# Basic GO Usage

#### R. Gentleman

November 8, 2008

### Introduction

In this vignette we describe some of the basic characteristics of the data available from the Gene Ontology (GO), (The Gene Ontology Consortium, 2000) and how these data have been incorporated into Bioconductor. We assume that readers are familiar with the basic DAG structure of GO and with the mappings of genes to GO terms that are provide by GOA (Camon et al., 2004). We consider these basic structures and properties quite briefly.

GO, itself, is a structured terminology. The ontology describes genes and gene products and is divided into three separate ontologies. One for cellular component (CC), one for molecular function (MF) and one for biological process (BP). We maintain those same distinctions were appropriate. The relationship between terms is a parent-child one, where the parents of any term are less specific than the child. The mapping in either direction can be one to many (so a child may have many parents and a parent may have many children). There is a single root node for all ontologies as well as separate root nodes for each of the three ontologies named above. These terms are structured as a directed acyclic graph (or a DAG).

GO itself is only the collection of terms; the descriptions of genes, gene products, what they do, where they do it and so on. But there is no direct association of genes to terms. The assignment of genes to terms is carried out by others, in particular the GOA project (Camon et al., 2004). It is this assignment that makes GO useful for data analysis and hence it is the combined relationship between the structure of the terms and the assignment of genes to terms that is the concern of the GO.db package.

The basis for child-parent relationships in GO can be either an *is-a* relationship, where the child term is a more specific version of the parent. Or, it can be a *has-a*, or *part-of* relationship where the child is a part of the parent. For example a telomere is a part-of a chromosome.

Genes are assigned to terms on the basis of their LocusLink ID. For this reason we make most of our mappings and functions work for LocusLink identifiers. Users of specific chips, or data with other gene identifiers should first map their identifiers to LocusLink before using *GOstats*.

A gene is mapped only to the most specific terms that are applicable to it (in each ontology). Then, all less specific terms are also applicable and they are easily obtained by traversing the set of parent relationships down to the root node. In practice many of these mappings are precomputed and easily obtained from the different hash tables provided by the GO.db package.

Mapping of a gene to a term can be based on many different things. GO and GOA provide an extensive set of evidence codes, some of which are given in Table 1, but readers are referred to the GO web site and the documentation for the GO.db package for a more comprehensive listing. Clearly for some investigations one will want to exclude genes that were mapped according to some of the evidence codes.

IMP	inferred from mutant phenotype
IGI	inferred from genetic interaction
IPI	inferred from physical interaction
ISS	inferred from sequence similarity
IDA	inferred from direct assay
IEP	inferred from expression pattern
IEA	inferred from electronic annotation
TAS	traceable author statement
NAS	non-traceable author statement
ND	no biological data available
IC	inferred by curator

Table 1: GO Evidence Codes

In some sense TAS is probably the most reliable of the mappings. IEA is a weak association and is based on electronic information, no human curator has examined or confirmed this association. As we shall see later, IEA is also the most common evidence code.

The sets of mappings of interest are roughly divided into three parts. First there is the basic description of the terms etc., these are provided in the GOTERMS hash table. Each element of this hash table is named using its GO identifier (these are all of the form GO: followed by seven digits). Each element is an instance of the GOTerms class. A complete description of this class can be obtained from the appropriate manual page (use class?GOTerms). From these data we can find the text string describing the term, which ontology it is in as well as some other basic information.

There are also a set of hash tables that contain the information about parents and children. They are provided as hash tables (the XX in the names below should be substituted for one of BP, MF, or CC.

- GOXXPARENTS: the parents of the term
- GOXXANCESTOR: the parents, and all their parents and so on.

- GOXXCHILDREN: the children of the term
- GOXXOFFSPRING: the children, their children and so on out to the leaves of the GO graph.

For the GOXXPARENTS mappings (only) information about the nature of the relationship is included.

#### > GDTERM\$"GO:0003700"

GOID: GO:0003700

Term: transcription factor activity

Ontology: MF

Definition: The function of binding to a specific DNA sequence in order to modulate transcription. The transcription factor may or may not also interact selectively with a protein or macromolecular complex.

Synonym: GO:0000130 Secondary: GO:0000130

> GOMFPARENTS\$"GO:0003700"

isa isa "GD:0003677" "GD:0030528"

> GOMFCHILDREN\$"GO:0003700"

isa "GD:0003705"

Here we see that the term GO:0003700 has two parents, that the relationships are is-a and that it has one child. One can then follow this chains of relationships or use the ANCESTOR and OFFSPRING hash tables to get more information.

The mappings of genes to GO terms is not contained in the GO package. Rather these mappings are held in each of the chip and organism specific data packages, such as hgu95av2GO and org.Hs.egGO are contained within packages hgu95av2.db and org.Hs.eg.db respectively. These mappings are from a Entrez Gene ID to the most specific applicable GO terms. Each such entry is a list of lists where the innermost list has these names:

- GOID: the GO identifier
- Evidence: the evidence code for the assignment
- Ontology: the ontology the GO identifier belongs to (one of BP, MF, or CC).

Some genes are mapped to a GO identifier based on two or more evidence codes. Currently these appear as separate entries. So you may want to remove duplicate entries if you are not interested in evidence codes. However, as more sophisticated use is made of these data it will be important to be able to separate out mappings according to specific evidence codes.

In this next example we consider the gene with Entrez Gene ID 4121, this corresponds to Affymetrix ID 39613\_at.

```
> 111 = hgu95av2G0[["39613_at"]]
> length(ll1)
[1] 14
> sapply(ll1, function(x) x$Ontology)
GD:0008152 GD:0000139 GD:0016020 GD:0016021 GD:0005783 GD:0005624 GD:0005793
      "BP"
                  "CC"
                             "CC"
                                         "CC"
                                                    "CC"
                                                                "CC"
                                                                            "CC"
GD:0005794 GD:0005509 GD:0005509 GD:0004571 GD:0004571 GD:0015923 GD:0016798
                                         "MF"
                                                    "MF"
                                                                "MF"
      "CC"
                  "MF"
                             "MF"
                                                                            "MF"
```

We see that there are 14 different mappings. We can get only those mappings for the BP ontology by using getOntology. We can get the evidence codes using getEvidence and we can drop those codes we do not wish to use by using dropECode.

```
> getOntology(111, "BP")
[1] "GO:0008152"
> getEvidence(111)
GD:0008152 GD:0000139 GD:0016020 GD:0016021 GD:0005783 GD:0005624 GD:0005793
     "IEA"
                "IEA"
                            "IEA"
                                       "IEA"
                                                   "TAS"
                                                               "TAS"
                                                                          "IDA"
GD:0005794 GD:0005509 GD:0005509 GD:0004571 GD:0004571 GD:0015923 GD:0016798
     "IEA"
                "IEA"
                            "TAS"
                                       "IEA"
                                                   "TAS"
                                                               "TAS"
                                                                          "IEA"
> zz = dropECode(111)
> getEvidence(zz)
GD:0005783 GD:0005624 GD:0005793 GD:0005509 GD:0004571 GD:0015923
     "TAS"
                "TAS"
                            "TDA"
                                       "TAS"
                                                   "TAS"
                                                               "TAS"
```

### A Basic Description of GO

We now characterize GO and some of its properties. First we list some of the specific GO IDs that might be of interest (please feel free to propose even more).

- GD:0003673 is the GO root.
- GO:0003674 is the MF root.
- GO:0005575 is the CC root.
- GO:0008150 is the BP root.
- GO:0000004 is biological process unknown
- GO: 0005554 is molecular function unknown
- GO:0008372 is cellular component unknown

We can find out how many terms are in each of the different ontologies by:

Or we can ask about the number of is-a and part of relationships in each of the three different ontologies.

```
> BPisa = eapply(GOBPPARENTS, function(x) names(x))
> table(unlist(BPisa))
                 isa negatively_regulates
                                                        part_of
                                       982
                                                           3171
positively_regulates
                                regulates
                 962
                                      1131
> MFisa = eapply(GOMFPARENTS, function(x) names(x))
> table(unlist(MFisa))
    isa part_of
> CCisa = eapply(GOCCPARENTS, function(x) names(x))
> table(unlist(CCisa))
    isa part_of
   3131
            881
```

## Working with GO

Finding terms that have specific character strings in them is easily accomplished using grep. In the next example we first convert the data from GOTERM into a character vector to make it easier to do multiple searches.

```
> goterms = unlist(eapply(GOTERM, function(x) x@Term))
> whmf = grep("molecular_function", goterms)
```

So we see that there are 1 terms with the string "molecular\_function" in them in the ontology. They can be accessed by subsetting the goterms object.

> goterms[whmf]

GO:0003674
"molecular function"

#### Working with chip specific meta-data

In some cases users will want to restrict their attention to the set of terms etc that map to genes that were assayed in the experiments that they are working with. To do this you should first get the appropriate chip specific meta-data file. Here we demonstrate some of the examples on the Affymetrix HGu95av2 chips and so use the package hgu95av2.db. Each of these packages has a data environment whose name is the basename of the package with a GO suffix, so in this case hgu95av2GO. Note that if there are many manufacturer ids that map to the same Entrez Gene identifier then these will be duplicate entries (with different keys).

We can get all the MF terms for our Affymetrix data.

```
> affyGO = eapply(hgu95av2GO, getOntology)
> table(sapply(affyGO, length))
   0
        1
                   3
                                   6
                                         7
                                                    9
                                                        10
                                                              11
                                                                   12
                                                                         13
                                                                              14
                                                                                    15
1863 1884 2287 2159 1592 1073 705
                                      486
                                            258
                                                  113
                                                        86
                                                              55
                                                                   30
                                                                         11
                                                                               9
                                                                                     4
  16
       17
   6
        4
```

How many of these probes have multiple GO terms associated with them? What do we do if we want to compare two genes that have multiple GO terms associated with them?

What about evidence codes? To find these we apply a similar function to the affyGO terms.

```
> affyEv = eapply(hgu95av2GO, getEvidence)
> table(unlist(affyEv, use.names = FALSE))
```

EXP IC IDA **IEA IEP** IGI IMP IPI ISS NAS ND RCA TAS 3010 330 10061 60358 180 75 1894 4243 3436 7154 1516 13 20447

- > test1 = eapply(hgu95av2GO, dropECode, c("IEA", "NR"))
- > table(unlist(sapply(test1, getEvidence), use.names = FALSE))

**EXP** RCA IC IDA **IEP** IGI ISS ND TAS IMP IPI NAS 3010 330 10061 180 75 1894 4243 3436 7154 1516 13 20447

These functions make is somewhat straightforward to select subsets of the GO terms that are specific to different evidence codes.

# References

E. Camon, M. Magrane, D. Barrell, V. Lee, E. Dimmer, D. Binns J. Maslen, N. Harte, R. Lopez, and R. Apweiler. The Gene Ontology annotation (GOA) database: sharing knowledge in Uniprot with Gene Ontology. *Nucleic Acids Research*, 32:D262–D266, 2004.

The Gene Ontology Consortium. Gene Ontology: tool for the unification of biology. *Nature Genetics*, 25:25–29, 2000.