# FlyMine

## An integrated database for *Drosophila* and *Anopheles* genomics
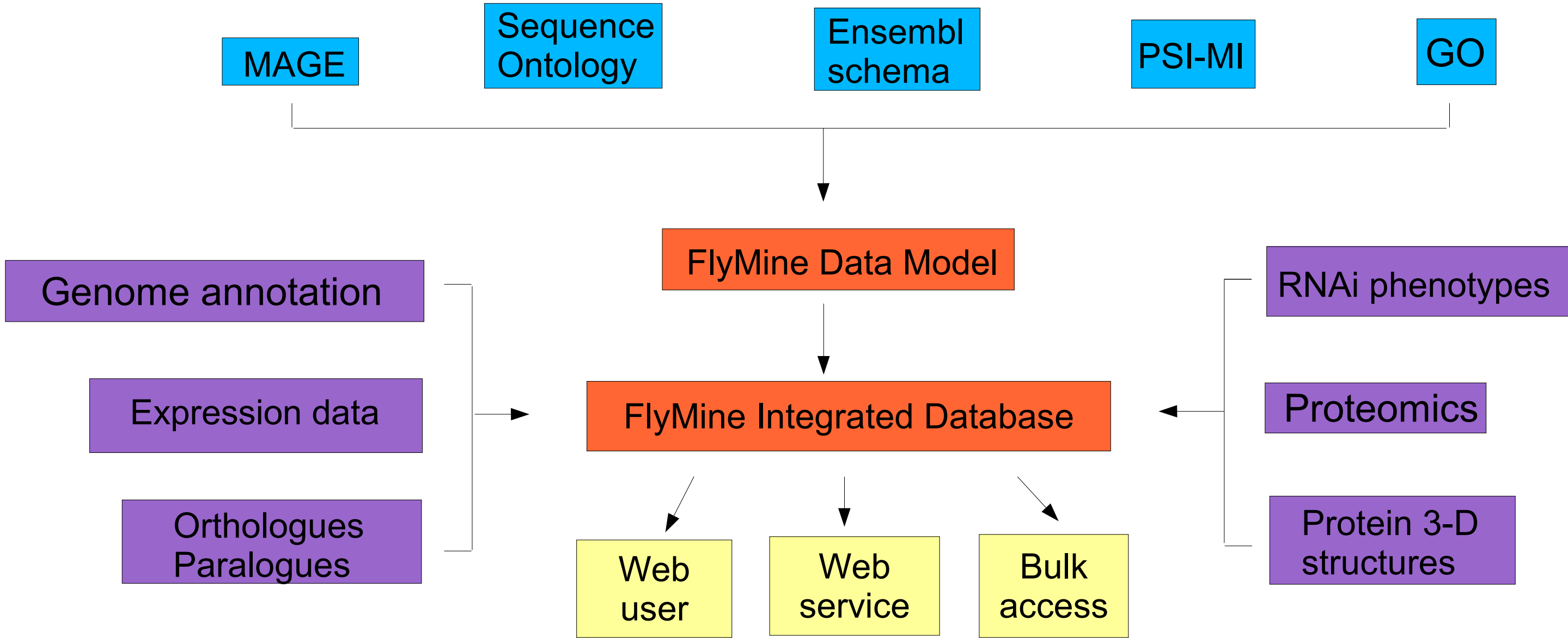
# www.flymine.org

## Introduction

FlyMine is an open-source project to build an integrated database of genomic, expression and proteomics data for *Drosophila, Anopheles* and other insects.

With the completion of increasing numbers of genomic sequences has come an explosion in the development of both computational and experimental high-throughput techniques for deciphering molecular functions and their interactions. These techniques include theoretical methods for deducing function such as analysis of protein homologies, structural domain predictions, phylogenetic profiling and protein domain fusions and empirical methods such as microarray gene expression studies, two-hybrid protein-protein interaction experiments, protein co-immunoprecipitation and large-scale RNAi screens.

The result is a huge amount of data, and the challenge is to extract meaniful knowledge and patterns of biological significance that can lead to new experimentally testable hypotheses. Many of the resulting high-throughput data sets, however, are noisy and the data quality can vary significanlty. Stronger inferences can be made if one is able to combine data from different sources. Currently, biological data is stored in a wide variety of formats in numerous different places, making complex queries across these data sources very difficult. One of the current challenges for bioinformatics is therefore to develop resources for the integration and combined analysis of such data.

## Technical Description

• FlyMine is built using InterMine (www.intermine.org), an open-source data warehouse system for integrating data from many sources into one object-based database, allowing complex queries to be performed on the integrated data.

• Although FlyMine primarily focuses on *Drosophila* and *Anopheles* data, the InterMine software is designed to be completely generic and is not tied to any particular datasource or application.

• The FlyMine model is generated from existing and emerging ontologies and standards (MAGE, SO, PSI and others). This enables data to be loaded from a variety of sources and provides the flexibility to incorporate new data as they emerge. A number of FlyMine specific classes help to link together those generated from the merged ontologies and standards.

• A powerful and flexible query interface enables users to perform arbitrary and complex queries across the data either via a web interface or a programming interface. Users are not confined to 'fill-in-the-blanks' query templates or predefined queries.
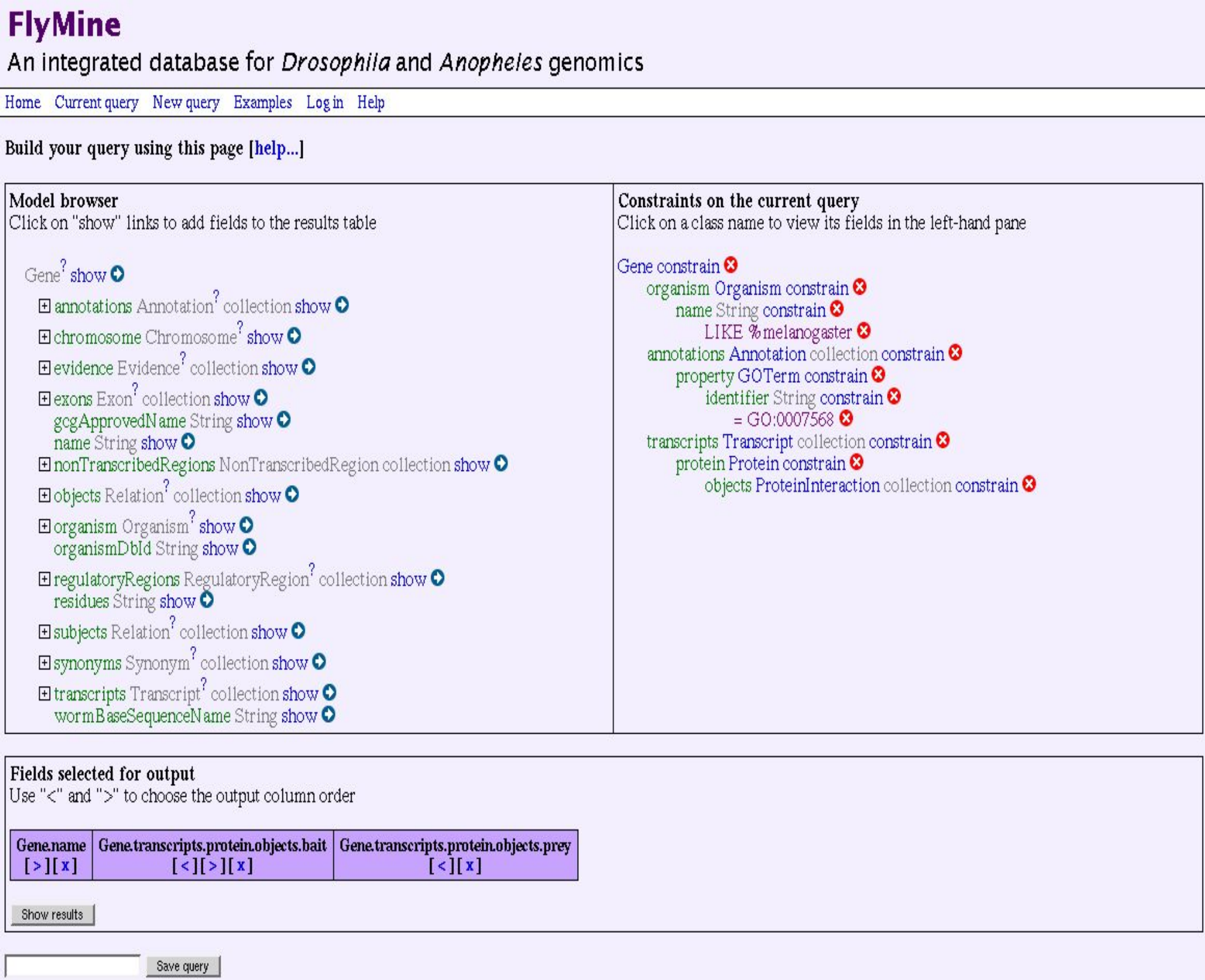


## The FlyMine Query Builder Page

• The query interface consists of three sections: a model browser, a constraints list and an output fields list.

• The model browser allows you to select the classes that contain the objects you require. You can navigate to other classes using the references between them.

• From the model browser classes and fields can be selected to constrain in the constraints list and/or show in the results output.

• Queries can be saved for future reference or futher refinement.



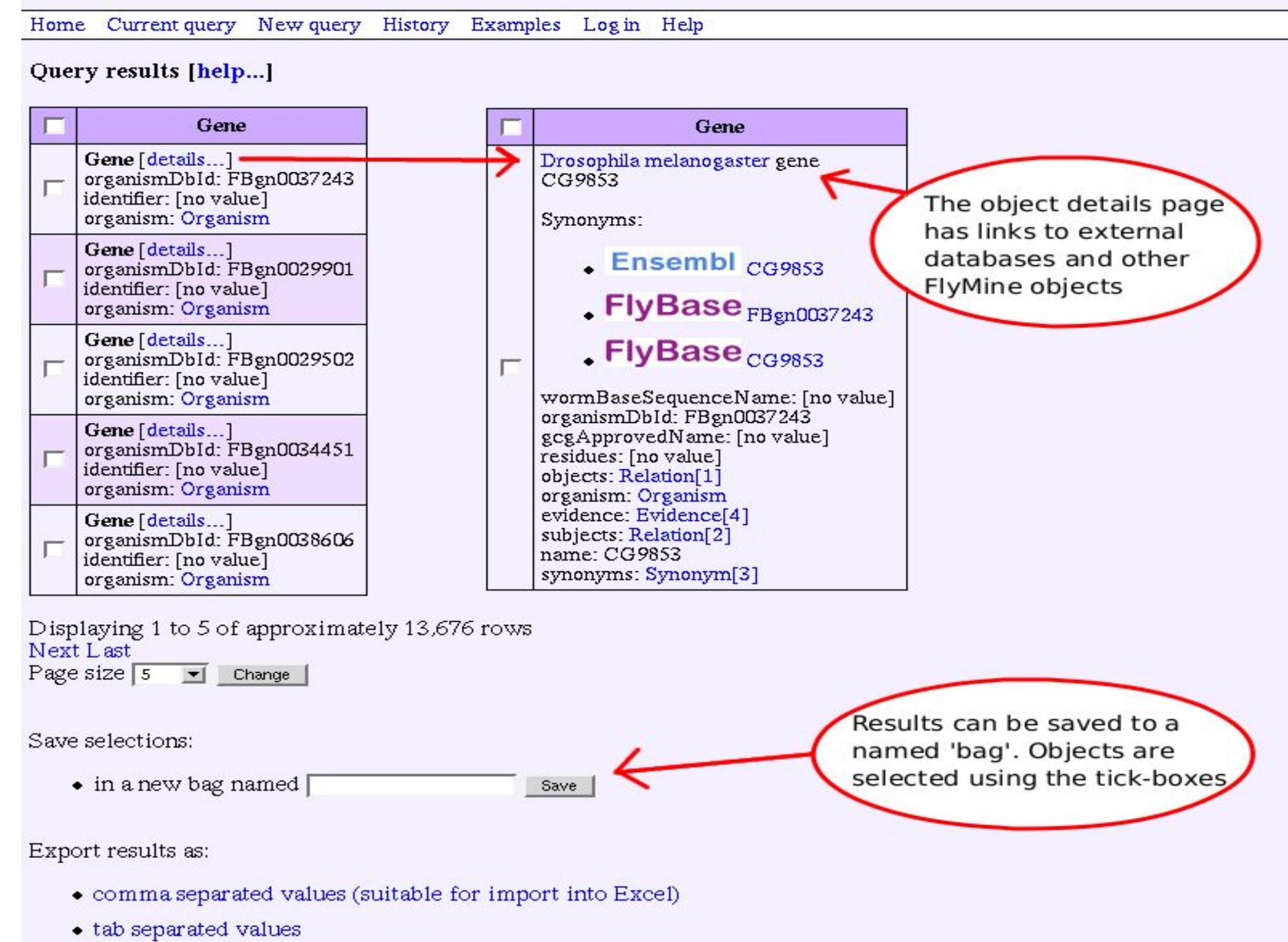## A Query Example: building a query across integrated data sources

The screenshot below illustrates how to construct the query:

**"For *Drosophila* genes that have been annotated with GO term GO:0007568, and whose proteins have protein interaction data, show the gene name and the proteins involved in the interactions".**



## The Result Display

• The result display includes links to external databases and class and attribute browsing capability. Thus the interface allows for both well-defined hypothesis driven analysis as well as more exploratory analysis.

• Integration of display and analysis tools is in progress, and will allow users to upload their results to a suitable viewing or analysis tool.

• Bags of data can be saved and used to constrain further queries and perform logical operations. A user log-in system enables queries and bags of data to be saved between sessions.

• Data can be downloaded in a variety of formats, for example tab-delimited text for use in spreadsheets.



## Data Sources First Release

| Data | Organism | Source | Publication | Website |
|---|---|---|---|---|
| Genome Annotation | *D. melanogaster* | Ensembl | | http://www.ensembl.org |
| | *A. gambiae* | Ensembl | | http://www.ensembl.org |
| 2-hybrid protein interactions | *D. melanogaster* | intAct | Giot et al (2003) Science 302 1727-1736 | http://www.ebi.ac.uk/index.html |
| | *D. melanogaster* | HybriGenics | | http://pim.hybrigenics.com/pimriderext/droso/index.html |
| | *C. elegans* | intAct | Li et al (2004) Science 303 540-543 | http://inparanoid.cgb.ki.se/index.html |
| Microarray expression data | *D. melanogaster* | ArrayExpress | | http://www.ebi.ac.uk/arrayexpress |
| RNAi phenotypes | *C. elegans* | WormBase | Kamath et al (2003) Nature 421 231-237 | http://www.wormbase.org |
| | | | Fraser et al (2000) Nature 408 325-330 | |
| | | | Simmer et al (2003) PloS Biol 1 (1) E12 | |
| Orthologues/ Paralogues | *D. melanogaster:A. gambiae* | Ensembl | | http://www.ensembl.org |
| | *D. melanogaster:C. elegans* | InParanoid | | http://inparanoid.cgb.ki.se/index.html |

## Further Information and Download

• Further information and documentation can be found on the FlyMine website, www.flymine.org, by joining one of the FlyMine electronic mailing lists (details on the website), or by email to info@flymine.org. The intermine code is available for download under the open-source LGPL licence from www.intermine.org. We welcome co-developers.

• Help is available in the form of a quick start guide, query examples, query tutorials and a user manual. There is also an analysis support help desk. The help pages are under continued development.

• The FlyMine team: Rachel Lyne, Richard Smith, Kim Rutherford, Matthew Wakeling, Mark Woodbridge, Wenyan Ji, François Guillier, and Gos Micklem.