# Function-centered approaches for finding and analyzing genes within FlyBase
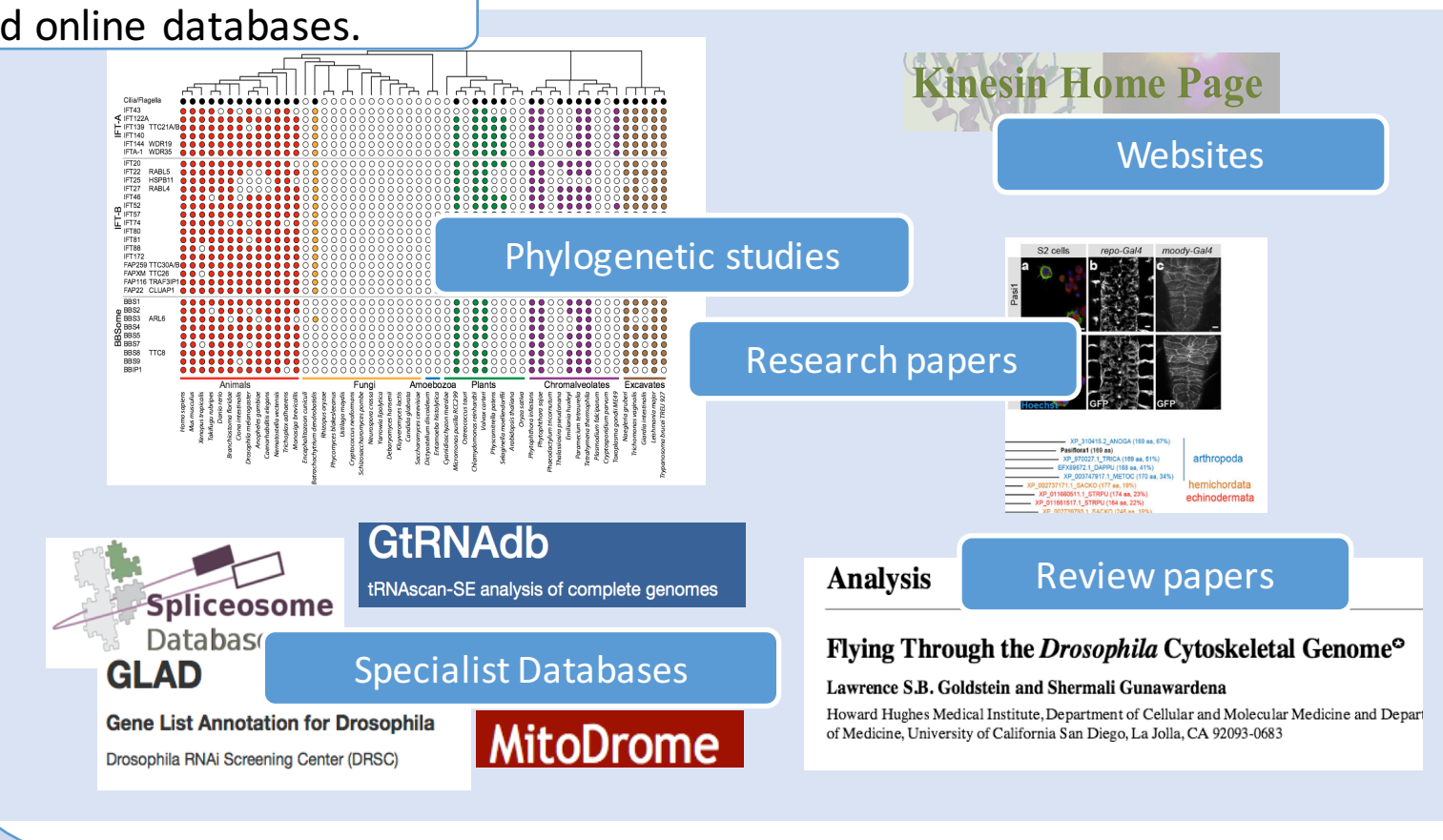
Helen Attrill[1], Giulia Antonazzo[1], Phani Garapati[1], Joshua L. Goodman[2], Steven J. Marygold[1], Alix J. Rey[1], Victor Strelets[2], Jim Thurmond[2] and the FlyBase Consortium

1. Department of Physiology, Development and Neuroscience, University of Cambridge, Cambridge, CB2 3DY, UK. 2. Dept. of Biology, Indiana University, Bloomington, IN 47405, USA. E-mail H.Attrill: hla28@cam.ac.uk

**Introduction:** It is often desirable to search for and view groups of genes whose products are related in some way, such as their known or predicted function. A list of functionally related genes may provide the starting point for a genetic/molecular screen, or be the basis for in silico analyses using associated data (phenotypes, reagents, genomic data etc.), or allow comparison with equivalent gene sets in other species. FlyBase provides two main ways to search for functionally related *Drosophila melanogaster* genes: via Gene Ontology (GO) annotations and our Gene Group resource. Here we show describe how we compile the Gene Groups resource, assign GO annotations to genes and how these two approaches can be used to find functionally-related genes.

## Gene Groups
The FlyBase Gene Group resource is a collection of sets of genes with shared attributes. These groups are primarily focused on well-defined, easily delimited gene sets, such as evolutionary-related gene families (e.g. actins, Wnts), subunits of macromolecular complexes (e.g. ribosome subunits) and sets of genes whose products share a common molecular function (e.g. ubiquitin ligases).

Gene Groups are manually curated based on published research papers, reviews and online databases.



**Gene Group Reports**: Gene Groups are presented in FlyBase report pages. Gene Group reports contain the list of member genes together with additional information organized into sections.



The 'Description' section gives an overview of the criteria used to compile a group. The 'Notes on Group' may contain justification for the inclusion or exclusion of particular genes. The 'Key Gene Ontology (GO) terms' – are those GO terms (or their children) that are associated with most/all of the member genes and are typical of that group.

The 'Members' table, displays all the genes within the group, frequently used synonyms and the source material for group membership of each individual gene within a group. Buttons at the top of this section allow the gene list to be downloaded or exported for analysis.

The 'External Data' section provides links to equivalent groups in other databases, facilitating navigation between different species databases.

All sources used to compile a Gene Group are displayed at the bottom of the page and, for individual genes, in the 'Members' table.

Many areas of biology are comprehensively covered by this resource such as post-translational modification, chromatin organisation, intracellular transport and transmembrane transport. There are 538 gene groups in FlyBase, representing 4487 unique genes. This growing resource now covers over 25% of all of sequence-localized genes and 30% of protein coding genes of the *D. melanogaster* genome.

## Gene Ontology (GO) Annotation
The GO is a widely used controlled vocabulary used to label gene products with biological attributes. The GO is arranged in a hierarchical structure, with more specific child terms nested under higher-level parent terms. For example, 'protein kinase activity' is a child of 'kinase activity'. GO annotations are displayed on individual gene pages in FlyBase.

GO annotation for Ork1 (Open rectifier K+ channel 1, FBgn0017561)



**GO annotation**
Within the gene report, GO annotations are split into the three aspects: **Molecular Function**, **Biological Process** and **Cellular Component**. These are further subdivided into annotations that have been inferred from **experimental** observations and those that have been inferred from **predictions** or **assertions** made by curators or from automated pipelines.

The combinatorial approach of manual and electronic annotation results in good coverage of GO annotation data over the *D. melanogaster* genome: 73% of sequence-localized genes and 88% of protein-coding genes have associated GO terms.

### GO ribbons - a graphical summary of gene function
In the new version of FlyBase, currently available as a beta release, GO summary ribbons are shown at the top of each gene page. The GO summary ribbons use the hierarchical structure of the ontology to group terms under generalized, high-level categories. The detailed annotation can be found in the 'Gene Ontology' section.

The GO summary ribbon for Ork1 (Open rectifier K+ channel 1, FBgn0017561)



Mousing over cell reveals the term(s) represented.

## Finding Gene Groups
Gene Groups can be accessed in three ways:

1. From the FlyBase homepage via QuickSearch



2. Links within individual Gene reports
3. A browsable, hierarchical list of Gene Groups (example below). Groups may be split into nested sub-groups, each of which represents a Gene Group:



Each level of the hierarchy has its own Gene Group report page, with subgroups and their component member genes displayed.



Gene Group members can be downloaded or exported to analysis tools using the buttons at the top of the 'Members' table within a Gene Group report.

## Finding Genes using the GO
The Gene Ontology can be searched in two main ways:

1. From the FlyBase homepage via QuickSearch



2. Or via the Vocabularies link at the top of the homepage.



GO queries via the QuickSearch or Vocabularies tools will take the user to a a Term Report. A Term Report displays information and data associated with a controlled vocabulary term.

Term Report



Clicking this number generate a hit-list of genes annotated with this term or its child terms.

Clicking this number generate a hit-list of genes annotated with this exact term.

GO terms are displayed in the GO hierarchy.

## Combinatorial Searches
It may be useful to combine searches of Gene Groups with GO data, or indeed any other data in FlyBase, to generate more refined gene lists. Here we show two ways in which this can be done.



**HitList Refinement tools**
A list of genes can be downloaded from a Gene Group report page or, when querying a controlled vocabulary, such as the GO, from a Term Report. The tools under 'Results Analysis/Refinement' button can be used to show simple intersections with other data in FlyBase. In the example shown, the hit-list is populated with genes exported from the ION CHANNELS Gene Group. A results analysis of the top 15 individual Biological Process GO terms associated with genes from the ION CHANNELS Gene Group are shown. These can be selected to generate a refined hit-list.



**QueryBuilder**
For more complex queries, the QueryBuilder tool can be used to combine individual query segments using Boolean operators (AND, OR, BUT NOT) in order to generate lists that combine or exclude the given criteria. In the example shown here, a 2-leg query has been performed with genes imported from the ION CHANNELS Gene Group and a second query for GO term search for 'sensory perception'.



## Looking to the future:
- The curation of Gene Groups is often linked with GO annotation curation. When a curator compiles a Gene Group the GO annotation of member genes is quickly reviewed to ensure that it is consistent with the known characteristics of the gene set. There are three main areas into which GO and Gene Group curation will extend in the future:
  - Non-coding RNAs and the components of non-coding RNA biogenesis.
  - Signalling pathways: new pathway and network pages are planned.
  - Protein Complexes with dedicated complex pages linked to physical interaction data.