# New Mines and New Features

Richard Smith
richard@flymine.org

**InterMine**

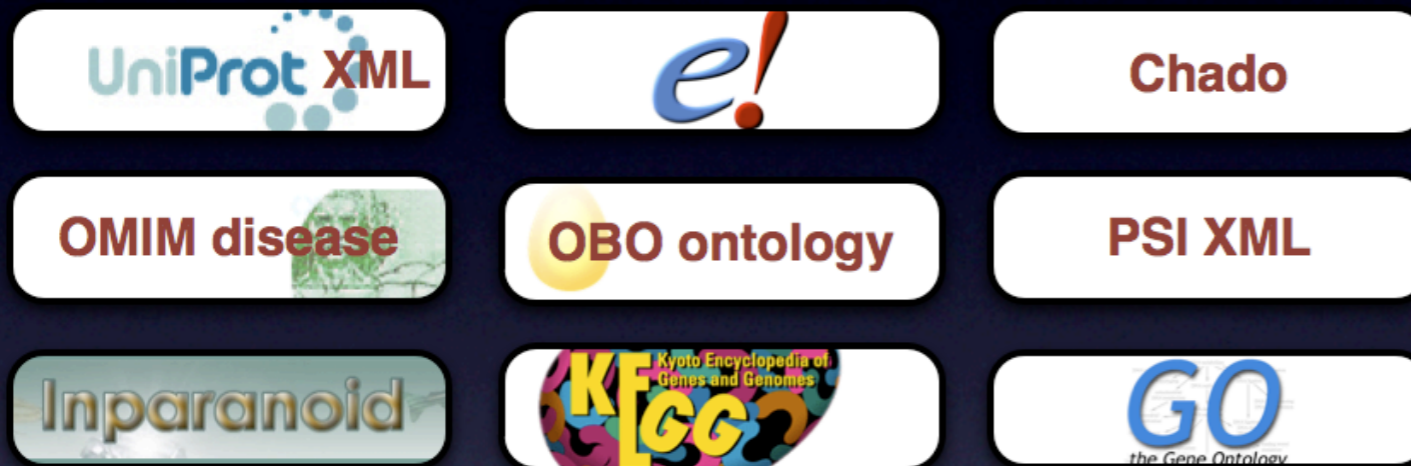**www.intermine.org**

# Overview

- Query-optimised data warehouse system
- Java, object-based data model
- PostgreSQL
- Free, open source (LGPL)

```
Integrate data  →  InterMine data warehouse  →  Web application
                                              →  Query API
                                              →  Web services
```
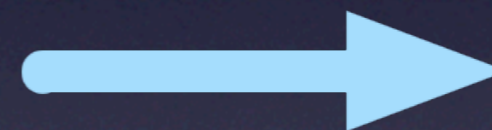
# Data Integration
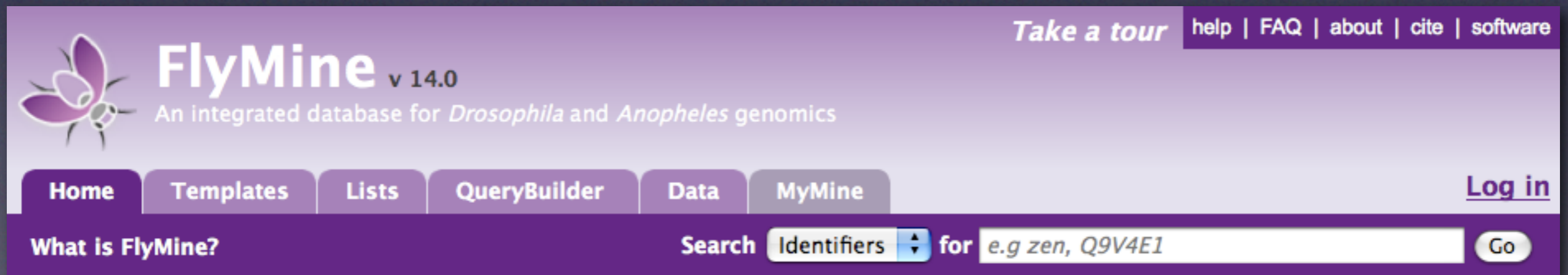
*Existing data sources*



*Configurable data integration*

*Java and Perl APIs*

# Web Application

- Works for any data model

- Advanced functionality for bench biologist

- Highly configurable

- Configuration from within web interface

# Web Interface

- **QueryBuilder -** custom queries, advanced
- **Template queries** - pre-defined queries
- **Report pages** - configurable, templates
- **Lists** - upload, use in queries, analysis widgets
- **Export & API**
- **MyMine**

# Mines for MODs

FlyMine

SGD   Yeast

Zebrafish

RGD   Rat

0.5 FTE each
started at different times

# Mines for MODs

FlyMine

SGD  Yeast

**public beta**

Zebrafish

RGD  Rat

**private beta**

**released**

MGI Mouse

WormBase  Worm

# Mines for MODs

- InterMine alongside existing web site
    - contains MOD data + other sources
    - *working on better embedding of InterMine*
- New features for MODs
- MODs replacing some older search tools
- Lots of input for InterMine development!

# Common Interface



*Your Mine here...*

**Orthologues**
Protein domains
Pathways
Gene Ontology

Images from Wikipedia and RGD

# gene list in FlyMine



141 up-regulated
in heart

**export to RGD**
(orthologues)

40 related to
Cardiovascular
disease

**export to FlyMine**
(orthologues)

- look for relevant
publications
- export sequences

- fetch results dynamically between Mines
- other developers can create mash-ups

# metabolicMine

- Metabolic diseases: *diabetes, obesity*

- Close collaboration with major labs

- Human, mouse, rat

- Lots of data!
  - SNPs, HapMap, expression, ENCODE...
  - Already have parsers for many sources

# metabolicMine

- We've identified a gene, what does it do?
  - currently looking at many sites per gene
- Prioritise candidate genes
- Compare lists, find common attributes
- Upload genome coordinates -> features

# InterMine 0.94

new features

# Search

Templates and QueryBuilder are structured searches

- Lucene search over whole database

- FlyMine:  45 mins,  2.5 GB index file

  - parallel fetching,  precomputes

- Each object is a document

  - attributes are fields

  - also related data  e.g. GO, pathways

# Faceted Search

- BOBO - created by LinkedIn

- Group results by 'facets'

- Filter by multiple facets

- Facet on any property

    e.g. pathway, expression

- *Display is hard*

**Hits by Type**
Allele: 134
TFBindingSite: 93
Stock: 29
CRM: 14
TransposableElementInsertionSite: 10
Intron: 8
Gene: 7
MRNA: 7
Protein: 7
ChromosomeStructureVariation: 4
PointMutation: 3
ChromosomalTranslocation: 2
Exon: 2
CDNAClone: 1
ChromosomalDeletion: 1
ProteinDomain: 1

**Hits by Organism**
D. melanogaster: 296
D. yakuba: 5

Integrates with lists and templates

# Configurable

```
index.references.BioEntity = synonyms organism
index.references.OntologyTerm = synonyms
index.references.Gene = pathways omimDiseases
proteins.proteinDomains goAnnotation.ontologyTerm
index.references.Protein = proteinDomains
index.ignore = Annotation AnalysisResult Comment Homologue
Interaction Image OntologyAnnotation OntologyRelation
OntologyAnnotation ProteinStructure Reporter Sequence
SymmetricalRelation Synonym
index.facet.Category = Category
index.facet.Organism = organism.shortName
index.boost.Gene = 2.0
index.boost.Protein = 1.5
```

# Templates

- Problem: too many similar templates

# Templates

- Select multiple values
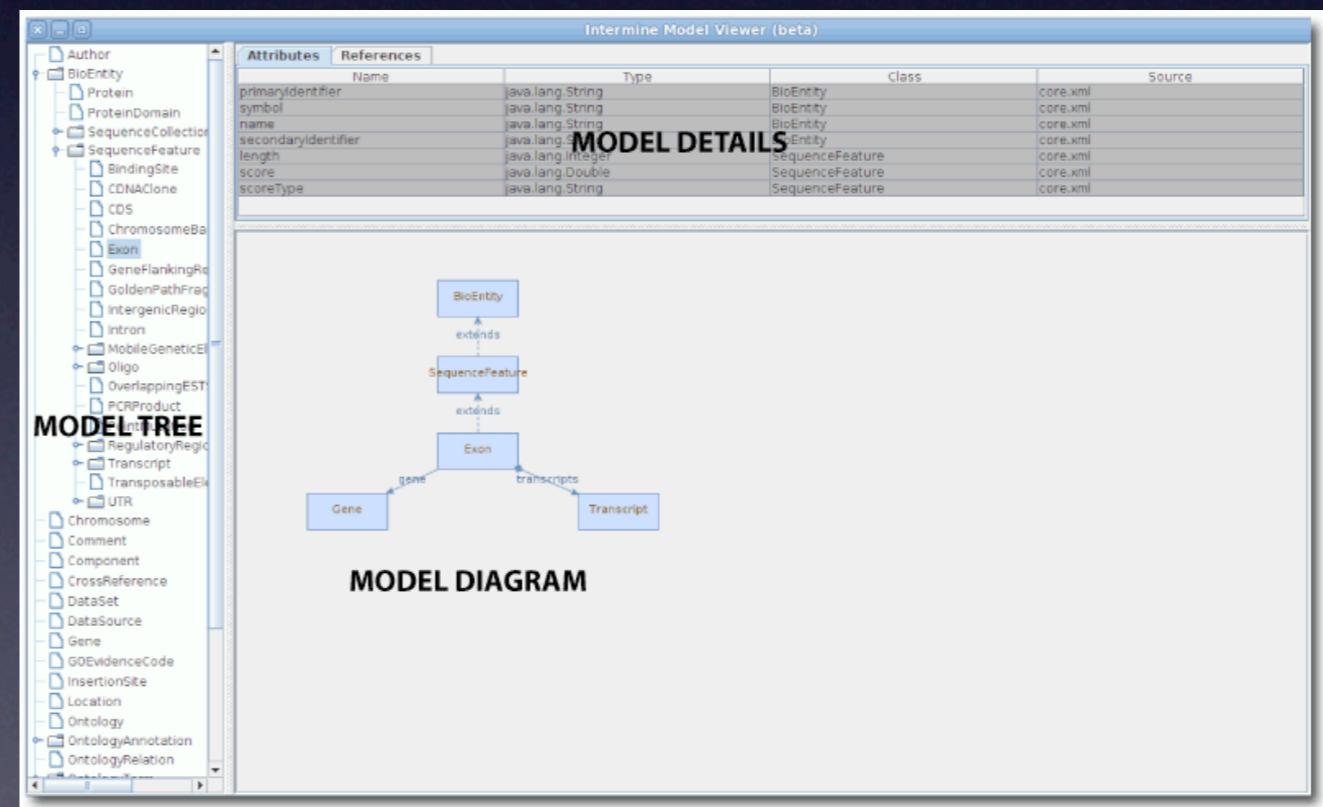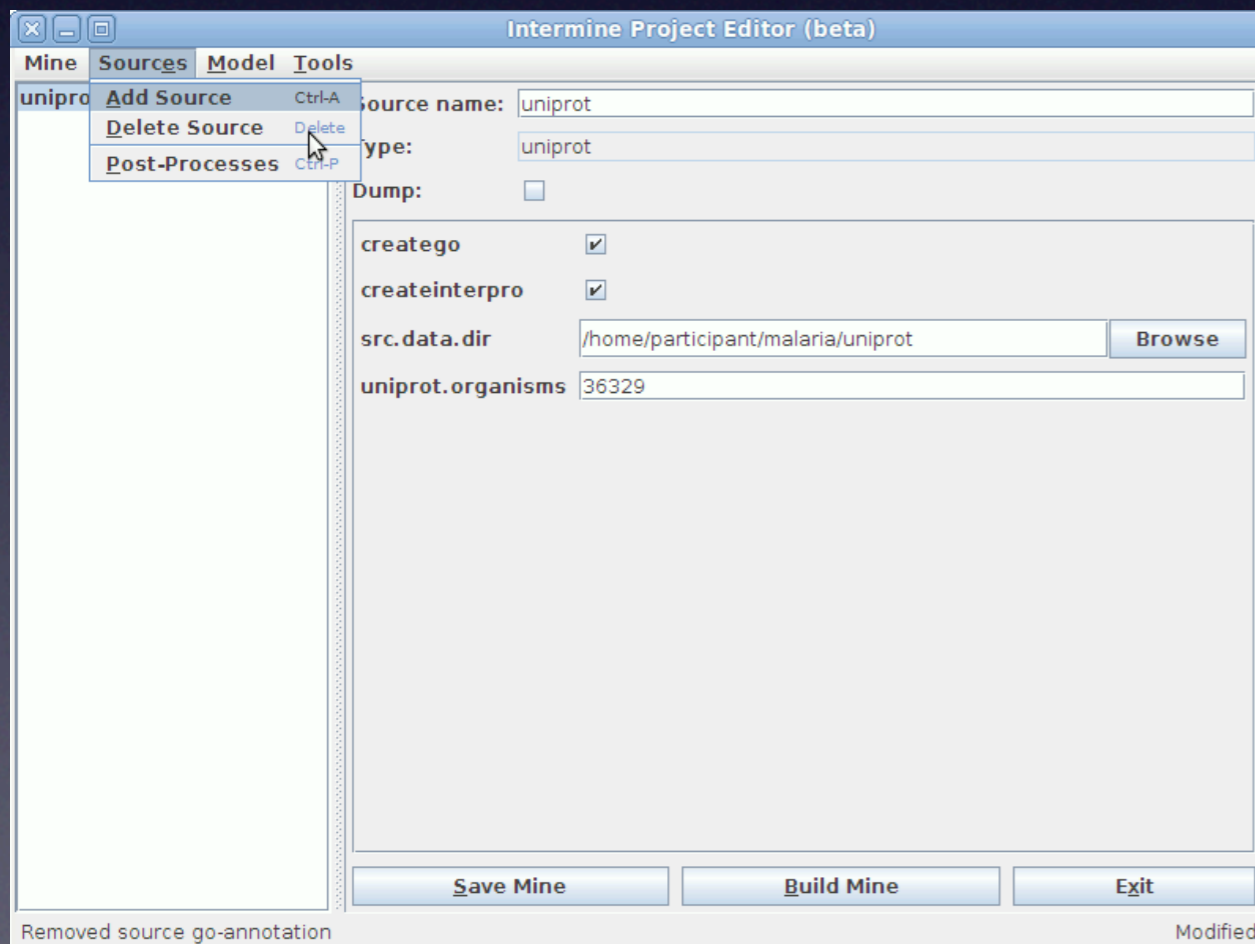
# Performance

- Faster builds:
  - FlyMine:  18 hours to 10 hours  (75GB)
- Better results cache
- More caches!
- Faster export

# MineManager

- GUI to simplify installation

- Simple selection of data sources

# And theres more...

- Galaxy integration

- Improved Perl web service API

- Automatic SO -> model generation

- Upload list of genome regions

- CytoscapeWeb plugin

# Workshop

Come and build your own InterMine

Tomorrow!

We're hiring: calling all
Java/Database/performance gurus.

# InterMine Team

**Biologists** Rachel Lyne, Adrian Carr, Mike Lyne

**Developers** Richard Smith, Sergio Contrino, Julie Sullivan, Matthew Wakeling, Alex Kalderimis, Fengyuan Hu, Daniela Butano

**Students** Nils Kolling, Richard B.

**Sys Admin** Dan Tomlinson

**PI** Gos Micklem

**www.intermine.org**

**richard@flymine.org**