



9. Bioconductor: annotation databases

Thomas Lumley
Ken Rice

Universities of Washington and Auckland

Liège, September 2011

Outline

One goal of Bioconductor is to provide efficient access inside R to the genome databases that are vital to interpreting associations.

We will look at a few of these

- annotate
- biomaRt
- genomeGraphs

The reason to have an R interface to these databases is to be able to analyze annotation data for many SNPs or RNA transcripts.

Online or stored data

Annotation data can be downloaded in a single file or retrieved for each query from an online database.

Local storage is faster, but may require too much space (eg Ensembl) or become obsolete too quickly.

Local storage is ideal for fixed annotation data such as gene names for a microarray or SNP chip.

Types of database

Translations of names: Affy probe 32972_at is the gene **NADPH oxidase 1** with symbol **NOX1** and Ensembl gene id **ENSG00000007952**

Location: NOX1 is on Xq22.1, from 99984969 to 100015990, coded on the negative strand. There are 120 known polymorphisms (SNPs or indels) in this range.

Homology: The mouse version of NOX1 is also on the X chromosome, starting at 130621066.

Structure and function: **NOX1** is a membrane protein (location), involved in voltage-gated ion channel activity (molecular function), and involved in signal transduction (biological process).

Annotate

Bioconductor distributes annotation packages for a wide range of gene expression microarrays. The `annotate` package is one way to use this annotation information.

```
> library("annotate")
> library("hgu95av2.db")
> library("GO.db")
```

loads the `annotate` package and the databases for the Gene Ontology and one of the Affymetrix human microarray chips.

Lookups

The databases are queried with `get()` or `mget()` for multiple queries

```
> mget(c("738_at", "40840_at", "32972_at"), hgu95av2GENENAME)
$'738_at'
[1] "5'-nucleotidase, cytosolic II"

$'40840_at'
[1] "peptidylprolyl isomerase F (cyclophilin F)"

$'32972_at'
[1] "NADPH oxidase 1"

> go<-get("738_at", hgu95av2GO)
> names(go)
[1] "GO:0009117" "GO:0005829" "GO:0005737" "GO:0000166" "GO:0000287"
[6] "GO:0008253" "GO:0008253" "GO:0016787"
```

Lookups

```
> get("GO:0009117", GOTERM)
```

OID: GO:0009117

Term: nucleotide metabolic process

Ontology: BP

Definition: The chemical reactions and pathways involving a nucleotide, a nucleoside that is esterified with (ortho)phosphate or an oligophosphate at any hydroxyl group on the glucose moiety; may be mono-, di- or triphosphate; this definition includes cyclic nucleotides (nucleoside cyclic phosphates).

Synonym: nucleotide metabolism

What lookups are there?

```
> library(help="hgu95av2.db")
```

```
hgu95av2ALIAS2PROBE  Map between Common Gene Symbol Identifiers and  
Manufacturer Identifiers
```

```
> get("NOX1",hgu95av2ALIAS2PROBE)  
[1] "32972_at"    "32973_s_at"
```

You can also reverse a lookup table with `revmap()`

```
> get("NOX1",revmap(hgu95av2SYMBOL))  
[1] "32972_at"    "32973_s_at"  
> get("X",revmap(hgu95av2CHR))  
[1] "1016_s_at"   "107_at"      "1100_at"     "112_g_at"    "1155_at"  
.... and lots more
```

BioMart

BioMart (www.biomart.org) is a query-oriented data management system developed jointly by the European Bioinformatics Institute (EBI) and Cold Spring Harbor Laboratory (CSHL).

`biomaRt` is an R interface to BioMart systems, in particular to Ensembl (www.ensembl.org). Ensembl is a joint project between EMBL - European Bioinformatics Institute (EBI) and the Wellcome Trust Sanger Institute (WTSI) to develop a software system which produces and maintains automatic annotation on selected eukaryotic genomes.

BioMart

We begin by choosing which BioMart to use

```
> library(biomaRt)
Loading required package: RCurl
> listMarts()
      biomart
1          ensembl
2              snp
3      functional_genomics
4              vega
5      bacteria_mart_10
6          fungi_mart_10
7      fungi_variations_10
8      metazoa_mart_10
9      metazoa_variations_10
...
60     ENSEMBL_MART_PLANT
61     ENSEMBL_MART_PLANT_SNP
62     GRAMENE_MARKER_30
63     GRAMENE_MAP_30
64     QTL_MART
65     salmosalar2_mart
66     trucha_mart
> ens <- useMart("ensembl")
```

biomart				
1	ensembl	ENSEMBL GENES	63	(SANGER U)
2	snp	ENSEMBL VARIATION	63	(SANGER U)
3	functional_genomics	ENSEMBL REGULATION	63	(SANGER U)
4	vega	VEGA	43	(SANGER U)
5	bacteria_mart_10	ENSEMBL BACTERIA	10	(EBI U)
6	fungi_mart_10	ENSEMBL FUNGI	10	(EBI U)
7	fungi_variations_10	ENSEMBL FUNGI VARIATION	10	(EBI U)
8	metazoa_mart_10	ENSEMBL METAZOA	10	(EBI U)
9	metazoa_variations_10	ENSEMBL METAZOA VARIATION	10	(EBI U)
...				
60	ENSEMBL_MART_PLANT	GRAMENE	30	ENSEMBL GENES (CSHL/CORNELL U)
61	ENSEMBL_MART_PLANT_SNP	GRAMENE	30	VARIATION (CSHL/CORNELL U)
62	GRAMENE_MARKER_30	GRAMENE	30	MARKERS (CSHL/CORNELL U)
63	GRAMENE_MAP_30	GRAMENE	30	MAPPINGS (CSHL/CORNELL U)
64	QTL_MART	GRAMENE	32	QTL DB (CSHL/CORNELL U)
65	salmosalar2_mart	UNIGENE SALMO SALAR DATABASE		(CMM CHIL)
66	trucha_mart	UNIGENE ONCORHYNCHUS MYKISS DATABASE		(CMM CHIL)

BioMart

We then choose a database to use

```
> listDatasets(ens)
      dataset                               description
1   oanatinus_gene_ensembl    Ornithorhynchus anatinus genes (OANA5)
2   tguttata_gene_ensembl    Taeniopygia guttata genes (taeGut3.2.4)
3   cporcellus_gene_ensembl   Cavia porcellus genes (cavPor3)
4   gaculeatus_gene_ensembl  Gasterosteus aculeatus genes (BROADS1)
5   lafricana_gene_ensembl   Loxodonta africana genes (loxAfr3)
6   mlucifugus_gene_ensembl  Myotis lucifugus genes (myoLuc2)
7   hsapiens_gene_ensembl    Homo sapiens genes (GRCh37.p3)
...
26  pvampyrus_gene_ensembl   Pteropus vampyrus genes (pteVam1)
27  mdomestica_gene_ensembl  Monodelphis domestica genes (monDom5)
28  vpacos_gene_ensembl     Vicugna pacos genes (vicPac1)
...
52  meugenii_gene_ensembl   Macropus eugenii genes (Meug_1.0)
53  cfamiliaris_gene_ensembl Canis familiaris genes (CanFam_2.0)

> ens <- useDataset("hsapiens_gene_ensembl",mart=ens)
```

BioMart

The `getGene` function queries the database for gene information. It accepts many forms of gene identifier, eg Entrez, HUGO, Affy transcript

```
> getGene(id=1440, type="entrezgene", mart=ens)
  entrezgene hgnc_symbol
1        1440          CSF3

1 Granulocyte colony-stimulating factor Precursor (G-CSF)(Pluripotin)
(Filgrastim)(Lenograstim) [Source:UniProtKB/Swiss-Prot;Acc:P09919]
  chromosome_name band strand start_position end_position ensembl_gene_id
1            17 q21.1       1      35425140      35427592 ENSG00000108342

> getGene(id=c("AGT","AGTR1"), type="hgnc_symbol", mart=ens)
  hgnc_symbol hgnc_symbol
1        AGTR1          AGTR1
2         AGT            AGT
1

2 Angiotensinogen Precursor (Serpine A8) [Contains Angiotensin-1(Angiotensin I)
(Ang I);Angiotensin-2(Angiotensin II)(Ang II);Angiotensin-3(Angiotensin III)
(Ang III)(Des-Asp[1]-angiotensin II)] [Source:UniProtKB/Swiss-Prot;Acc:P01019]
  chromosome_name band strand start_position end_position ensembl_gene_id
1            3   q24       1     149898348     149943478 ENSG00000144891
2           1  q42.2      -1     228904897     228916564 ENSG00000135744
```

BioMart

getBM is more general than getGene. It specifies a list of **filters** for selecting genes or SNPs and **attributes** to return from the database.

```
> affyids <- c("202763_at", "209310_s_at", "207500_at")
> getBM(attributes = c("affy_hg_u133_plus_2", "hgnc_symbol", "chromosome_name",
  "start_position", "end_position", "band"), filters = "affy_hg_u133_plus_2",
  values = affyids, mart = ens)
      affy_hg_u133      hgnc      chromosome_name start_position end_position band
1      202763_at      CASP3          4        185785844    185807623 q35.1
2      207500_at      CASP5         11        104370180    104384957 q22.3
3     209310_s_at      CASP4         11        104318804    104344535 q22.3
```

listAttributes(ens) and listFilters(ens) list the available attributes and filters (hundreds)

BioMart

```
> getBM(mart=ens,c("band","hgnc_symbol"),
         filters=c("band_start","band_end","chromosome_name"),
         values=list("p21.33","p21.33",6))
      band hgnc_symbol
1   p21.33
2   p21.33    SNORD117
3   p21.33    SNORA38
4   p21.33    SNORD48
5   p21.33    SNORD52
6   p21.33    MIR877
7   p21.33    MIR1236
8   p21.33    GTF2H4
9   p21.33    VARS2
10  p21.33    SFTA2
11  p21.33    DPCR1
12  p21.33    MUC21
...
121 p21.33    HSPA1A
122 p21.33    TNXB
123 p21.33    STK19
124 p21.33    C4A
125 p21.33    C4B
```

Homology

`getLDS()` combines two data marts, for example to homologous genes in other species. We can look up the mouse equivalents of a particular Affy transcript, or of the NOX1 gene.

```
> human = useMart("ensembl", dataset = "hsapiens_gene_ensembl")
> mouse = useMart("ensembl", dataset = "mmusculus_gene_ensembl")
> getLDS(attributes = c("hgnc_symbol", "chromosome_name", "start_position"),
  filters = "hgnc_symbol", values = "NOX1", mart = human,
  attributesL = c("chromosome_name", "start_position"),
  martL = mouse)
V1 V2      V3 V4      V5
1 NOX1  X 99984969  X 130621066
```

Homology

The `getSequence` function looks up DNA or protein sequences by chromosome position or gene identifiers

```
> agt<-getSequence(id="AGT",type="hgnc_symbol", seqType="peptide",mart=ens)
> agt
```

```
1 MRKRAPQSEMAPAGVSLRATILCLLAWAGLAAGDRVYIHPFHLVIHNESTCEQLAKANAGKPKDPTFIPAPIQAKTS
PVDEKALQDQLVLVAAKLDTEDKLRAAMVGMLANFLGFRIYGMHSELWGVVHGATVLSPTAVFGTLASLYLGALDHTAD
RLQAILGVPWKDKNCTSRLDAHKVLSALQAVQGLLVAQGRADSQAQLLLSTVVGVFTAPGLHLKQPFVQGLALYTPVVL
PRSLDFTELDVAAEKIDRFMQAVTGWTGCSLMGASVDSTLAFNTYVHFQGKMKGSLLAEPQEFWVDNSTSVSPMLS
GMGTFQHWSDIQDNFSVTQVPFTESACLLLIQPHYASDLDKVEGLTFQQNSLNWMKKLSPRTIHLTMPQLVLQGSYDLQ
DLLAQAEELPAILHTELNLQKLSNDRIRVGEVLNSIFFELEADEREPTESTQQLNKPEVLEVTLNRPFLFAVYDQSATAL
HFLGRVANPLSTA*
```

Example: finding chromosomes

We had a set 1524 SNPs, of which 409 did not have their chromosome listed.

I needed to know which SNPs were on the X chromosome, to estimate sex from DNA intensity and heterozygous X-chromosome loci, for QC.

```
> head(unknown)
[1] "UGT1A3-001449-0_B_R_1538822" "LIPC-002761-0_B_R_1538453"
[3] "CETP-001265-0_B_R_1538254"     "F8-165293-0_T_F_1538626"
[5] "CPB2-051208-0_B_F_1539402"    "VDRDIL-1355-0_T_F_1539404"
```

A hand-search would be easy but tedious, so we want an automated approach

Example: finding chromosomes

First extract the gene names

```
genes <- sapply(unknown, function(snp) strsplit(snp,"-")[[1]][1])
ugenes <- unique(genes)
```

Now call out to Ensembl

```
getBM(attributes="chromosome_name", filters="hgnc_symbol", values=ugenes,
      mart=ensembl)
```

works for all except VRDIL, which isn't recognized.

Finding SNPs

Human SNPs are in a separate database from gene information. We can look up known SNPs and other polymorphisms for the NOX1 gene

```
> snpmart = useMart("snp", dataset = "hsapiens_snp")
Checking attributes ... ok
Checking filters ... ok
> getBM(c("refsnp_id", "allele", "chrom_start", "chrom_strand"),
  filters = c("chr_name", "chrom_start", "chrom_end"),
  values = list("X", 99984969, 100015990), mart = snpmart)
      refsnp_id          allele chrom_start chrom_strand
1      rs35477500        AAG/-    99985012            1
2      rs41310727        G/T     99985014            1
3      rs61639376        GAA/-    99985014            1
4      rs34181451        G/A     99985489            1
5      rs34009592        G/A     99985531            1
```

GenomeGraphs

This package makes pretty pictures from the annotation data.

For example, a pictures showing the standard and alternative splices for the NOX1 gene and the location of the gene on the X chromosome

```
> library(GenomeGraphs)
> gene <- makeGene(id = "NOX1", type = "hgnc_symbol",
  biomart = ensembl)
> transcript<-makeTranscript(id="NOX1",type="hgnc_symbol",
  biomart=ensembl)
> ideogram <- makeIdeogram(chromosome ="X")
> gdPlot(list(ideogram,gene,transcript))
```

GenomeGraphs

