

InterMine

Open source data warehouse and web interface

Richard Smith
University of Cambridge

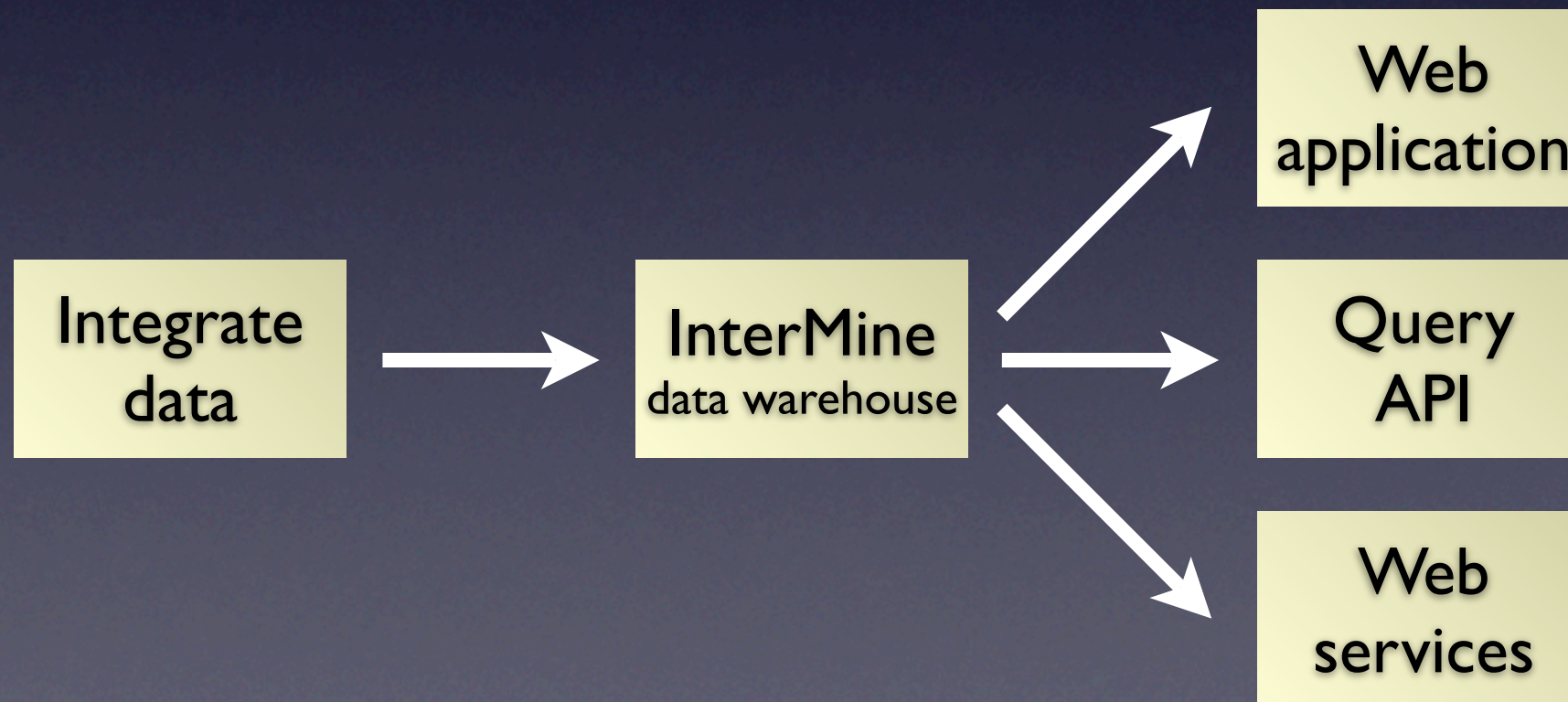
Demo: *Monday 12:15 room 701B*

Poster: *E34 (Monday)*

www.intermine.org

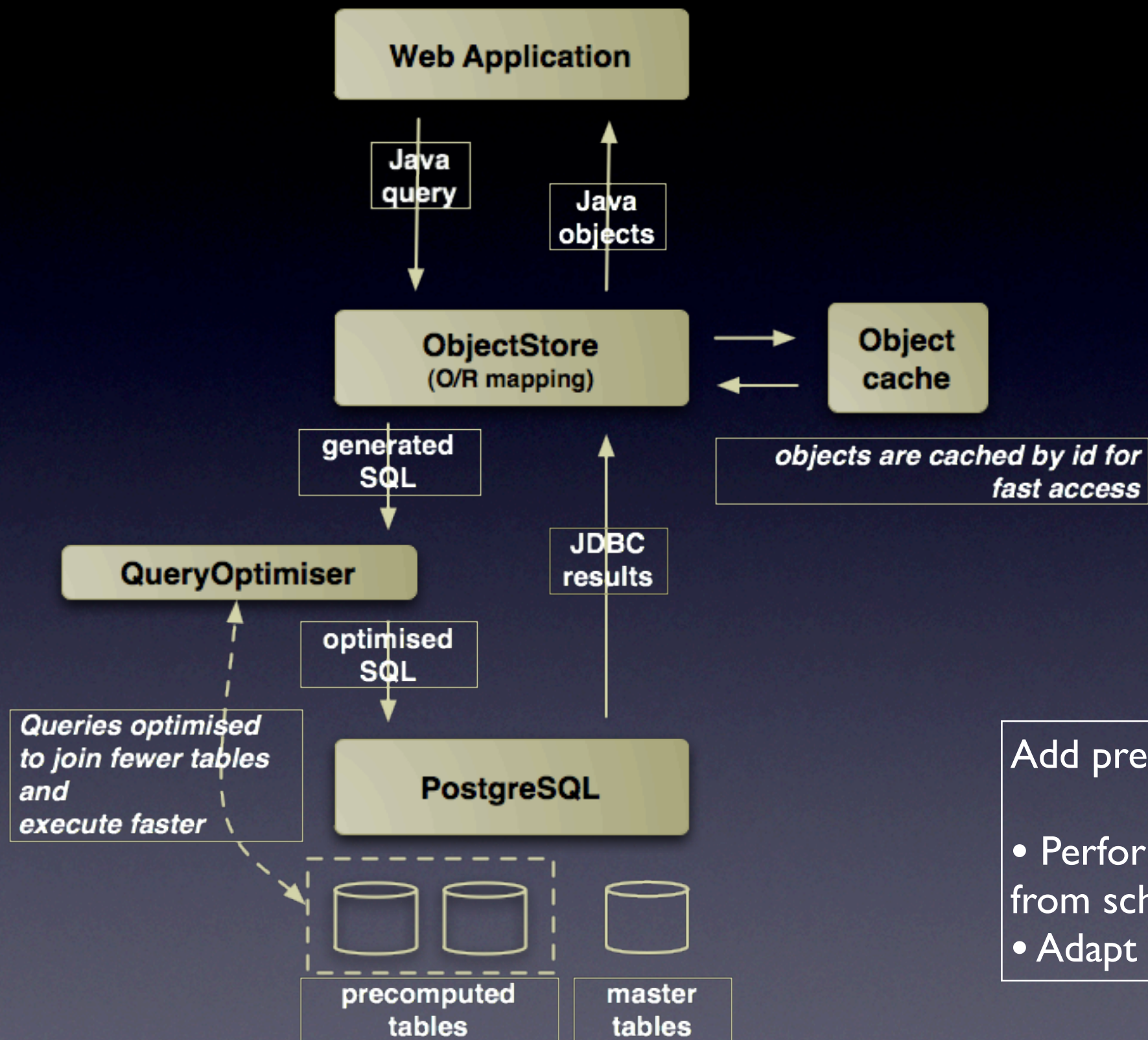
Overview

- Query-optimised data warehouse system
- Java, object-based data model
- Flexible querying



Projects

- Five developers, since 2002
- FlyMine - www.flymine.org
 - 30+ data sources, *Drosophila* & *Anopheles*
- modENCODE - www.modencode.org
 - *C. elegans*/*D. melanogaster* high throughput
- BOKU & IMP - Vienna
- MitoMiner - mitochondria
- MilkMine - milk proteins
- *Yeast, Rat, Zebrafish*

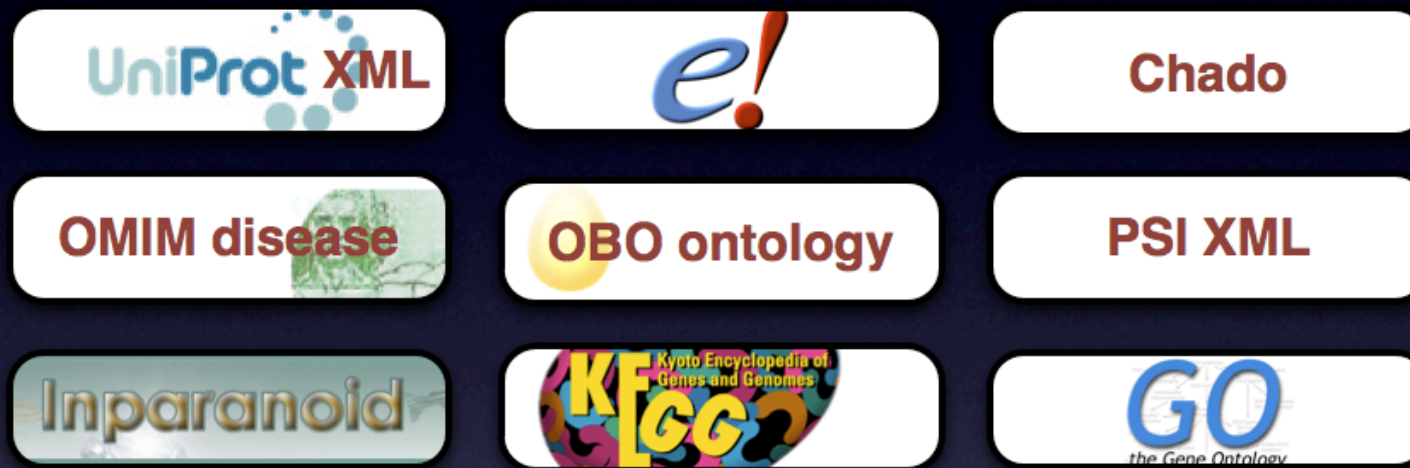


Add precomputed tables at any time:

- Performance optimisation separate from schema design
- Adapt performance to actual use

Data Integration

Existing data sources



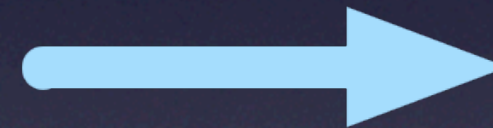
+

Custom Data Sources



Java and Perl APIs

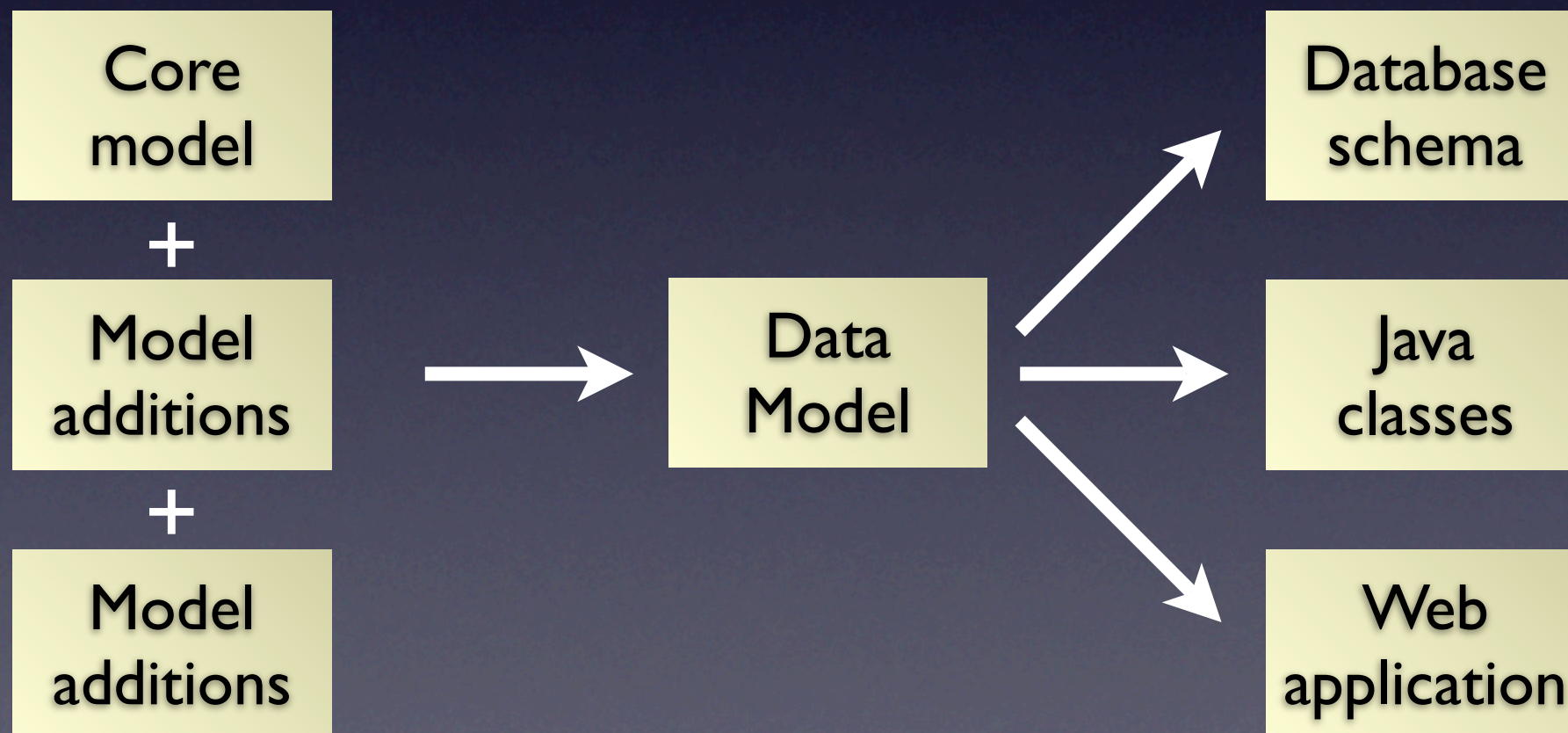
Configurable data integration



**InterMine
data warehouse**

Auto-generation

- Object model defined by XML file
- Low overhead to extending data model



Custom data

- Any source can add to the data model:

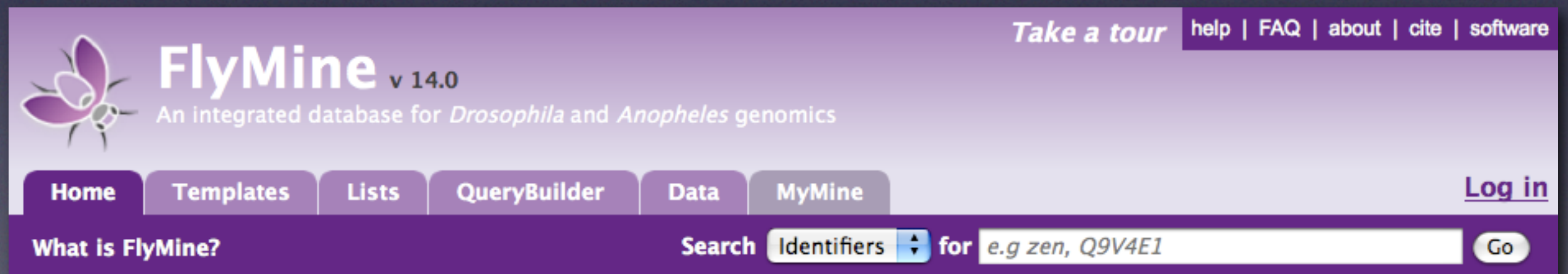
```
<class name="Pathway" is-interface="true" >  
  <attribute name="name" type="String"/>  
  <collection name="genes" referenced-type="Gene"  
    reverse-reference="pathways"/>  
</class>  
<class name="Gene" is-interface="true">  
  <collection name="pathways" referenced-type="Pathway"  
    reverse-reference="genes"/>  
</class>
```


Configure a new Mine

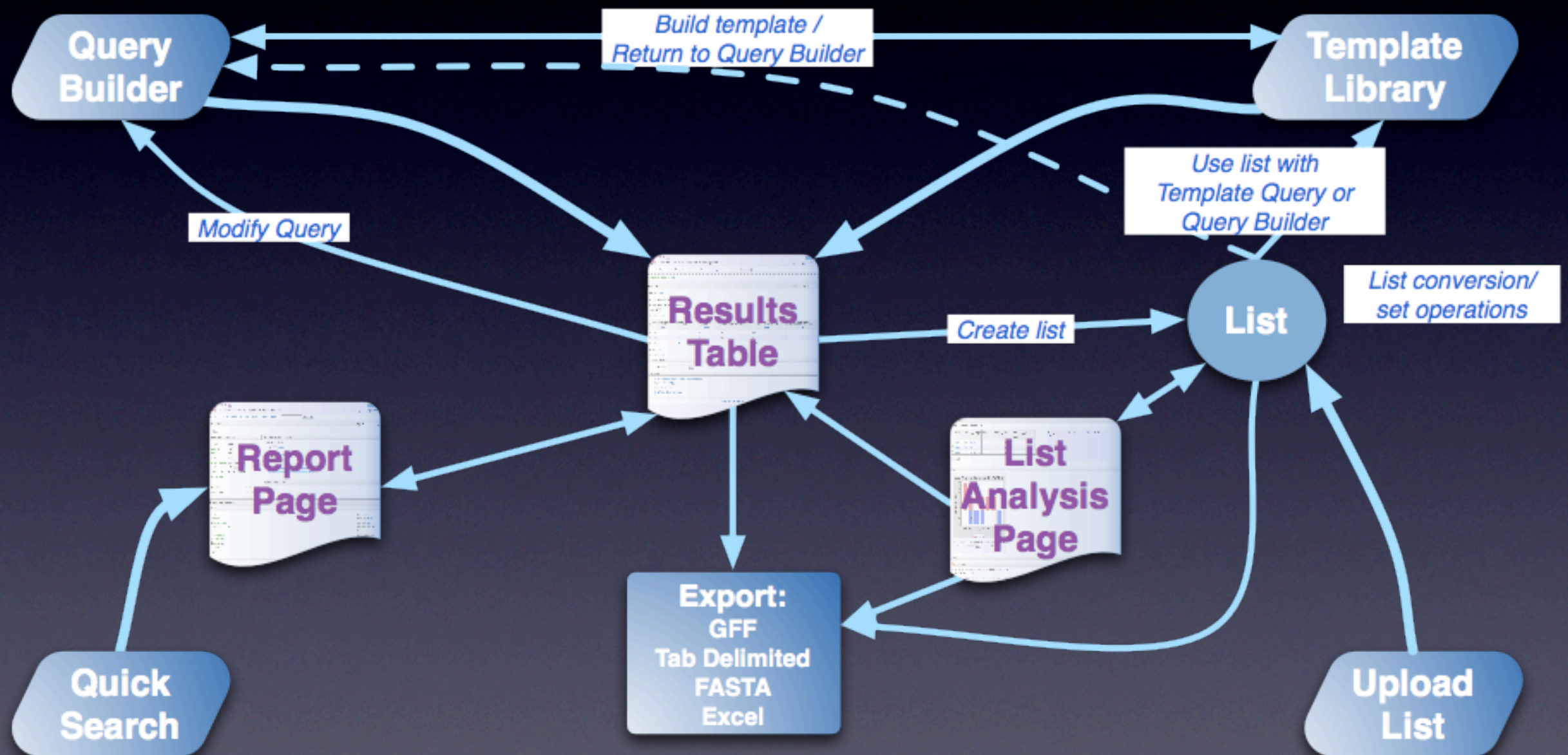
```
<project type="bio">
...
<sources>
  <source name="uniprot" type="uniprot" dump="true">
    <property name="uniprot.organisms"
      value="7227 6239" />
    <property name="src.data.dir"
      location="/data/uniprot" />
  </source>
  <source name="my-source" type="pathways">
    <property name="src.data.dir"
      location="/data/pathways" />
  </source>
  ...
</sources>
</project>
```


Web Application

- Works for any data model
- Advanced functionality for bench biologist
- Highly configurable
- Configuration from within web interface



Webapp Overview



QueryBuilder

Model browser ?

Browse through the classes and attributes. Click on [SUMMARY +](#) links to add summary of fields to the results table or on [SHOW +](#) links to add individual fields to the results. Use [CONSTRAIN +](#) links to constrain a value in the query.

- TFBindingSite ?
 - curated Boolean [SHOW +](#) [CONSTRAIN +](#)
 - identifier [SHOW +](#) [CONSTRAIN +](#)
 - length [SHOW +](#) [CONSTRAIN +](#)
 - name [SHOW +](#) [CONSTRAIN +](#)
 - [-] annotations Annotation collection [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] chromosome Chromosome ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] chromosomeLocation Location ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] evidence Evidence ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] factor Gene ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] gene Gene ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] motif Motif ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] objects Relation ? collection [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] organism Organism ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] overlappingFeatures LocatedSequenceFeature ? collection [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] relations SymmetricalRelation ? collection [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] sequence Sequence ? [SUMMARY +](#) [CONSTRAIN +](#)
 - [-] subjects Relation collection [SHOW +](#) [CONSTRAIN +](#)
 - [-] synonyms Synonym ? collection [SUMMARY +](#) [CONSTRAIN +](#)

Constraints on the current query ?

Click on a class name below to view its fields

- TFBindingSite ?
- factor Gene ?
- IN demo bag
- chromosomeLocation Location
- object Chromosome
- evidence DataSet collection

Constraint logic ?

one constraint

Fields selected for output ?

Columns to Display

Use the [SHOW +](#) or [SUMMARY +](#) links to add fields to the results table. Click and drag the blue output boxes to choose the output column order.

TFBindingSite > factor > identifier

TFBindingSite > chromosomeLocation > start

TFBindingSite > chromosomeLocation > end

TFBindingSite > gene > identifier

TFBindingSite > gene > symbol

Sort Results By Column

To sort the results by a specific field, click on [SORT +](#) in that field's blue box. Use the button in the purple box below to reverse the direction of the sort. Click [A-Z](#) to sort in ascending order. Click [Z-A](#) to sort the results in descending order.

TFBindingSite > factor > identifier

Query summary

Model
browser

Constraint
editor

Order output columns

Set sort order

Template Queries

Protein domain ➡ proteins from a specific organism

For a particular domain (or list of domains) give the proteins (from a specific organism) which have these domains.

[1] Search for proteins in the following organism:

Organism name: =

[2] Which contain the domain(s):

ProteinDomain name: =

- Homeobox
- Homeodomain Cdx
- Homeodomain-like
- Homeodomain-related
- Homeobox, Hox9
- Homeobox Pitx/unc30
- Abl-interactor, homeo-domain homologous region
- Homeodomain engrailed related
- Homeodomain Lbx related
- Homeodomain protein CUT

Show Results

Edit Query

NEW: [Embed this query.](#) [Help](#)

You are not logged in. [Log in](#) to mark item

XML

Template Library

[Home](#)[Templates](#)[Lists](#)[QueryBuilder](#)[Data](#)[MyMine](#)[Log in](#)Search for [FlyMine](#) > [Templates](#)

Templates

Templates are predefined queries, each has a simple form and a description. Click on a template to run it, you can search for templates by keyword and filter them by category.

Search: Filter: Actions: ☒ [Show descriptions](#)

You are not logged in. [Log in](#) to mark items as favourites

☐ [Gene --> Chromosomal location.](#)

Show the **chromosome** and the **chromosome** location of a particular **gene**

☐ [Chromosome --> All genes.](#)

Show **genes** located on a particular **chromosome**

☐ [All genes in organism --> All chromosomal locations.](#)

Show the **chromosomal** location and sequence for all **genes** from a particular organism

☐ [Chromosomal location --> All genes.](#)

Show the **genes** located between two points on a **chromosome**. (Data Source: FlyBase, Ensembl)

Results

Results for template: **Protein domain --> proteins from a specific organism**

For a particular domain (or list of domains) give the proteins (from a specific organism) which have these domains.

CREATE LIST
 ADD TO LIST
 EXPORT
 //
 PAGE SIZE
 << FIRST < PREVIOUS | NEXT > LAST >>

<input type="checkbox"/>	Protein > primaryIdentifier	Protein > primaryAccession	<input type="checkbox"/>	Gene > primaryIdentifier	Gene > symbol	Protein domain > name	<input type="checkbox"/>	Protein domain > primaryIdentifier	Protein domain > type
<input type="checkbox"/>	Q9Y1P6_DROME	Q9Y1P6	<input type="checkbox"/>	FBgn0027364	Six4	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9Y0Z9_DROME	Q9Y0Z9	<input type="checkbox"/>	FBgn0003896	tup	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9XZU0_DROME	Q9XZU0	<input type="checkbox"/>	FBgn0023489	Pph13	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9XZC2_DROME	Q9XZC2	<input type="checkbox"/>	FBgn0019650	toy	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9W4B3_DROME	Q9W4B3	<input type="checkbox"/>	FBgn0053980	CG33980	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9W4B2_DROME	Q9W4B2	<input type="checkbox"/>	FBgn0029775	CG4136	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9W423_DROME	Q9W423	<input type="checkbox"/>	FBgn0040918	Lag1	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9W3C6_DROME	Q9W3C6	<input type="checkbox"/>	FBgn0030058	CG11294	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	RX_DROME	Q9W2Q1	<input type="checkbox"/>	FBgn0020617	Rx	Homeobox	<input type="checkbox"/>	IPR001356	Domain
<input type="checkbox"/>	Q9W2P8_DROME	Q9W2P8	<input type="checkbox"/>	FBgn0008636	hbn	Homeobox	<input type="checkbox"/>	IPR001356	Domain

Selected:

<< First < Previous | Next > Last >> | Displaying rows 1 to 10 | Total rows: 190

Lists

- Lists of any type - e.g. *genes, protein domains, organisms*
- Use a list in any query
- Save lists from results pages
- Upload identifiers

Gene [D. melanogaster] -> KEGG Pathway.
Show the KEGG pathway identifier and name for the selected gene.

Gene: ?

or ☐ constrain to be list

NEW: [Embed this query.](#) [Help](#)
You are not logged in. [Log in](#) to mark items

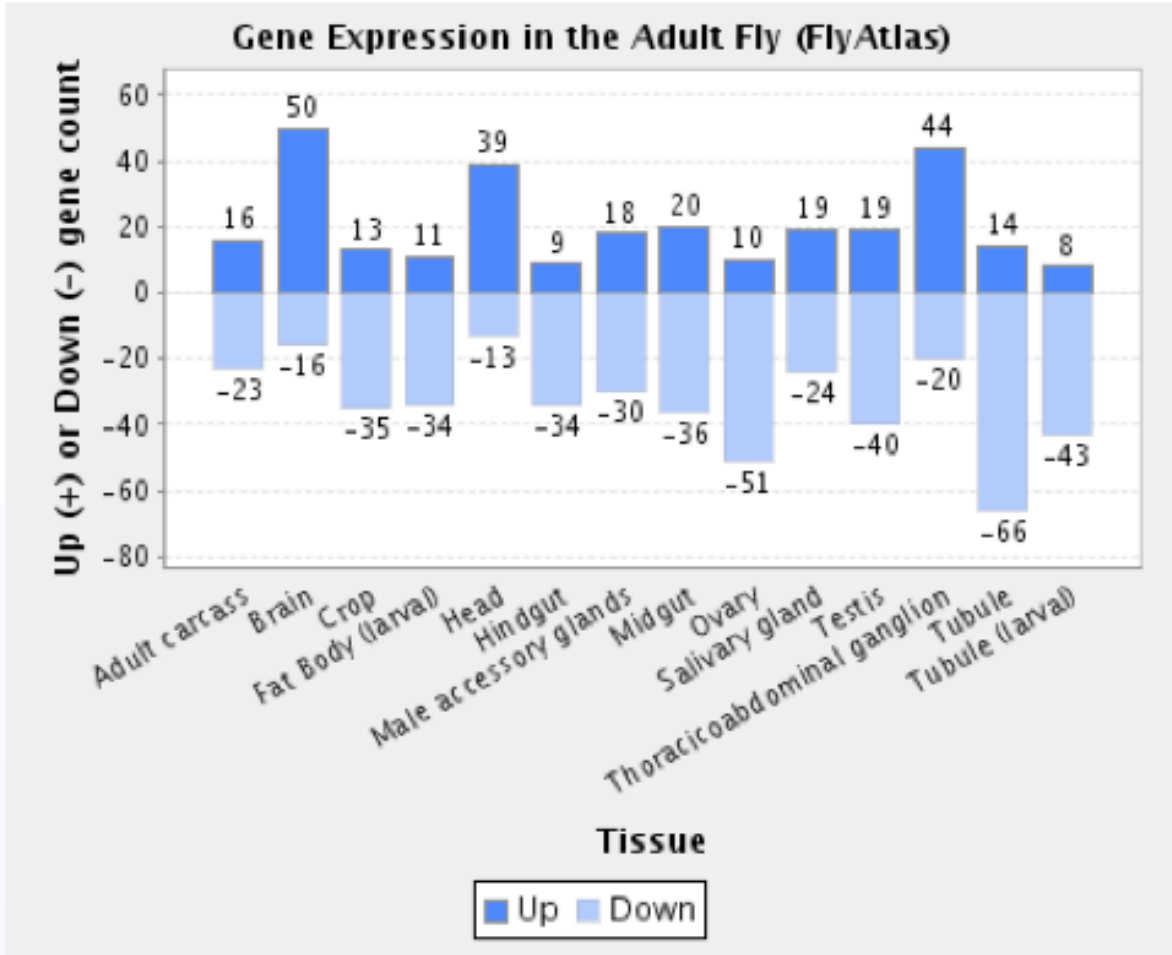
- PL FlyAtlas_brain_top
- PL FlyAtlas_head_top
- PL FlyAtlas_hindgut_top
- PL FlyAtlas_maleglands_top
- PL FlyAtlas_midgut_top
- PL FlyAtlas_ovary_top
- PL FlyAtlas_testis_top
- PL FlyAtlas_tubule_top

Widgets displaying properties of 'homeobox genes'

Click to select widgets you would like to display: [Chromosome Distribution](#) | [Gene Expression in the Adult Fly \(FlyAtlas\)](#) | [mRNA subcellular localisation \(fly-FISH\)](#) | [BDGP expression patterns](#) | [Pathway Information \(KEGG\)](#) | [Orthologues](#) | [Genetic Interactions](#) | [Gene Ontology Enrichment](#) | [Protein Domain Enrichment](#) | [Publication Enrichment](#) | [BDGP Enrichment](#)

Gene Expression in the Adult Fly (FlyAtlas)

For each tissue in the adult fly, the number of genes from this list for which the levels of expression are significantly high (Up) or low (Down) according to [FlyAtlas AffyCall](#). Number of Genes in this list not analysed in this widget: 10



Gene Ontology Enrichment

GO terms enriched for items in this list. Smaller p-values show greater enrichment. Method: [Hypergeometric test](#). Number of Genes in this list not analysed in this widget: 1

Options

Multiple Hypothesis Test Correction: [Benjamini and Hochberg](#)




Ontology: [biological_process](#) Maximum value to display: [0.01](#)

[Display](#) [Export](#)

<input type="checkbox"/>	GO Term	p-Value	
<input type="checkbox"/>	regulation of transcription, DNA-dependent [GO:0006355]	4.7508E-121	104
<input type="checkbox"/>	regulation of RNA metabolic process [GO:0051252]	8.6166E-116	104
<input type="checkbox"/>	RNA biosynthetic process [GO:0032774]	2.6877E-115	104
<input type="checkbox"/>	transcription, DNA-dependent [GO:0006351]	4.0436E-115	104
<input type="checkbox"/>	regulation of transcription [GO:0045449]	2.1217E-112	104
<input type="checkbox"/>	regulation of macromolecule biosynthetic process [GO:0010556]	5.2157E-108	104
<input type="checkbox"/>	regulation of biosynthetic process [GO:0009889]	1.3931E-107	104
<input type="checkbox"/>	regulation of nucleobase, nucleoside, nucleotide and nucleic acid metabolic process [GO:0019219]	1.6359E-107	104
<input type="checkbox"/>	transcription [GO:0006350]	9.474E-107	104

Superuser

- Non-programmer can configure
- Public template queries
- Public lists
- Templates on report pages
- Tagging

	<div>Chromosome_Gene </div> <div> aspect:Genomics [x] im:public [x]</div> <div><input type="text" value="im:frontpage"/> <input type="button" value="Add"/> <input type="button" value="Done"/></div>	<div>Chromosome --></div> <div>All genes.</div>
---	---	--

Query API

```
PathQuery q = new PathQuery(model);  
  
q.setView("Protein.primaryIdentifier",  
          Protein.genes.primaryIdentifier");  
  
q.addConstraint("Protein.proteinDomains.name",  
                Constraints.eq("Homeobox"));  
q.addConstraint("Protein.organism.shortName",  
                Constraints.eq("D. melanogaster"));
```


Web Services

- RESTful web service
- Run queries or templates
 - export XML from web app or use query API
- Java client API
- *Perl client API coming soon*
- *Lists, widgets, logins, tags to be added*

Embedding Templates

- Web service can return HTML

Your web page
e.g. *gene report*

Trail: Query > Results

Results for template: Protein domain --> proteins from a specific organism
For a particular domain (or list of domains) give the proteins (from a specific organism) which have these domains.

CREATE LIST | GO TO LIST | IMPORT // PAGE SIZE: 10 | << FIRST | PREVIOUS | NEXT | LAST >>

Protein > primaryAccession	Protein > primaryIdentifier	Gene > primaryIdentifier	Protein domain > name	Protein domain > primaryIdentifier	Protein domain > type
<input type="checkbox"/> A0AVV3	A0AVV3_DROME	<input type="checkbox"/> FBgn0085369	Paired-like homeodomain protein, OAR	<input type="checkbox"/> IPR003654	Domain
<input type="checkbox"/> A2RVG7	A2RVG7_DROME	<input type="checkbox"/> FBgn0085396	Paired-like homeodomain protein, OAR	<input type="checkbox"/> IPR003654	Domain
<input type="checkbox"/> A8DYQ2	A8DYQ2_DROME	<input type="checkbox"/> FBgn0085369	Paired-like homeodomain protein, OAR	<input type="checkbox"/> IPR003654	Domain
<input type="checkbox"/> A8DYR8	A8DYR8_DROME	<input type="checkbox"/> FBgn0085396	Paired-like homeodomain protein, OAR	<input type="checkbox"/> IPR003654	Domain

Call template
with parameters

InterMine
web
service

HTML

- 'Embed this template' link
- Saves remote site from integrating data
- *Widgets coming soon*

Acknowledgments

Biologists Hilde Jannsens, Rachel Lyne

Developers Richard Smith, Jakub Kulaviak, Julie Sullivan, Matthew Wakeling, Xavier Watkins

Sys Admin Dan Tomlinson

modENCODE Sergio Contrino, Kim Rutherford

PI Gos Micklem

www.intermine.org