("38912_at", "38187_at")])
> unlist(as.list(hgu95av2SYMBOL[c("38912_at", "38187_at")]))
> unlist2(as.list(hgu95av2SYMBOL[c("38912_at", "38187_at")]))
> ##but what happens when there are
> ##repeating values for the left key?
> unlist(as.list(revmap(hgu95av2SYMBOL)[c("STAT1", "PTGER3")]))
> ##unlist2 can help with this
> unlist2(as.list(revmap(hgu95av2SYMBOL)[c("STAT1", "PTGER3")]))
```

toggleProbes

How to hide/unhide ambiguous probes.

- `toggleProbes`: hides or displays the probes that have multiple mappings to genes.

```
> ## How many probes?
> dim(hgu95av2ENTREZID)
> ## Make a mapping with multiple probes exposed
> multi <- toggleProbes(hgu95av2ENTREZID, "all")
> ## How many probes?
> dim(multi)
> ## Make a mapping with ONLY multiple probes exposed
> multiOnly <- toggleProbes(multi, "multiple")
> ## How many probes?
> dim(multiOnly)
> ## Then make a mapping with ONLY single mapping probes
> singleOnly <- toggleProbes(multiOnly, "single")
> ## How many probes?
> dim(singleOnly)
```

Annotation exercise 3

Using the knowledge that Entrez IDs are good IDs that can be used to define genes uniquely, find the probe that maps to the largest number of different genes on hgu95av2.

Annotation exercise 3 solution

```
> mult <- toggleProbes(hgu95av2ENTREZID, "multi")
> dim(mult)
> multRank <- lapply(as.list(mult), length)
> tail(sort(unlist(multRank)))
```

Outline

- 1 Bioconductor Annotation Packages
- 2 AnnotationDbi
- 3 AnnotationDbi Basics
- 4 Working with GO.db**

GO

Some important considerations about the Gene Ontology

- GO is actually 3 ontologies (CC, BP and MF)
- Each ontology is a directed acyclic graph.
- The structure of GO is maintained separately from the genes that these GO IDs are usually used to annotate.

GO to gene mappings are stored in other packages

Mapping Entrez IDs to GO

- Each ENTREZ ID is associated with up to three GO categories.
- The objects returned from an ordinary GO mapping are complex.

```
> go <- org.Hs.egGO[["1000"]]
> length(go)
> go[[2]]$GOID
> go[[2]]$Ontology
```

Annotation exercise 4

Use what they have learned to write a function that 1) gets the GOIDs for a particular entrez gene ID, and then returns only their GOID as a named vector.

Annotation exercise 4 solution

```
> ##get GOIDs from Hs package.  
> getGOIDs <- function(ids){  
+   require(org.Hs.eg.db)  
+   GOs = mget(ids, org.Hs.egGO, ifnotfound=NA)  
+   unlist2(sapply(GOs,names))  
+ }  
> ##usage example:  
> getGOIDs(c("1","10"))
```

Working with GO.db

- Encodes the hierarchical structure of GO terms.
- The mapping between GO terms and individual genes is maintained in the GO mappings from the other packages.
- the difference between children and offspring is how many generations are represented. Children only nets you one step down the graph.

```
> library(GO.db)
> ls("package:GO.db")
> ## find children
> as.list(GOMFCHILDREN["GO:0008094"])
> ## all the descendants (children, grandchildren, and so on)
> as.list(GOMFOFFSPRING["GO:0008094"])
```

GO helper methods

Using the GO helper methods

- The GO terms are described in detail in the [GOTERM](#) mapping.
- The objects returned by GO.db are GOTerms objects, which can make use of helper methods like `GOID`, `Term`, `RfunctionOntology` and `Definition` to retrieve various details.
- You can also pass GOIDs to these helper methods.

```
> ##Mapping a GOTerms object
> go <- GOTERM[1]
> GOID(go)
> Term(go)
> ##OR you can supply GO IDs
> id = c("GO:0007155", "GO:0007156")
> GOID(id)
> Term(id)
> Ontology(id)
> Definition(id)
```

Annotation exercise 5

Use what they have learned to write a function that calls your previous function so that it can get GOIDs and then returns the GO definitions.

Annotation exercise 5 solution

```
> ##get GOIDs from Hs package.
> getGODEfs <- function(ids){
+   GOids <- getGOIDs(ids)
+   defs <- Definition(GOids)
+   names(defs) <- names(GOids)
+   defs
+ }
> ##usage example:
> getGODEfs(c("1", "10"))
```