

OrganismDbi: A meta framework for Annotation Packages

Marc Carlson

October 14, 2013

OrganismDbi is a software package that helps tie together different annotation resources. It is expected that users may have previously made seen packages like *org.Hs.eg.db* and *TxDb.Hsapiens.UCSC.hg19.knownGene*. Packages like these two are very different and contain very different kinds of information, but are still about the same organism: Homo sapiens. The *OrganismDbi* package allows us to combine resources like these together into a single package resource, which can represent ALL of these resources at the same time. An example of this is the *homo.sapiens* package, which combines access to the two resources above along with others.

This is made possible because the packages that are represented by *homo.sapiens* are related to each other via foreign keys.

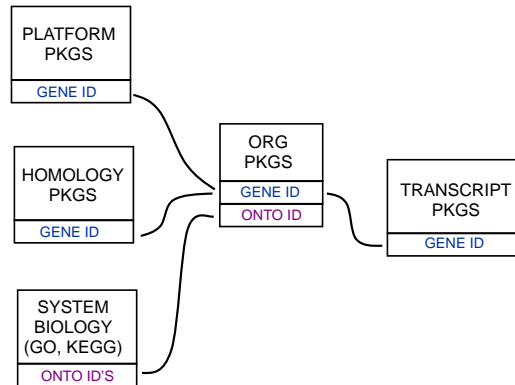


Figure 1: Relationships between Annotation packages

1 Getting started with OrganismDbi

Usage of a package like this has been deliberately kept very simple. The methods supported are the same ones that work for all the packages based on *AnnotationDb* objects. The methods that can be applied to these new packages are `columns`, `keys`, `keytypes` and `select`.

So to learn which kinds of data can be retrieved from a package like this we would simply load the package and then call the `columns` method.

```
> library(Homo.sapiens)
> columns(Homo.sapiens)

 [1] "GOID"           "TERM"           "ONTOLOGY"       "DEFINITION"     "ENTREZID"
 [6] "PFAM"          "IPI"            "PROSITE"        "ACCNUM"         "ALIAS"
[11] "CHR"           "CHRLOC"         "CHRLOCEND"      "ENZYME"         "MAP"
[16] "PATH"          "PMID"           "REFSEQ"         "SYMBOL"         "UNIGENE"
[21] "ENSEMBL"       "ENSEMBLPROT"   "ENSEMBLTRANS"  "GENENAME"       "UNIPROT"
[26] "GO"            "EVIDENCE"       "GOALL"          "EVIDENCEALL"    "ONTOLOGYALL"
[31] "OMIM"          "UCSCCKG"        "CDSID"          "CDSNAME"        "CDSCHROM"
[36] "CDSSTRAND"    "CDSSTART"       "CSEND"          "EXONID"         "EXONNAME"
[41] "EXONCHROM"    "EXONSTRAND"     "EXONSTART"     "EXONEND"        "GENEID"
[46] "TXID"          "EXONRANK"       "TXNAME"         "TXCHROM"        "TXSTRAND"
[51] "TXSTART"      "TXEND"
```

To learn which of those kinds of data can be used as keys to extract data, we use the `keytypes` method.

```
> keytypes(Homo.sapiens)

 [1] "GOID"           "TERM"           "ONTOLOGY"       "DEFINITION"     "ENTREZID"
 [6] "PFAM"          "IPI"            "PROSITE"        "ACCNUM"         "ALIAS"
[11] "CHR"           "CHRLOC"         "CHRLOCEND"      "ENZYME"         "MAP"
[16] "PATH"          "PMID"           "REFSEQ"         "SYMBOL"         "UNIGENE"
[21] "ENSEMBL"       "ENSEMBLPROT"   "ENSEMBLTRANS"  "GENENAME"       "UNIPROT"
[26] "GO"            "EVIDENCE"       "GOALL"          "EVIDENCEALL"    "ONTOLOGYALL"
[31] "OMIM"          "UCSCCKG"        "GENEID"         "TXID"           "TXNAME"
[36] "EXONID"        "EXONNAME"       "CDSID"          "CDSNAME"
```

To extract specific keys, we need to use the `keys` method, and also provide it a legitimate keytype:

```
> head(keys(Homo.sapiens, keytype="ENTREZID"))
```

```
[1] "1"          "10"         "100"        "1000"       "10000"      "100008586"
```

And to extract data, we can use the `select` method. The `select` method depends on the values from the previous three methods to specify what it will extract. Here is an example that will extract, UCSC transcript names, and gene symbols using Entrez Gene IDs as keys.

```
> k <- head(keys(Homo.sapiens, keytype="ENTREZID"),n=3)
> select(Homo.sapiens, keys=k, columns=c("TXNAME","SYMBOL"), keytype="ENTREZID")
```

	ENTREZID	SYMBOL	TXNAME
1	1	A1BG	uc002qsd.4
2	1	A1BG	uc002qsf.2
3	10	NAT2	uc003wyw.1
4	100	ADA	uc002xmj.3

In Addition to `select`, some of the more popular range based methods have also been updated to work with an *AnnotationDb* object. So for example you could extract transcript information like this:

```
> transcripts(Homo.sapiens, columns=c("TXNAME","SYMBOL"))
```

GRanges with 82960 ranges and 2 metadata columns:

	seqnames	ranges	strand	TXNAME
	<Rle>	<IRanges>	<Rle>	<CharacterList>
[1]	chr1	[11874, 14409]	+	uc001aaa.3
[2]	chr1	[11874, 14409]	+	uc010nxq.1
[3]	chr1	[11874, 14409]	+	uc010nxr.1
[4]	chr1	[69091, 70008]	+	uc001aal.1
[5]	chr1	[321084, 321115]	+	uc001aaq.2
...
[82956]	chrY	[27605645, 27605678]	-	uc004fwx.1
[82957]	chrY	[27606394, 27606421]	-	uc022cpc.1
[82958]	chrY	[27607404, 27607432]	-	uc004fwz.3
[82959]	chrY	[27635919, 27635954]	-	uc022cpd.1
[82960]	chrY	[59358329, 59360854]	-	uc011ncc.1

	SYMBOL
	<CharacterList>
[1]	DDX11L1
[2]	DDX11L1
[3]	DDX11L1

```

      [4]          OR4F5
      [5]           NA
      ...
[82956]           NA
[82957]           NA
[82958]           NA
[82959]           NA
[82960]           NA
----
seqlengths:
              chr1                chr2 ...      chrUn_g1000249
              249250621            243199373 ...      38502

```

And the *GRanges* object that would be returned would have the information that you specified in the `columns` argument. You could also have used the `exons` or `cds` methods in this way.

The `transcriptsBy`, `exonsBy` and `cdsBy` methods are also supported. For example:

```
> transcriptsBy(Homo.sapiens, by="gene", columns=c("TXNAME","SYMBOL"))
```

```
GRangesList of length 23459:
```

```
$1
```

```
GRanges with 2 ranges and 3 metadata columns:
```

	seqnames	ranges	strand	tx_name	TXNAME
	<Rle>	<IRanges>	<Rle>	<character>	<CharacterList>
[1]	chr19	[58858172, 58864865]	-	uc002qsd.4	uc002qsd.4
[2]	chr19	[58859832, 58874214]	-	uc002qsf.2	uc002qsf.2

```
SYMBOL
```

```
<CharacterList>
```

```
[1]          A1BG
[2]          A1BG
```

```
$10
```

```
GRanges with 1 range and 3 metadata columns:
```

	seqnames	ranges	strand	tx_name	TXNAME	SYMBOL
[1]	chr8	[18248755, 18258723]	+	uc003wyw.1	uc003wyw.1	NAT2

```
$100
```

```
GRanges with 1 range and 3 metadata columns:
```

```

      seqnames          ranges strand |      tx_name      TXNAME SYMBOL
[1]   chr20 [43248163, 43280376]    - | uc002xmj.3 uc002xmj.3   ADA
...
<23456 more elements>
---
seqlengths:
      chr1          chr2 ...      chrUn_gl000249
      249250621      243199373 ...      38502

```

2 Making your own OrganismDbi packages

So in the preceding section you can see that using an *OrganismDbi* package behaves very similarly to how you might use a *TranscriptDb* or an *OrgDb* package. The same methods are defined, and they behave similarly except that they now have access to much more data than before. But before you make your own *OrganismDbi* package you need to understand that there are few logical limitations for what can be included in this kind of package.

- The 1st limitation is that all the annotation resources in question must have implemented the four methods described in the preceding section (`columns`, `keys`, `keytypes` and `select`).
- The 2nd limitation is that you cannot have more than one example of each field that can be retrieved from each type of package that is included. So basically, all values returned by `columns` must be unique across ALL of the supporting packages.
- The 3rd limitation is that you cannot have more than one example of each object type represented. So you cannot have two *org* packages since that would introduce two *OrgDb* objects.
- And the 4th limitation is that you cannot have cycles in the graph. What this means is that there will be a graph that represents the relationships between the different object types in your package, and this graph must not present more than one pathway between any two nodes/objects. This limitation means that you can choose one foreign key relationship to connect any two packages in your graph.

With these limitations in mind, let's set up an example. Let's show how we could make *Homo.sapiens*, such that it allowed access to *org.Hs.eg.db*, *TxDb.Hsapiens.UCSC.hg19.knownGene* and *GO.db*.

The 1st thing that we need to do is set up a list that expresses the way that these different packages relate to each other. To do this, we make a `list` that contains short two element long character vectors. Each character vector represents one relationship between a pair of packages. The names of the vectors are the package names and the values are the foreign keys. Please note that the foreign key values in these vectors are the same strings that are returned by the `columns` method for the individual packages. Here is an example that shows how *GO.db*, *org.Hs.eg.db* and *TxDb.Hsapiens.UCSC.hg19.knownGene* all relate to each other.

```
> gd <- list(join1 = c(GO.db="GOID", org.Hs.eg.db="GO"),
+           join2 = c(org.Hs.eg.db="ENTREZID",
+           TxDb.Hsapiens.UCSC.hg19.knownGene="GENEID"))
```

So this `data.frame` indicates both which packages are connected to each other, and also what these connections are using for foreign keys.

Once this is finished, we just have to call the `makeOrganismPackage` function to finish the task.

```
> destination <- tempfile()
> dir.create(destination)
> makeOrganismPackage(pkgname = "Homo.sapiens",
+ graphData = gd,
+ organism = "Homo sapiens",
+ version = "1.0.0",
+ maintainer = "Package Maintainer<maintainer@somewhere.org>",
+ author = "Some Body",
+ destDir = destination,
+ license = "Artistic-2.0")
```

`makeOrganismPackage` will then generate a lightweight package that you can install. This package will not contain all the data that it refers to, but will instead depend on the packages that were referred to in the `data.frame`. Because the end result will be a package that treats all the data mapped together as a single source, the user is encouraged to take extra care to ensure that the different packages used are from the same build etc.