

Advances in Big Data Bio Analytics

Nicos Angelopoulos

University of Essex, Colchester, UK
nicos.angelopoulos@essex.ac.uk

Jan Wielemaker

Centrum voor Wiskunde en Informatica (CWI), Amsterdam, The Netherlands
J.Wielemaker@cwi.nl

Delivering effective data analytics is of crucial importance to the interpretation of the multitude of biological datasets currently generated by an ever increasing number of high throughput techniques. Logic programming has much to offer in this area. Here, we detail advances that highlight two of the strengths of logical formalisms in developing data analytic solutions in biological settings: access to large relational databases and building analytical pipelines collecting graph information from multiple sources. We present significant advances on the *bio_db* package which serves biological databases as Prolog facts that can be served either by in-memory loading or via database backends. These advances include modularising the underlying architecture and the incorporation of datasets from a second organism (mouse). In addition, we introduce a number of data analytics tools that operate on these datasets and are bundled in the analysis package: *bio_analytics*. Emphasis in both packages is on ease of installation and use. We highlight the general architecture of our components based approach. An experimental graphical user interface via *SWISH* for local installation is also available. Finally, we advocate that biological data analytics is a fertile area which can drive further innovation in applied logic programming.

1 Introduction

Logic programming (LP) provides a powerful platform for analytics in data-rich areas such as biology. The ability to view data as part of the background knowledge allows LP a unique position in reasoning with relational data. In addition, the level of abstraction of Prolog's constructs along with its mathematical properties, allow for knowledge-based approaches to biological data analytics. Although LP is not a computational biology mainstream approach, still early pioneering approaches exist, such as *blipkit*, a general bioinformatics toolkit, [10] and *P\FDM* a functional data model system for biological databases [8]. These have been followed by a biological data centric library [3] and Reactome Pengine [13, 12] a modern approach exploiting the Pengine [9] web architecture to deliver the Reactome [7] database as an API that can be queried by a host of programming languages.

Recent availability of a number of interface libraries connecting Prolog to widely used software systems can play a critical role in the deployment of Prolog applications in bioinformatics. *Real* [2, 1] is an industrial strength interface to *R* [15] allowing Prolog code to access the vast array of statistical reasoning and graphics generation libraries of *R*. *Real* was also the blueprint for *Rserve* a Prolog interface to *Rserve* that provides access to *R* for web services. *Rserve* is used in the *SWISH* on-line collaborative environment. Finally, and of key importance to big data analytics, is access to external databases as the data are ever increasing in size and multiple data tables are often needed to be accessed concurrently;

B. Bogaerts, E. Erdem, P. Fodor, A. Formisano,
G. Ianni, D. Incezan, G. Vidal, A. Villanueva,
M. De Vos, F. Yang (Eds.): International
Conference on Logic Programming 2019 (ICLP'19).
EPTCS 306, 2019, pp. 309–322, doi:10.4204/EPTCS.306.36

making loading to memory impossible. *SWI-Prolog* has a plethora of options, including a generic *ODBC* interface and direct interfaces to database systems such *RocksDB*, *Berkeley DB* and *SQLite* [5].

In this contribution we extend previous work [3] on providing access to biological data. We highlight the longevity of *bio_db* and its improvement through iterative development cycles. Key to this evolving approach is the use of the library on daily basis for a variety of projects in computational biology. Furthermore, we present a new library (*bio_analytics*) that encodes a number of analytics tools which enable users to analyse their experimental data often in relation to data in *bio_db*. Both presented libraries depend on a hierarchy of independent packages which are part of an actively maintained code-base.

Critical to the successful deployment of such intricate dependencies between libraries is the *SWI-Prolog* package manager¹. Herein we use the terms library and pack to refer to single Prolog package that is managed by the *SWI-Prolog* package manager. The software described here is implemented in the *SWI-Prolog* [22] but its core parts should also be executable in Yap [6] with small changes, thanks to previous compatibility work between the two systems [21].

The remainder of this paper is as follows. Section 2 details the extension to the structure and data served by *bio_db*. Section 3 introduces an analytics library that enables the analysis of users' experimental data and which extensively uses to the datasets in *bio_db*. Section 4 provides details about the availability and usage of the presented libraries, and Section 5 presents the concluding remarks.

2 Biological data tables as Prolog facts

We bring forward our earlier work on biological databases [3] driven by the desire to incorporate data for new organisms. Originally, all data served by *bio_db* were from human (*homo sapiens*) databases. In the version described here (*bio_db* v3.0) we also provide extensive coverage of mouse (*Mus musculus*) data. In accommodating an additional organism, we radically changed the architecture of the data library and extended a number of the infrastructure packages which facilitate a compositional approach to module interfaces. Having made these adjustments, *bio_db* can now easily and intuitively be extended to include data from additional organisms.

A key observation that drove our design was that often a particular analysis of experimental data only involves access to data from a single organism. Occasionally, however, it might still be the case that we need access to data for more than one organism in a single analysis. A possibility would be to have a library for each organism, however, a single library is preferred as all code making the data available is common, with only the actual data tables being different between the two organisms. Furthermore, the module interface even for one organism, human, was becoming long, making the navigation to specific data table difficult. With the addition of mouse data it became clear that the conventional Prolog approach of listing all module predicates was ill suited. In our approach a module is comprised of a number of organism sub-*cells* each of which is further subdivided to a number of database-specific compartments. The user can at load time declare which components are to be incorporated.

2.1 Architecture

In this paper we primarily discuss two Prolog libraries: *bio_db* and *bio_analytics*. The former is the biological data provider with the latter utilising this data to perform analytic tasks on results from biological experiments. The two libraries are marked as *iface*, short for interface, in Table 1. They both depend on a number of libraries that provide functionality of some complexity marked as *mid* (short for middleware)

¹<http://www.swi-prolog.org/pack/list>

Layer	Package	Description
iface	<i>bio_analytics</i>	Computational biology data analytics.
iface	<i>bio_db</i>	Access, use and manage big, biological datasets.
	<i>bio_db_repo</i>	Data package for <i>bio_db</i> .
mid	<i>Real</i>	Integrative statistics with R.
mid	<i>mtx</i>	Working with data matrices.
mid	<i>wgraph</i>	Weighted graphs, with plotting via Real.
<i>R</i>		<i>wgraph</i> dependencies: <i>igraph</i> , <i>qgraph</i> , <i>GGally</i>
infra	<i>lib</i>	Predicate based code development.
infra	<i>options</i>	Options handling.
infra	<i>os_lib</i>	Operating system interaction predicates.
infra	<i>stoics_lib</i>	A medley of library predicates for stoics packs.

Table 1: Package architecture

and a number of libraries providing basic functionality, marked on Table 1 as *infra* (for infrastructure). All the Prolog libraries are developed in-house and provided as independent *SWI-Prolog* packages. Libraries *Real*, *options*, and *os_lib* were already mature libraries that saw marginal improvements from the current work. Whereas, libraries *lib*, *mtx* and *wgraph* were substantially extended within the work described here.

Our approach differs to other similar size projects in that substantial and important parts are made available as stand alone libraries. Past projects with comparable size and some biological hue such as *blipkit* [10, 20] and *P\FDM*, would have all components within a single code-base. *blipkit* was designed as a general bioinformatics toolkit with emphasis on database querying and particular strengths in ontologies encoding biological knowledge. *P\FDM* [8] was a complete functional database system implemented in Prolog with a number of applications in storing biological data.

bio_db contains a number of useful features. Its primary functionality is to serve biological data from high quality biological databases. Data are served as predicates which hold the source database tables as facts. The information pertains to biological products, their features and relationships amongst products. For instance, query

```
?- map_unip_hgnc_unip( Hgnc, Unip ).
```

interrogates the relation between genes and proteins in human.

2.1.1 Incremental

At installation, *bio_db* comes with no actual data for its data predicates (such a predicate is the one introduced above: *map_unip_hgnc_unip(Hgnc,Unip)*). The data tables can be downloaded either en-mass via library *bio_db_repo*, or the user will be prompted for the library to automatically download the specific table. So, in our running example, the first time *bio_db* attempts the previous query, the interaction with system will be as follows:

```
?- map_unip_hgnc_unip( Hgnc, Unip ).
% prolog DB:table unip:map_unip_hgnc_unip/2 is_not_installed,
```

```

do you want to download it (Y/n) ?

% Continuing with: yes
% Downloading dataset from server:
      http://stoics.org.uk/~nicos/sware/packs/bio_db_repo/data
% Delete the zip file:  '../map_unip_hgnc_unip.pl.zip'    (y/N) ?
% Continuing with: yes
Hgnc =5,
Unip = 'MOR009' ...

```

2.1.2 Backends

The library's native representation of each data table is as plain Prolog fact bases. The data are stored as compressed Prolog files on the server² and can also be accessed directly. In addition to loading the fact bases as Prolog facts *bio_db* allows to store them in a variety of external databases: *RocksDB*, *Berkeley DB* and *SQLite*. The user can easily control the database backend but the underlying data are accessed in identical manner irrespective of the backend being used. The availability of a range of backends provides alternatives that are suitable to different use cases. The zero configuration approach of *SQLite* along with each pervasiveness might best suit packaged distributions of code whereas the performance of *RocksDB* might be best suited for web services scenarios. As shown previously [3] when memory is available and loading time is not a significant overhead, Prolog is by far the fastest way to interrogate the data tables.

2.1.3 Rebuilds

The scripts for building the whole of *bio_db_repo* are provided in *bio_db* (*auxil/build_repo*). These depend on the libraries in Table 1 as well as a few more packages also publicly available through the *SWI-Prolog* package manager. The benefits to the user from having the build scripts are twofold. First, they can rebuild any of the data tables at arbitrary time points without having to wait for the official half yearly updates of *bio_db_repo*. In addition, it is quite straight forward to adapt the scripts to both (a) derive new relations from the biological databases already accessed through *bio_db* and (b) to incorporate additional tables from new biological sources.

2.1.4 Hot-swapping

At loading time, data predicates in *bio_db* are merely place holder code snippets that will be hot swapped at run time for either the fact bases, if the backend is Prolog, or database serving code. In the former case the swapping code for the user query *Call* essentially reads as:

```

bio_db_load_call(true,Pname,Arity,Iface,PIFile,Call):-
  functor(Phead, Pname, Arity),
  atom_concat(Pname, '_info', InfoPname),
  dynamic(bio_db:InfoPname/2),
  functor(InfoHead, InfoPname, 2),
  abolish(bio_db:Pname/Arity),
  retractall(bio_db:InfoHead),

```

²http://stoics.org.uk/~nicos/sware/packs/bio_db_repo/data/

```

bio_db_ensure_loaded(Iface,Pname/Arity,PFile,Handle,From),
assert(bio_db_handle(Pname/Arity,Iface,File,Handle,From)),
call(Call).

```

This predicate first manages the *_info* postfixed predicate that is associated with each table, before retracting the current stub for the table predicate and loading the serving code (for the Prolog backed *bio_db_ensure_loaded/6* simply calls *ensure_loaded(PFile)*). Finally, the original user call is called in the last line of the above code (*call(Call)*).

2.1.5 Compositional

bio_db v3.0 introduces compositional loading based on recent extensions in library *lib*. As per normal convention in *SWI-Prolog bio_db* defines a single module. However the user can specify from a number of hierarchical sub-components organised on two tiers. The first tier is organised across organism lines and it is further sub-divided according to database of origin lines. Each component at either tier can be an independent entry point. Components can be loaded incrementally.

```
?- lib(& bio_db(hs)).
```

will load all table predicates that are related to human (*hs*), whereas

```
?- lib(& bio_db(mouse/mgim)).
```

will load the mouse data that are originating in the mouse genomic initiative (*MGI*) database. Either of the following two queries will load the full library

```
?- lib(bio_db).
```

```
?- use_module(library(bio_db)).
```

Library *lib* locates sub-components relative to the *cell/* directory of packs. The main module points to the first tier components (here at the organism level), each signposted by a source file possibly pointing to lower level components, thus creating a hierarchy of dependencies. As the interface has expanded, *bio_db* provides *bio_db_data_predicate/4* for interrogating the available data tables:

```
?- bio_db_data_predicate(Pname,Parity,Org,File),
   write(Pname/Parity:Org:File), nl, fail.
```

```

edge_strg_hs_symb/3:hs:hs/strg.pl
map_unip_unip_hgnc/2:hs:hs/unip.pl
map_mgim_mouse_mgim_unip/2:mouse:mouse/mgim.pl
...

```

2.2 Source databases

bio_db harvests data from curated biological databases that provide accurate, up to date information: *HGNC* [4], *UniProt* [19], *Ensembl NCBI* [11], *MGI* [16] and *STRING* [17].

Central to both the organisms supported, are gene identifiers and gene names. *bio_db* sources its identifiers and unique short gene name (also known as symbols) from the most respected and well curated

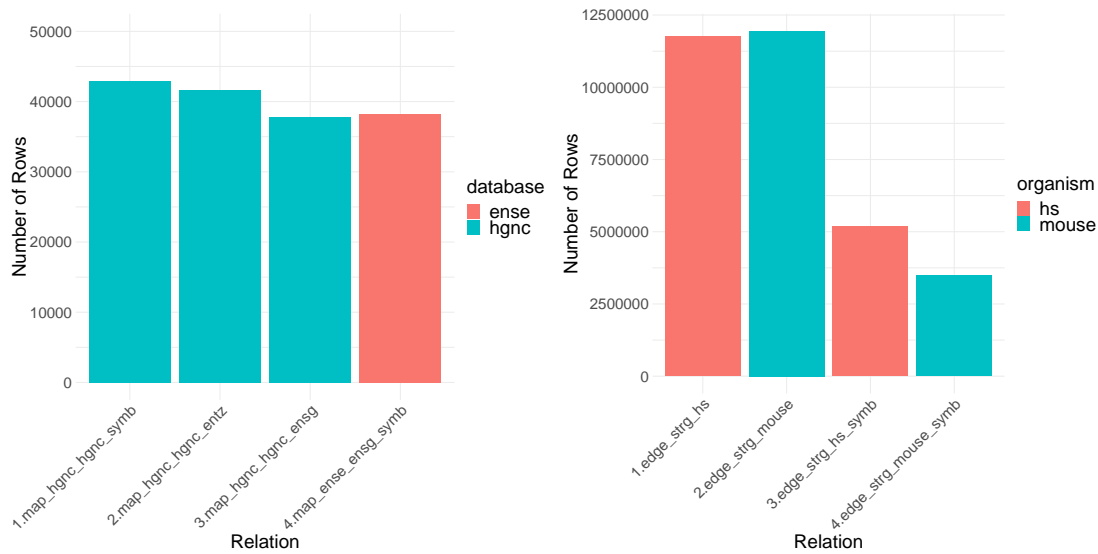


Figure 1: Population of database tables served as Prolog facts. (LEFT) Slight variations on the number of entries between different gene name identifiers and contrast of same relation derived from different databases (hgnc versus ense, last two bars). (RIGHT) Population statistics for the *STRING* database illustrating the differences between protein and gene entries and across organisms.

database for each organism: HUGO Gene Nomenclature Committee (HGNC) for human genes and Mouse Genome Informatics (*MGI*) for mouse. Each provide a unique integer as identifier and short alphanumeric as gene symbols. Human gene symbols are formed from capital letters and numbers while for the most part mouse symbols are formed from one capital letter followed by lower case letters and numbers. Human symbols are standardised to a larger extent whereas mouse symbols are less well curated, with a larger proportion of genes still having non standardised symbols. For the majority of studied genes that have been identified in both organisms, there is a correspondence in their symbols.

```
?- map_hgnc_hgnc_symb(19295, HsSymb).
HsSymb = 'LMTK3'.
```

```
?- map_mgim_mouse_mgim_symb(3039582, MmSymb).
MmSymb = 'Lmtk3'.
```

The left panel of Figure 1 shows the row population of tables translating between human gene identifiers. Although all databases are well curated and research in human is intense, there still are some discrepancies in the number of genes we can map from one db identifier to another. Ideally the four bars of the plot should be of equal height.

Data predicates follow a uniform naming convention, each consisting of underscore separated parts. The first one, declares the type of data which can be either a *map* or an *edge*. Then comes the source database: *ense*, *gont*, *hgnc*, *mgim*, *ncbi*, *pros*, *strg*, and *unip*. The third part is the organism: *hs* or *mouse*. However in the case of human data, the token (*hs*) is not included in the interest of backward compatibility and brevity. The remaining parts of the data predicate name depend on the type of relation. When this is a map, there are two tokens: one indicating the object biologic product and the other the attributes

or subject biological product. The main identifier field in a database is named by the same token as the database itself. Tokens for databases and tokens for map predicates comprise of four letters.

Proteins are the basic functional units of organisms. They are built from the blueprint encoded by genes stored in the DNA. In the general case the relation between genes and proteins is a many-to-many one. However, it is more often the case that each gene maps to a number of proteins; some of which are well studied and some that are less so. For example:

```
?- map_hgnc_symb_hgnc('LMTK3', Hgnc ), map_unip_hgnc_unip( Hgnc, Prot ).
Hgnc = 19295, Prot = 'A0A0A0MQW5' ;
Hgnc = 19295, Prot = 'A0A3B3IRV9' ;
Hgnc = 19295, Prot = 'A0A3B3ISL5' ;
Hgnc = 19295, Prot = 'A0A3B3ITQ7' ;
Hgnc = 19295, Prot = 'Q96Q04'.
```

```
?- map_unip_hgnc_unip( Hgnc, 'Q96Q04' ).
Hgnc = 19295.
```

where *Q96Q04* is a well studied protein. Uniprot is the primary source of information on proteins and it contains two parts: a curated section of high quality information and a non-curated part that is less well studied. In the above example *Q96Q04* is the only entry from the curated part. The example shows that although *LMTK3* maps to possibly many structures, protein *Q96Q04* is associated with a single gene.

Effectively, *bio_db* implements a straight forward functional data model with simple mapping between product and their functional elements, or between products, while ensuring there is a path between any two entities that should be connected. The above example finds proteins of the symbol *LMTK3* by first navigating to its unique *HGNC* id and then find proteins of that gene identifier. The library does not have to provide a specific predicate to map between symbols and proteins.

Gene ontology (GO) [18] codifies a restricted language for describing biological processes. *bio_db* has two types of *GO* predicates. First, a predicate for the membership of genes to gene ontology terms

```
?- map_gont_symb_gont('LMTK3', Gont),
   map_gont_gont_gonm(Gont,Gonm), write(Gont:Gonm), nl, fail.
```

```
GO:0000139:Golgi membrane
GO:0003674:molecular function
GO:0004674:protein serine/threonine kinase activity
GO:0005515:protein binding
GO:0005524:ATP binding
...
GO:0030424:axon
GO:0030425:dendrite
GO:0046872:metal ion binding
false.
```

and relations between gene ontology terms, such as

```
?- edge_gont_is_a(139,Gont), map_gont_gont_gonm(139, GnmI),
```

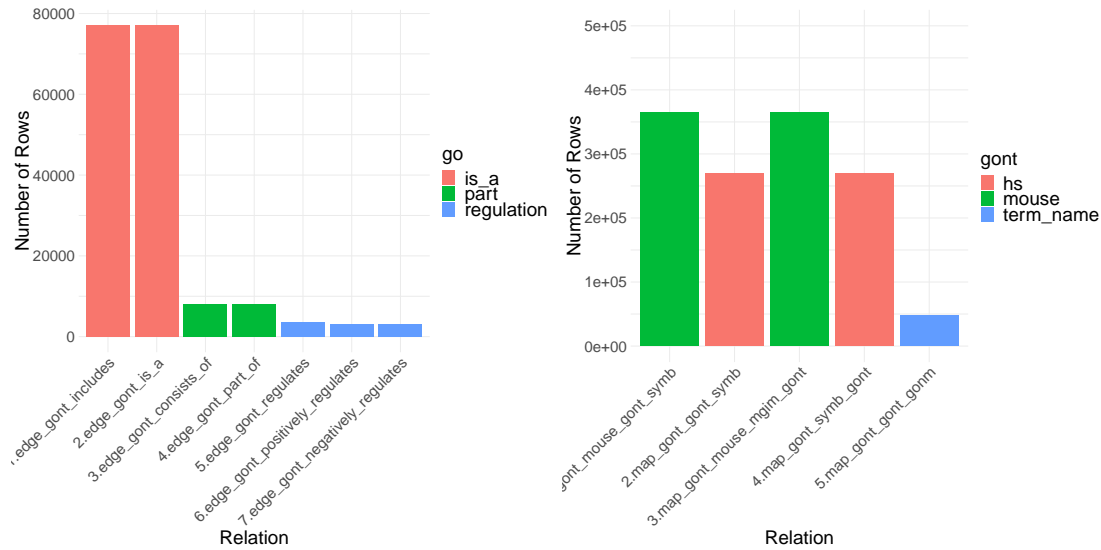


Figure 2: Gene ontology relations in *bio_db*. (LEFT) Plots number of rows for tables that encode the ontological relationships in GO. (RIGHT) Comparative plot of genes and proteins, in each GO term for human and mouse (last bar shows the number of GO terms).

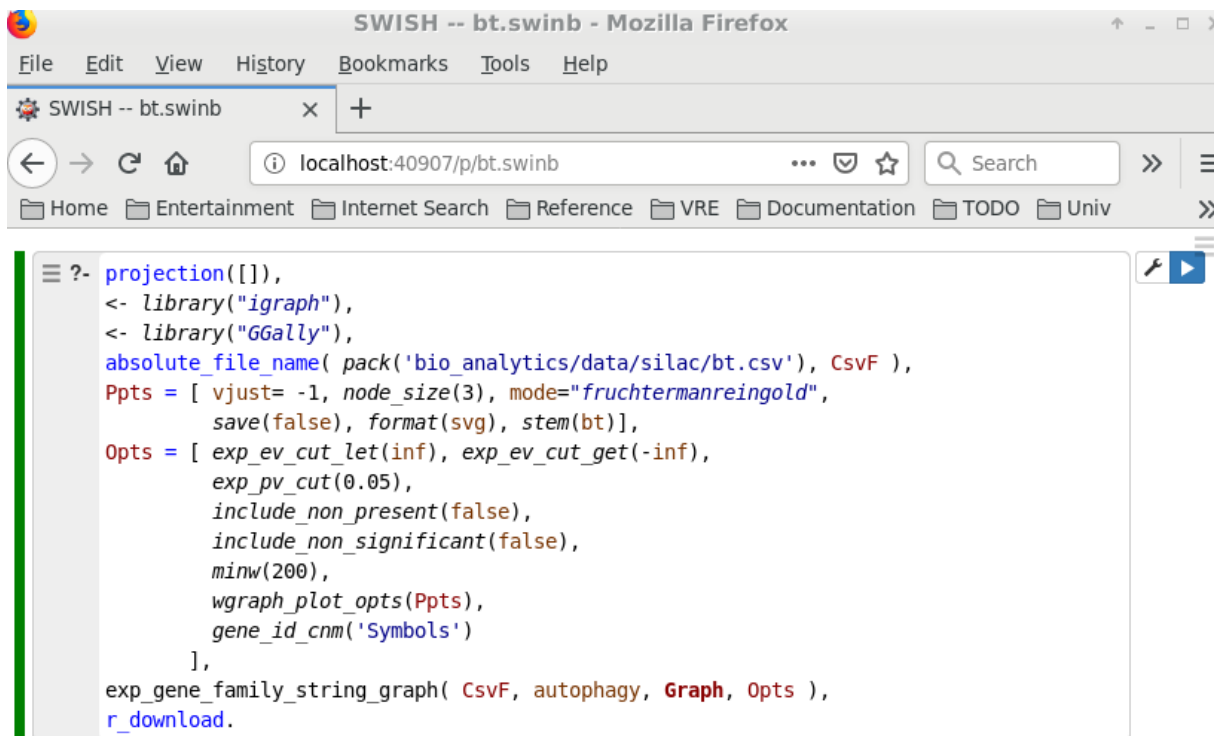
```
map_gont_gont_gonm(Gont,Gnm2).
Gont = 44431, Gnm1 = 'Golgi_membrane',
Gnm2 = 'Golgi_apparatus_part' ;
Gont = 98588, Gnm1 = 'Golgi_membrane',
Gnm2 = 'bounding_membrane_of_organelle'.
```

Figure 2 plots the populations for data predicates from the gene ontology database: hierarchical relations between ontology terms (left) and membership relations of genes and proteins to GO terms (right) for human (red), mouse (green) and ontology terms (blue).

STRING [17] is a protein-protein interaction database that provides weighted relations between proteins. Each weight is in range (0,1000) where the highest the number the strongest the evidence that there exists a direct physical interaction between two proteins. This relation can be seen as encoding a graph amongst proteins.

3 Analytics

Building on the wealth of data provided in *bio_db*, library *bio_analytics* allows experimental results to be analysed against data in *bio_db*. It implements a number of low level predicates that (a) interface with predicates responsible for reading experimental results from *csv* files, and (b) allow the filtering the read-in results according to thresholds on significance of differential expression and level of fold change. The end result of applying such filters will be a list of genes or proteins that are significantly up or down regulated in a condition against a control. We will refer to these sets as an experiment's hits-list. If the experimental reads are proteins, as is the case in our example, they will then need to be mapped to genes, as these are often what experimentalists are interested in. Each thus identified gene will commonly have attached to it an arithmetic value indicating its fold change or relative expression.



```
projection([],  
<- library("igraph"),  
<- library("GGally"),  
absolute_file_name( pack('bio_analytics/data/silac/bt.csv'), CsvF ),  
Ppts = [ vjust= -1, node_size(3), mode="fruchtermanreingold",  
save(false), format(svg), stem(bt)],  
Opts = [ exp_ev_cut_let(inf), exp_ev_cut_get(-inf),  
exp_pv_cut(0.05),  
include_non_present(false),  
include_non_significant(false),  
minw(200),  
wgraph_plot_opts(Ppts),  
gene_id_cnm('Symbols')  
],  
exp_gene_family_string_graph( CsvF, autophagy, Graph, Opts ),  
r_download.
```

Figure 3: Code exemplifying the use of options in controlling the behaviour of predicate `exp_gene_family_string_graph/4`.

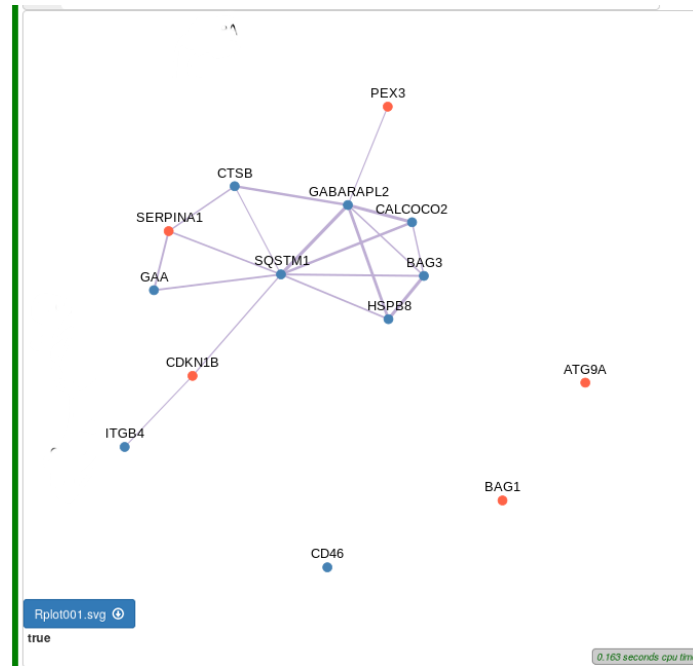


Figure 4: STRING edges for genes in experiment data/silac/bt.csv that are present in the autophagy family. Red nodes are up regulated genes and blue nodes are for down-regulated genes. The graph is displayed via the SWISH interface.

Higher level predicates in *bio_analytics* build on these primitives to enable complex analysis tasks. Through the use of options lists, *pack(options)*, behaviour and specific cut offs can be cascaded and influence all code participating in an analysis. *exp_gene_family_string_graph/4* implements such a complex analysis. A call of the form:

?- *exp_gene_family_string_graph(CsvF, Family, Graph, Opts)*.

will, first, identify a number of proteins with significant fold change in file *CsvF* and then select those that participate in a specified gene family (*Family*). *Graph* is the constructed weighted *STRING* graph of genes in *Family* which also appear in the experimental significant hits-list. The example in Figure 3 is a call to *exp_gene_family_string_graph/4* with a number of options set as to control the process of graph and plot generation (also see source file *examples/bt.pl*). In this example we see the effect of the experiment on *autophagy* pathway/family (file *data/families/autophagy.csv*). In Figure 4, significantly up-regulated genes are shown in red and down regulated gene are shown in blue.

Biologists can thus visually inspect the effects of their experiment on a gene family of interest. The family can either be a collection of genes forming a pathway or it could be a gene ontology term selected from data-driven interrogation of the experimental proteins against ontology terms (over-representation analysis). The list of option terms *Opts*, controls aspects of both low level (e.g. significance level cut off) and high level behaviour (e.g. include non identified genes in the family). In addition to returning the graph, the predicate can also produce images of the generated graph in a variety of formats (*pack(wgraph)* and *pack(Real)*). Plot options are collected in list *Ppts*, controlling aspects such as the size of nodes (3) and the format of the output image (svg). The above example is included in the *bio_analytics* sources: *examples/bt.pl*. The experimental data of this example are also included in file: *data/silac/bt*.

GOBPID	Pvalue	OddsRatio	Count	Size	Term
GO:0061387	0.00018	Inf	8	8	regulation of extent of cell growth
GO:0030516	0.00055	Inf	7	7	regulation of axon extension
GO:0048675	0.00055	Inf	7	7	axon extension
GO:0007275	0.00056	1.59	132	315	multicellular organism development
GO:0032501	0.00067	1.54	166	411	multicellular organismal process
GO:0008361	0.00087	7.17	11	14	regulation of cell size
...					
GO:0048729	0.00682	2.27	23	43	tissue morphogenesis
...					

Table 2: Table of statistical results showing *GO* terms that are significantly over-represented in the hit-list of an experiment.

csv [14].

Building on the analysis described so far, predicate *exp_go_over_string_graphs/4* combines gene ontology term over representation with graphs highlighting the *STRING* interactions. Each *GO* term is examined in turn, and the number of genes that significantly up or down regulated within an experiment are counted against the background population. Statistics on the likelihood that a particular term is significantly over represented in the gene hit-list for the experiment are gathered in tables similar to that shown in Table 2. For each *GO* term that is significant, the predicate generates the *STRING* network of the genes identified in the experiment and fall within the *GO* term. The colour and intensity of each node shows the direction of dysregulation. Figure 5 shows such a network for the *GO* term, *GO:0048729*. An executable example can be found in the library sources (`examples/bt_over_go_string.pl`).

4 Availability

The software described herein comprises of a number of layered packages. Each is available as an *SWI-Prolog* package that can be installed from with Prolog by a single, simple query. The two main packages can be installed by:

```
?- pack_install(bio_db).
?- pack_install(bio_analytics).
```

Loading of libraries can be either via:

```
?- use_module(library(bio_db)).
or
?- use_module(library(lib)).
?- lib(bio_db).
```

The user will be asked if each missing Prolog and *R* libraries should be installed at loading time. This is done via library *lib*. This library also allows loading *suggested* code. The main idea is that these are softer dependencies that implement fringe features of the package or require substantial external libraries which the user might want to avoid loading until it is certain that their functionality is required.

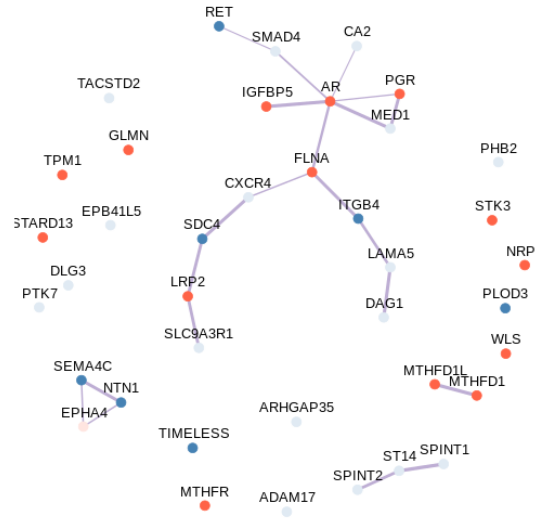


Figure 5: The graph of *STRING* edges for all genes that appear in experiment data/silac/bt.csv and belong to gene ontology term GO:0048729, which is an over-represented term (Table 2) for the experiment's hits-list.

?- use_module(library(lib)).

?- current_prolog_flag(lib_suggests_warns, WarnFlag).

WarnFlag = auto.

A value of *debug* will inform of all suggested code a particular query depends on (both Prolog and *R* dependencies which are appropriately marked) and when the flag is set to *true* the library will interrogate user whether they wish to install each specific dependency. Although not demonstrated here, library *lib* also implements *promised* code³, where the source of a predicate is not loaded until the predicate is called. At call time a stub is used to (a) delete itself from the memory, (b) load the appropriate sources and (c) re-evaluate the call.

The packages are available from our server (<http://stoics.org.uk/~nicos/sware>) and the source code is also available from github (<https://github.com/nicos-angelopoulos>). We also make available an experimental *SWISH* based graphical interface⁴ that currently allows the user to run the provided examples and interrogate the data tables by graphical means. We plan to expand this interface to both allow users to run similar queries to their own data. We will also extend this interface for secure hosting on public-facing web servers.

The data in *bio_db_repo* are updated at least twice a year, traditionally in early autumn and early spring. The library is entering its fifth year of availability and it has been recently substantially extended to incorporate a new organism. Furthermore, we are committed to providing and extending the data and building Prolog-based analytical tools on top of this library.

Although developers can be side tracked by infrastructure packages such as *lib* and *options* (1) our strategy of breaking code to independent packages has the distinct benefit of making the libraries available to other users and projects. It also means that our libraries can depend on pre-existing stable pack-

³<http://stoics.org.uk/~nicos/sware/lib/>

⁴http://stoics.org.uk/~nicos/sware/bio_db/swish_install.html

ages that have a strong user base: *Real* (354 downloads) and *proSQLite* (464). The *SWI-Prolog* package manager⁵ has been an invaluable tool for managing the whole hierarchy of libraries and enabling wider audiences for the code. The package manager currently provides access to 290 packages and has 14,094 registered downloads.

5 Conclusions

We have presented substantial recent advances to the availability of high-quality, big, biological datasets in logic programming along with easily accessible data-analytics workflows for addressing important biological questions. The architecture of the data providing package introduced a number of innovations to module loading. Data on new organisms can now be added easily with no disruption to existing data or the structure of the library. We also expect to soon incorporate more complex databases such as the Reactome database [7]. In the future we also plan to provide collaborative web presences for *bio_analytics* via *SWISH* by extending the current *SWISH* interface. Such tools will enable experimental scientists to easily analyse their data without learning Prolog. In this direction, it will also be interesting to explore integration with existing web services such as the Reactome Penguins [13, 12].

Although breaking big projects to small independent libraries is more time consuming, we have found to be an extremely useful exercise as each unit of code becomes more independent and re-usable. Even for basic functionality such as loading Prolog code which is often best left to the system, it is useful in terms of research to sometimes re-examine and consider extensions to these. In our case, pack *lib* started as a simple tool that went against the *SWI-Prolog* mantra of single-file, single-module libraries. It was then substantially extended to provide new functionalities some of which are described in this paper.

Taken together, the presented libraries provide logic based building blocks for constructing high-quality data analytic tools that are useful to biologists. Prolog can both play an important role in representation and reasoning with biological knowledge but also computational biology can play a crucial role in the future of logic programming as a driver application area.

References

- [1] N. Angelopoulos, S. Abdallah & G. Giamas (2016): *Advances in integrative statistics for logic programming*. *Int. J. of Approximate Reasoning* 78, pp. 103–115, doi:10.1016/j.ijar.2016.06.008.
- [2] N. Angelopoulos, V. S. Costa, J. Azevedo, J. Wielemaker, R. Camacho & L. Wessels (2013): *Integrative functional statistics in logic programming*. In: *Proc. of Prac. Asp. of Decl. Lang., LNCS 7752*, Springer, pp. 190–205, doi:10.1007/978-3-642-45284-0_13.
- [3] Nicos Angelopoulos & Jan Wielemaker (2017): *Accessing biological data as Prolog facts*. In: *19th International Symposium on Principles and Practice of Declarative Programming (PPDP 2017)*, ACM digital library., Namur, Belgium, pp. 29–38, doi:10.1145/3131851.3131857.
- [4] B Braschi, P Denny, K Gray, T Jones, R Seal, S Tweedie, B Yates & E Bruford (2019): *Genenames.org: the HGNC and VGNC resources in 2019*. *Nucleic Acids Research* 47, pp. D786–92, doi:10.1093/nar/gky930.
- [5] S. Canisius, N. Angelopoulos & L. Wessels (2013): *ProSQLite: Prolog file based databases via an SQLite interface*. In: *Proc. of Prac. Aspects of Decl. Lang., LNCS 7752*, Springer, pp. 222–7, doi:10.1007/978-3-642-45284-0_15.
- [6] Vítor Santos Costa, Ricardo Rocha & Luís Damas (2012): *The YAP Prolog system. Theory and Practice of Logic Programming* 12, pp. 5–34, doi:10.1017/S1471068411000512.

⁵<http://www.swi-prolog.org/pldoc/man?section=prologpack>

- [7] David Croft, Antonio Fabregat Mundo, Robin Haw, Marija Milacic et al. (2014): *The Reactome pathway Knowledgebase*. *Nucleic Acids Research* 42(D1), pp. D472–D477, doi:10.1093/nar/gkt1102.
- [8] G. J. L. Kemp, J. Dupont & P. M. D. Gray (1996): *Using the functional data model to integrate distributed biological data sources*. In: *Proceedings of 8th International Conference on Scientific and Statistical Data Base Management*, ACM digital library, Budapest, Hungary, pp. 176–185, doi:10.1109/SSDM.1996.506060.
- [9] Torbjorn Lager & Jan Wielemaker (2014): *Pengines: Web Logic Programming Made Easy*. *Theory and Practice of Logic Programming* 14(4-5), pp. 539–552, doi:10.1017/S1471068414000192.
- [10] Chris Mungall (2009): *Experiences Using Logic Programming in Bioinformatics*. In Patricia M. Hill & David S. Warren, editors: *Int. Conf. on Logic Programming, LNCS 5649*, Springer, Berlin, pp. 1–21, doi:10.1007/978-3-540-39718-2_41.
- [11] NCBI Resource Coordinators (2013): *Database resources of the National Center for Biotechnology Information*. *Nucleic Acids Research* 41(Database issue), pp. D8–D20, doi:10.1093/nar/gks1189.
- [12] Sam Neaves (2019): *Explorations in Logic Programming for Bioinformatics*. Ph.D. thesis, Kings’s College. Department of Computer Science, London, UK.
- [13] Samuel R Neaves, Sophia Tsoka & Louise A C Millard (2018): *Reactome Pengine: a web-logic API to the Homo sapiens reactome*. *Bioinformatics* 34, p. 28562858, doi:10.1093/bioinformatics/bty181.
- [14] J. Nunes, H. Zhang, N. Angelopoulos, J. Chetri, C. Osipo, J. Stebbing & G. Giamas (2016): *ATG9A loss confers resistance to trastuzumab via c-Cbl mediated Her2 degradation*. *Oncotarget* 7, pp. 27599–27612, doi:10.18632/oncotarget.8504.
- [15] R Core Team (2019): *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- [16] C L Smith, J A Kadin et al. (2018): *Mouse Genome Database (MGD) 2019*. *Nucleic Acids Research* 47(D1), pp. D801–D806, doi:10.1093/nar/gky1056. Database issue.
- [17] D Szklarczyk, A Franceschini et al. (2015): *STRING v10: protein-protein interaction networks, integrated over the tree of life*. *Nucleic Acids Research* 43(D1), pp. D447–D452, doi:10.1093/nar/gku1003.
- [18] The Gene Ontology Consortium (2000): *Gene ontology: tool for the unification of biology*. *Nat. Genet.* 25(1), pp. 25–9, doi:10.1091/mbc.9.12.3273.
- [19] The UniProt Consortium (2015): *UniProt: a hub for protein information*. *Nucleic Acids Research* 43, pp. D204–D212, doi:10.1093/nar/gku989.
- [20] Vangelis Vassiliadis, Jan Wielemaker & Chris Mungall (2009): *Processing OWL2 Ontologies using Thea: an Application of Logic Programming*. In: *Proceedings of the 6th International Conference on OWL: Experiences and Directions (OWLED’09)*, 529, pp. 89–98.
- [21] Jan Wielemaker & Vítor Santos Costa (2011): *On the Portability of Prolog Applications*. In: *Practical Aspects of Declarative Lang.*, LNCS 6539, Springer, Berlin, pp. 69–83, doi:10.1007/978-3-642-18378-2_8.
- [22] Jan Wielemaker, Tom Schrijvers, Markus Triska & Torbjörn Lager (2012): *SWI-Prolog*. *Theory and Practice of Logic Programming* 12(1-2), pp. 67–96, doi:10.1017/S1471068411000494.