

rols: an R interface to the Ontology Lookup Service

Laurent Gatto

lg390@cam.ac.uk

Cambridge Center for Proteomics
University of Cambridge

October 14, 2013

Abstract

The rols package provides a common interface to 87 different ontologies through EBI's Ontology Lookup Service. This vignette provides a brief overview of the available interface and functionality as well as a short use case.

Keywords: infrastructure, bioinformatics, ontology.

1 Introduction

The Ontology Lookup Service¹ (OLS) (Côté et al., 2006, 2008) is a spin-off of the PRoteomics IDentifications database (PRIDE) service, located at the EBI. OLS provides a unified interface to 87 ontologies (see below).

rols makes use of the SOAP service at the EBI to post XML requests. The SOAP XML responses are then parsed and returned in an R friendly data structure. This is achieved using Duncan Temple Lang's SSOAP package (Temple Lang, 2011).

2 Brief rols overview

2.1 Ontologies

There are 87 ontologies available in the OLS, listed in the table 1 below. Their name is to be used to define which ontology to query.

¹<http://www.ebi.ac.uk/ontology-lookup/>

Name	Description
AAO	Amphibian Gross Anatomy
APO	Yeast phenotypes
ATO	Amphibian Taxonomy
BFO	Basic Formal Ontology
BS	Biosapiens Annotations
BSPO	Spatial Reference Ontology
BTO	BRENDA tissue / enzyme source
CARO	Common Anatomy Reference Ontology
CCO	Cell Cycle Ontology
CHEBI	Chemical Entities of Biological Interest
CL	Cell Type
DDANAT	Dictyostelium discoideum Anatomy
DOID	Human Disease
DPO	Drosophila Phenotype Ontology
ECO	Evidence Codes
EDAM	EMBRACE Data and Methods
EFO	Experimental Factor Ontology
EHDA	Human Developmental Anatomy, timed version
EMAP	Mouse Gross Anatomy and Development, timed
EMAPA	Mouse Gross Anatomy and Development, abstract
ENA	European Nucleotide Archive Submission Ontology
ENVO	Environmental Ontology
EO	Plant Environmental Conditions
EV	eVOC (Expressed Sequence Annotation for Humans)
FAO	Fungal Gross Anatomy
FBbi	Biological Imaging Methods
FBbt	Drosophila Gross Anatomy
FBcv	Flybase Controlled Vocabulary
FBdv	Drosophila Development
FBsp	Fly taxonomy
FIX	Physico-Chemical Methods and Properties
FMA	Foundational Model of Anatomy Ontology
FYPO	Fission Yeast Phenotype Ontology
GAZ	Gezetteer ontology
GO	Gene Ontology
GRO	Cereal Plant Development
HAO	Hymenoptera Anatomy Ontology
HOM	Homology Ontology
HP	Human phenotype ontology
IDO	Infectious Disease Ontology
IEV	Event (INOH)
IMR	Molecule Role (INOH)
LSM	Leukocyte Surface Markers
MA	Mouse Adult Gross Anatomy
MAT	Minimal Information About Anatomy ontology
MFO	Medaka Fish Anatomy and Development

MI	Molecular Interaction (PSI MI 2.5)
MIAA	Minimal Information About Anatomy ontology
MIRO	Mosquito Insecticide Resistance
MOD	Protein Modifications (PSI-MOD)
MP	Mammalian Phenotype
MPATH	Mouse Pathology
MS	PSI Mass Spectrometry Ontology
NEWT	NEWT UniProt Taxonomy Database
OBO_REL	OBO Relationship Types
PAR	Protein Affinity Reagents
PATO	Phenotypic qualities (properties)
PM	Phenotypic manifestation (genetic context)
PO	Plant Ontology (Structure, Growth and Developmental Stage)
PRIDE	PRIDE Controlled Vocabulary
PRO	Protein Ontology
PW	Pathway Ontology
REX	Physico-Chemical Process
RO	Multiple Alignment
SBO	Systems Biology Ontology
SEP	Separation Methods
SO	Sequence Types and Features
SPD	Spider Comparative Biology Ontology
SYMP	Symptom Ontology
TADS	Tick Gross Anatomy
TAIR	Arabidopsis Development
TAO	Teleost Anatomy and Development Ontology
TAXRANK	Taxonomic rank vocabulary
TGMA	Mosquito Gross Anatomy
TO	Cereal Plant Trait
TRANS	Pathogen transmission
TTO	Teleost taxonomy
UBERON	Uber anatomy ontology
UO	Unit Ontology
VariO	Variation Ontology
WBPhenotype	C. elegans phenotype
WBbt	C. elegans gross anatomy
WBls	C. elegans Development
XAO	Xenopus anatomy and development
ZDB	Zebrafish Anatomy and Development
ZEA	Maize Gross Anatomy
ZFA	Zebrafish Anatomy and Development

Table 1: Available ontologies in the OLS and rols package.

2.2 Interface

Table 2 summarised the common interface available for the 87 ontologies of table 1. More information is provided in the respective manual pages.

Function	Description
olsVersion	Returns the OLS version
ontologies	Returns all available ontologies
ontologyNames	Returns all ontologyNames
ontologyLoadDate	Returns the ontology load date
isIdObsolete	Is the ontology id obsolete
term	Returns the term of a given identifier
termMetadata	Returns an identifier's metadata
termXrefs	Returns the identifier's ontology cross references
rootId	Returns the root identifiers of an ontology
allIds	Returns all identifiers and terms of an ontology
olsQuery	Returns matching identifiers
parents	Returns the parent(s) of a term.
childrenRelations	Returns the children relation type(s).

Table 2: Functions available to query the ontologies.

2.3 Use case

A researcher might be interested in the trans-Golgi network and interested in knowing in which ontologies his favourite organelle is referenced. This can be done by querying all ontologies with a relevant pattern. The code below describes how to achieve this.

```
> library("rols")
> alltgn <- olsQuery("trans-golgi network")
```

As shown below, 4 different ontologies have matched the query string.

```
> alltgn
                                CCD:C0000738
"CC0:trans-Golgi network transport vesicle membrane"
                                GO:0005802
                                "GO:trans-Golgi network"
                                FMA:61756
                                "FMA:Trans Golgi network"
                                CCD:C0001468
                                "CC0:trans-Golgi network membrane"
                                CCD:C0000984
                                "CC0:trans-Golgi network transport vesicle"
                                CCD:C0000975
                                "CC0:clathrin coat of trans-Golgi network vesicle"
                                CCD:C0000381
                                "CC0:trans-Golgi network"
                                PR:000016299
"PRO:trans-Golgi network integral membrane protein 2"
                                PR:000016925
"PRO:trans-Golgi network integral membrane protein 1"
```

```

GO:0032588
  "GO:trans-Golgi network membrane"
GO:0044795
"GO:trans-Golgi network to recycling endosome transport"
GO:0012510
  "GO:trans-Golgi network transport vesicle membrane"
GO:0030130
  "GO:clathrin coat of trans-Golgi network vesicle"
GO:0030140
  "GO:trans-Golgi network transport vesicle"

> allonts <- sapply(strsplit(names(alltgns), ":"), "[", 1)
> onto.tab <- table(allonts)
> onto.tab

allonts
CCO FMA  GO  PR
  5  1   6   2

```

The description of the 4 ontologies of interest can then be used to subset the ontology description:

```

> ontologies()[names(onto.tab), ]

      Name                Description
CCO  CCO                Cell Cycle Ontology
FMA  FMA Foundational Model of Anatomy Ontology
GO   GO                 Gene Ontology
NA   <NA>                <NA>

```

To restrict the search to a specific ontology of interest, one can specify the ontology name as a parameter to `olsQuery`.

```

> gotgns <- olsQuery("trans-golgi network", "GO")
> gotgns

GO:0005802
  "trans-Golgi network"
GO:0030130
  "clathrin coat of trans-Golgi network vesicle"
GO:0032588
  "trans-Golgi network membrane"
GO:0030140
  "trans-Golgi network transport vesicle"
GO:0044795
"trans-Golgi network to recycling endosome transport"
GO:0012510
  "trans-Golgi network transport vesicle membrane"

```

Details about relevant terms can be retrieved with the `term` and `termMetadata` functions. This functionality provides on-line access to the same data that is available in the GO.db, and can be extended to any of the 87 available ontologies.

```
> term("GO:0005802", "GO")

[1] "trans-Golgi network"

> mtd <- termMetadata("GO:0005802", "GO")
> names(mtd)

[1] "related_synonym_3" "related_synonym_2" "related_synonym_1"
[4] "definition"       "related_synonym_4" "exact_synonym_1"
[7] "comment"          "exact_synonym_2"

> mtd["comment"]

                                                                 comm
"The TGN is not considered part of the Golgi apparatus but is a separate organelle"

> mtd["definition"]

"The network of interconnected tubular and cisternal structures located at the side of the Golgi apparatus distal to the endoplasmic reticulum, from which secretory vesicles emerge. The trans-Golgi network is important in the later stages of protein secretion where it is thought to play a key role in the sorting and targeting of secreted proteins to the correct destination."

> ## same as from GO.db
> GOTERM[["GO:0005802"]]

GOID: GO:0005802
Term: trans-Golgi network
Ontology: CC
Definition: The network of interconnected tubular and
           cisternal structures located at the side of the Golgi
           apparatus distal to the endoplasmic reticulum, from
           which secretory vesicles emerge. The trans-Golgi
           network is important in the later stages of protein
           secretion where it is thought to play a key role in
           the sorting and targeting of secreted proteins to the
           correct destination.
Synonym: Golgi trans face
Synonym: Golgi trans-face
Synonym: late Golgi
Synonym: maturing face
Synonym: TGN
Synonym: trans Golgi network
```

2.4 On-line vs. off-line data

It is possible to observe different results with rols and GO.db (Carlson et al.), as a result of the different ways they access the data. rols or biomaRt (Durinck et al., 2005) perform direct online queries, while GO.db and other annotation packages use database snapshot that are updated every release.

Both approaches have advantages. While online queries allow to obtain the latest up-to-date information, such approaches rely on network availability and quality. If reproducibility is a major issue, the version of the database to be queried can easily be controlled with off-line approaches. In the case of rols, although the load date of a specific ontology can be queried with `olsVersion`, it is not possible to query a specific version of an ontology.

Session information

- R version 3.0.2 (2013-09-25), x86_64-unknown-linux-gnu
- Locale: LC_CTYPE=en_US.UTF-8, LC_NUMERIC=C, LC_TIME=en_US.UTF-8, LC_COLLATE=C, LC_MONETARY=en_US.UTF-8, LC_MESSAGES=en_US.UTF-8, LC_PAPER=en_US.UTF-8, LC_NAME=C, LC_ADDRESS=C, LC_TELEPHONE=C, LC_MEASUREMENT=en_US.UTF-8, LC_IDENTIFICATION=C
- Base packages: base, datasets, grDevices, graphics, methods, parallel, stats, utils
- Other packages: AnnotationDbi 1.24.0, Biobase 2.22.0, BiocGenerics 0.8.0, DBI 0.2-7, GO.db 2.10.1, RSQLite 0.11.4, codetools 0.2-8, knitr 1.5, rols 1.4.0, xtable 1.7-1
- Loaded via a namespace (and not attached): IRanges 1.20.0, RCurl 1.95-4.1, SSOAP 0.8-0, XML 3.98-1.1, XMLSchema 0.7-2, evaluate 0.5.1, formatR 0.9, highr 0.2.1, stats4 3.0.2, stringr 0.6.2, tools 3.0.2

References

Marc Carlson, Seth Falcon, Herve Pages, and Nianhua Li. *GO.db: A set of annotation maps describing the entire Gene Ontology*. R package version 2.7.1.

Richard G Côté, Philip Jones, Rolf Apweiler, and Henning Hermjakob. The ontology lookup service, a lightweight cross-platform tool for controlled vocabulary queries. *BMC Bioinformatics*, 7:97, 2006. doi: 10.1186/1471-2105-7-97.

Richard G Côté, Philip Jones, Lennart Martens, Rolf Apweiler, and Henning Hermjakob. The ontology lookup service: more data and better tools for controlled vocabulary queries. *Nucleic Acids Res.*, 36(Web Server issue):372–376, 2008. doi: 10.1093/nar/gkn252.

S Durinck, Y Moreau, A Kasprzyk, S Davis, B De Moor, A Brazma, and W Huber. Biomart and bioconductor: a powerful link between biological databases and microarray data analysis. *Bioinformatics*, 21(16):3439–40, Aug 2005. doi: 10.1093/bioinformatics/bti525.

Duncan Temple Lang. *SSOAP: Client-side SOAP access for S*, 2011. URL <http://www.omegahat.org/SSOAP>, <http://www.omegahat.org>, <http://www.omegahat.org/bugs>. R package version 0.8-1.