

# InterMine Data Warehouse

Gos Micklem



# InterMine used to build FlyMine

## An integrated database of Drosophila and Anopheles data

**Main**

- Home
- Mailing lists
- Links
- Contact us

**Release info**

- Data sources
- Release notes
- Archived releases

**News**

- Server move complete (Thur, 4th Jan 2007 11:52 GMT)
- Service interruption (Tue, 2nd Jan 2007 17:02 GMT)
- FlyMine 6.0 Released (Fri, 24th Nov 2006 17:00 GMT)
- All news

**Documentation**

- Get Started
- Tutorials
- User Manual
- Presentations
- Link to FlyMine
- Cite FlyMine

**Download software**

- SVN checkout
- Browse source

**Project**

- About
- Team
- Funding
- Jobs
- Contact

**Related Projects**

- InterMine
- StemCellMine

**Internal**

Done

**Introduction**

FlyMine is an integrated database of genomic, expression and protein data for *Drosophila*, *Anopheles* and *C. elegans*. Integrating data makes it possible to run sophisticated data mining queries that span domains of biological knowledge.

[What can I do with FlyMine?](#)

FlyMine is now updated to version 6.0. See [release notes](#) for information on new content and an important note about saved bags of objects.



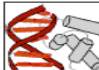
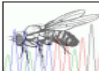







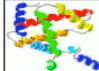

**Quick Search**

Search for genes, proteins or any other annotation. [Help.](#)

Enter an identifier or synonym (e.g. 'zen', 'Q9V4E1'):

**Aspects**

Choose an aspect to explore the different types of data in FlyMine and start queries.

 <p><a href="#">Genomics</a> Genome annotation</p>	 <p><a href="#">Gene Expression</a> ArrayExpress</p>	 <p><a href="#">Transcriptional Regulation</a> Regulatory regions and transcription factor binding sites</p>
 <p><a href="#">DrosDel</a> P-element insertions and deletions</p>	 <p><a href="#">Tiling Path</a> Microarray Tiling Primers</p>	 <p><a href="#">INDAC</a> Long oligos from the International <i>Drosophila</i> Array Consortium</p>
 <p><a href="#">RNAi</a> RNA interference</p>	 <p><a href="#">Disease</a> Human disease matches from Homophila</p>	 <p><a href="#">Comparative Genomics</a> Orthologues and paralogues</p>
 <p><a href="#">Proteins</a> Protein and proteomics data</p>	 <p><a href="#">Protein Interactions</a> IntAct</p>	 <p><a href="#">Protein Structure</a> 3-D protein structures</p>
 <p><a href="#">Gene Ontology</a></p>		

**Template search**

Search:

Enter a keyword to find pre-defined template queries relating to a certain type of data. Select 'Public templates' to search pre-defined templates, 'My templates' to search templates you have created yourself or Everything to search all templates.

[View all templates...](#)

**FlyMine QueryBuilder**

Advanced users can use a flexible query interface to construct their own data mining queries. See [user manual](#) and [tutorials](#).

[List all classes...](#) [Browse model...](#) [Import a query from XML...](#) [Help.](#)



# Project Stats

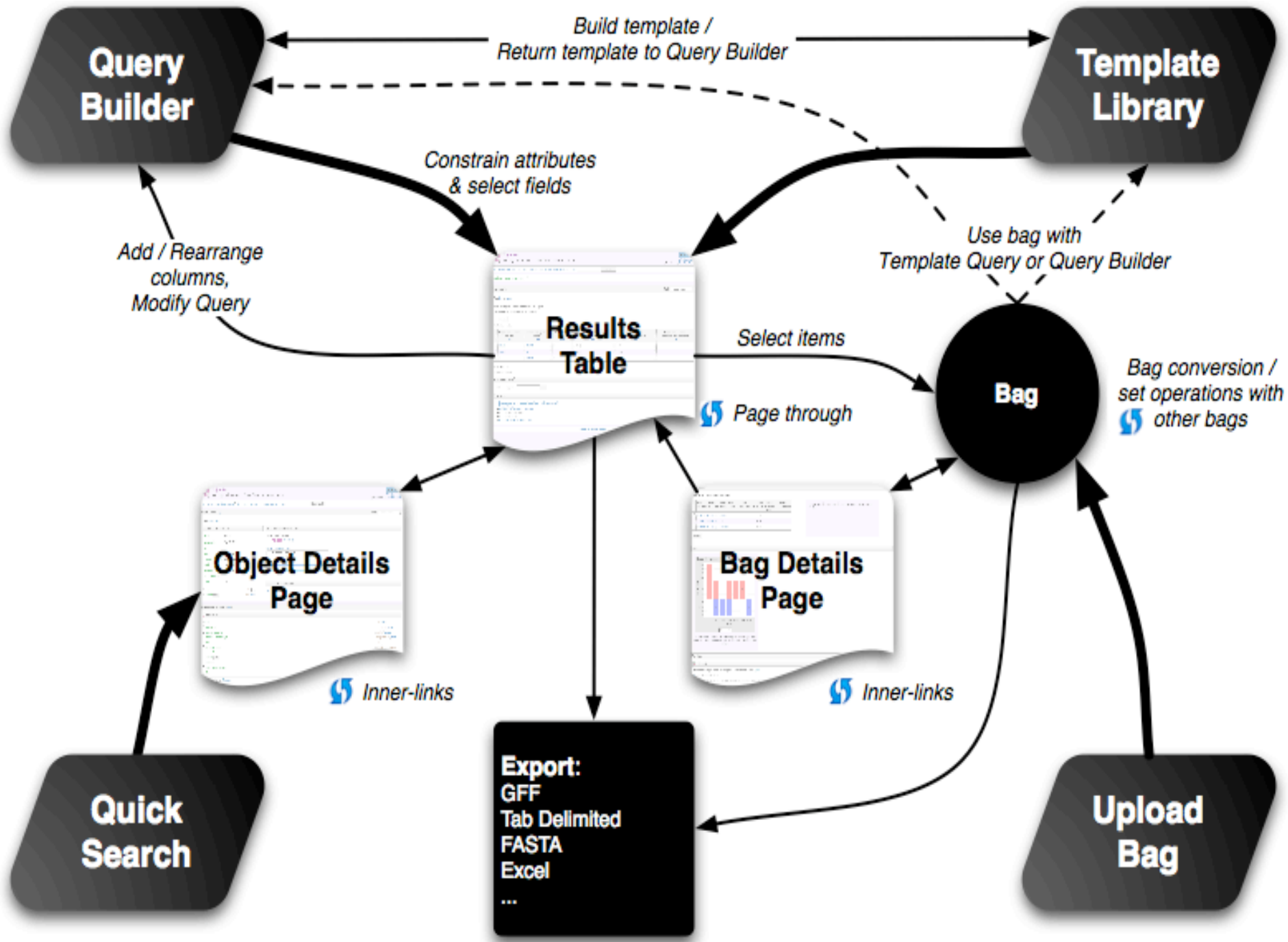
- \_ Team of 7 FTE
  - \_ 5 developers, one sys admin,
  - \_ 1 biologist/ bioinformatician
- \_ Java/ postgresSQL
- \_ SVN: 125,000 lines of code + 57,000 lines of tests
- \_ Started in 2002
- \_ In use by others in Cambridge, Edinburgh, Vienna... + modENCODE DCC if funded
- \_ [modENCODE/ Chado](#)



# InterMine Guiding Principles

- Open Source
- Generic
- Maximum use of automatic code generation/ testing
- Integrate diverse types of data - standard formats
- Load order independent (fine-grained source priorities)
- Query re-writing: uncouple data model and tuning
- Allow performance tuning of the live database
- Allow domain-specific configuration of the user interface
  - By domain experts not developers
  - Release independently of database
- Scope: small project databases to large scale data warehouses



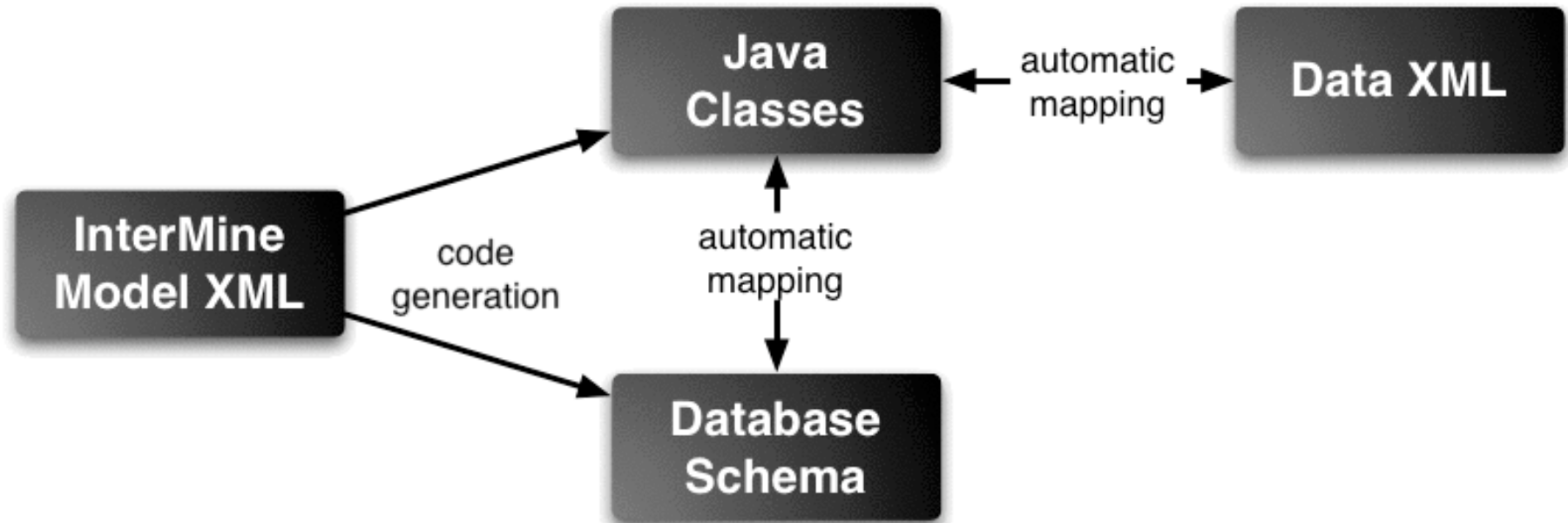


# Technical Overview

- \_ Custom Java object/relational system
  - \_ When we started, couldn't select from multiple classes at one time using hibernate.
- \_ Optimised for read-only performance
- \_ Designed for big, complex queries
- \_ Performance optimisation
  - \_ Transparent query re-writing
- \_ Web application - Struts/JSP/Ajax



# Code Generation



We use SOFA as core data model - similar to Chado. Added Gene.description, compiled, loaded data (here XML + FASTA), released webapp.



# Loading Data

- \_ Read-only in production environment
- \_ Load data from InterMine XML
- \_ Parsers from standard formats
  - \_ e.g. UniProt, GFF3, PSI, FASTA
- \_ Powerful integration system





# Example InterMine XML

```
<items>
  <item id="0_3" class=""
implements="http://www.flymine.org/model/genomic#Gene">
  <attribute name="identifier" value="xfile" />
  <attribute name="description" value="A test gene for GMOD
meeting" />
  <reference name="organism" ref_id="0_1" />
  <collection name="transcripts">
    <reference ref_id="0_9" />
  </collection>
</item>
<item id="0_1" class=""
implements="http://www.flymine.org/model/genomic#Organism">
  <attribute name="taxonId" value="7227" />
</item>
...
```



Trail: Query &gt; Gene

▼ Summary for selected Gene

**identifier** <sup>?</sup> **xfile**  
**organismDbId** <sup>?</sup>  
**symbol**  
**name** <sup>?</sup>  
**chromosomeLocation** **fake\_chromosome : 13691-14720**  
**length** **1030** [FASTA...](#)  
**organism.name** **Drosophila melanogaster**  
**wildTypeFunction** <sup>?</sup>  
 description **A test gene for GMOD meeting**

**organism** **1** Organism [\[details...\]](#)  
 **proteins** **0** Protein

Other details by aspect [help...]

▼ Genomics

**CDSs** **0** CDS  
 **chromosome** **1** Chromosome [\[details...\]](#)  
 **chromosomeLocation** **1** Location [\[details...\]](#)

Class	object.identifier	subject.identifier	start	end	
Location	fake_chromosome	xfile	13691	14720	<a href="#">[details...]</a>

[\[show in table...\]](#)

**downstreamIntergenicRegion** **1** IntergenicRegion  
 **exons** **2** Exon

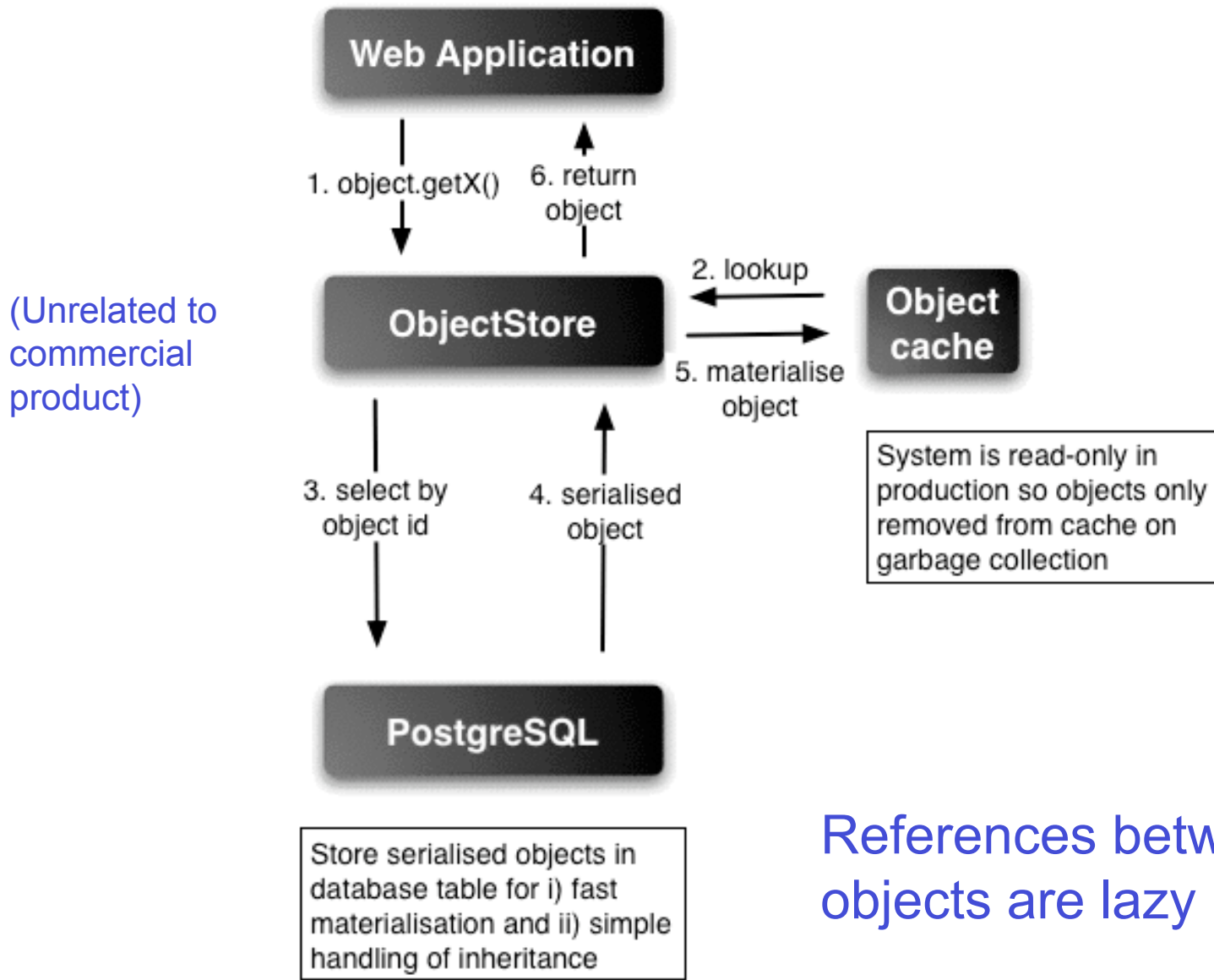
Class	identifier	symbol	chromosomeLocation	length	organism.name	
Exon	xfile-exon-2		<b>fake_chromosome : 14687-14720</b>	<b>34</b> <a href="#">FASTA...</a>	Drosophila melanogaster	<a href="#">[details...]</a>
Exon	xfile-exon-1		<b>fake_chromosome : 13691-13767</b>	<b>77</b> <a href="#">FASTA...</a>	Drosophila melanogaster	<a href="#">[details...]</a>

[\[show in table...\]](#)

**overlappingFeatures** **0** LocatedSequenceFeature



# Retrieving Objects



References between objects are lazy



Trail: Query &gt; Gene

▼ Summary for selected Gene

**identifier?** xfile  
**organismDbId?**  
**symbol**  
**name?**  
**chromosomeLocation** fake\_chromosome : 13691-14720  
**length** 103 FASTA...  
**organism.name** Drosophila melanogaster  
**wildTypeFunction?**  
**description** A test gene for GMOD meeting

**organism** 1 Organism [details...]  
 **proteins** 0 Protein

Other details by aspect [help...]

▼ Genomics

**CDSs** 0 CDS  
 **chromosome** 1 Chromosome [details...]  
 **chromosomeLocation** 1 Location [details...]

Class	object.identifier	subject.identifier	start	end	
Location	fake_chromosome	xfile	13691	14720	[details...]

[show in table...]

**downstreamIntergenicRegion** 1 IntergenicRegion

**exons** 2 Exon

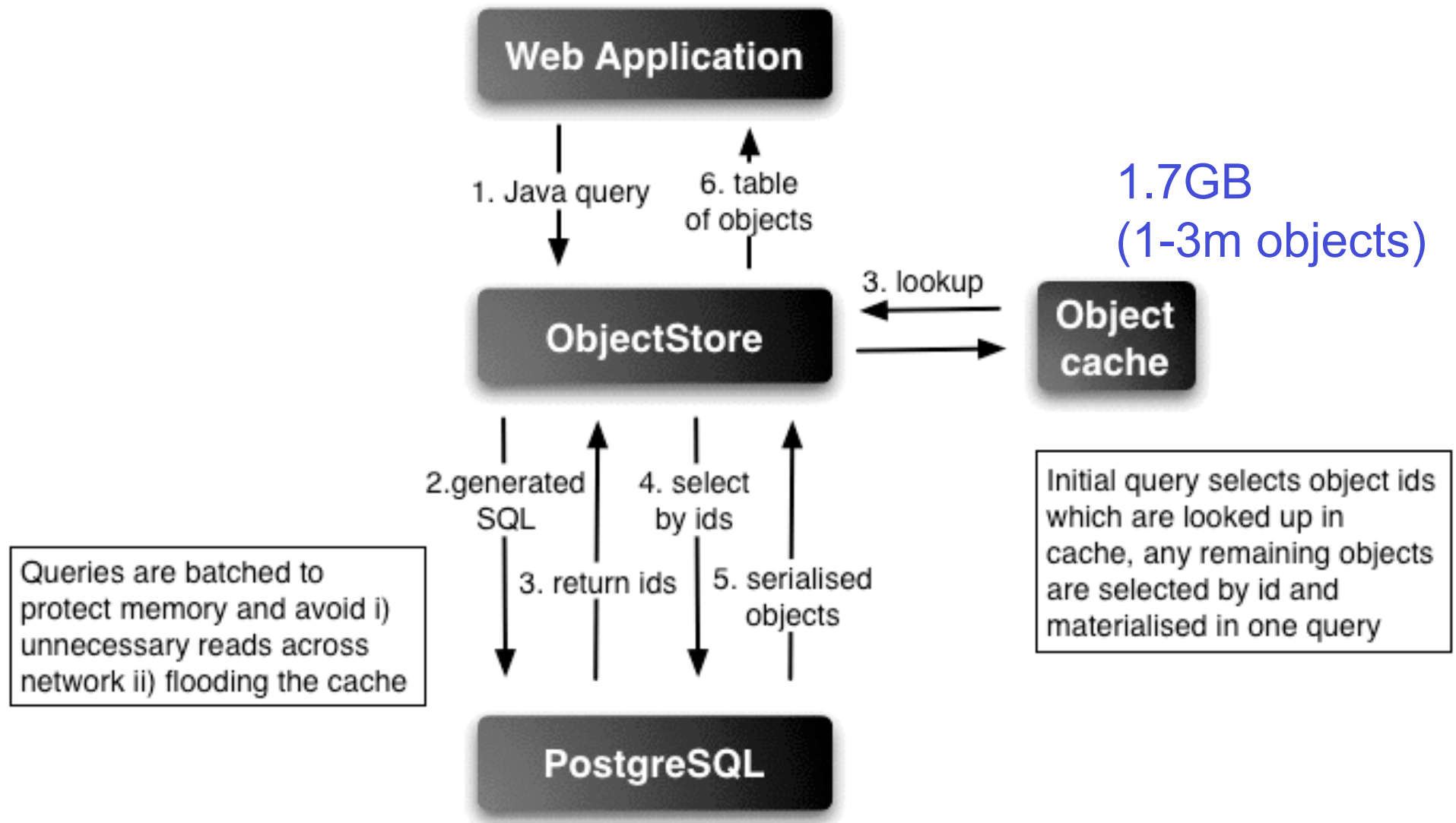
Class	identifier	symbol	chromosomeLocation	length	organism.name	
Exon	xfile-exon-2		fake_chromosome : 14687-14720	3 <span style="border: 2px solid red; border-radius: 50%; padding: 2px;">FASTA...</span>	Drosophila melanogaster	[details...]
Exon	xfile-exon-1		fake_chromosome : 13691-13767	7 <span style="border: 2px solid red; border-radius: 50%; padding: 2px;">FASTA...</span>	Drosophila melanogaster	[details...]

[show in table...]

**overlappingFeatures** 0 LocatedSequenceFeature



# Executing a Query





# InterMine Queries

## Quicksearch

### Java API:

```
Query q = new Query();
QueryClass qcObj = new QueryClass(Gene.class);
q.addFrom(qcObj);
q.addToSelect(qcObj);

QueryField qf = new QueryField(qcObj, "identifier");

SimpleConstraint sc = new SimpleConstraint(qf, ConstraintOp.MATCHES, new
    QueryValue("x-%"));
q.setConstraint(sc);
```

### IQL:

```
SELECT DISTINCT a1_.identifier AS a2_ FROM org.flymine.model.genomic.Gene
AS a1_ WHERE a1_.identifier LIKE 'x-%'
```

### Perl API:

```
my $genes = InterMine::Gene::Manager->get_genes(query => [
    identifier => { like => 'x-%' },,);
```



# Larger Query

Results for template: **Gene --> Transcripts and exons + chromosomal locations and lengths.**

*For a particular gene, list the transcripts and exons and show the chromosomal locations and lengths of the transcripts and exons.*

Page size  ▾

First Previous Next Last

<input type="checkbox"/>	Gene > identifier ▶▶	<input type="checkbox"/>	Gene > transcripts > identifier ◀▶▶	<input type="checkbox"/>	Gene > transcripts > chromosome > identifier ◀▶▶	<input type="checkbox"/>	Gene > transcripts > chromosomeLocation > start ◀▶▶	<input type="checkbox"/>	Gene > transcripts > chromosomeLocation > end ◀▶▶	<input type="checkbox"/>	Gene > transcripts > length ◀▶▶	<input type="checkbox"/>	Gene > transcripts > exons > identifier ◀▶▶
<input type="checkbox"/>	CG2328	<input type="checkbox"/>	CG2328-RA [MRNA]	<input type="checkbox"/>	2R	<input type="checkbox"/>	5491054	<input type="checkbox"/>	5492592	<input type="checkbox"/>	1468	<input type="checkbox"/>	CG2328:1
<input type="checkbox"/>	CG2328	<input type="checkbox"/>	CG2328-RA [MRNA]	<input type="checkbox"/>	2R	<input type="checkbox"/>	5491054	<input type="checkbox"/>	5492592	<input type="checkbox"/>	1468	<input type="checkbox"/>	CG2328:2

Displaying all 2 rows

First Previous Next Last

- Joins 7 classes - all are on select list of query
- Here, one gene (two rows) ~2 seconds
- All genes, all organisms ~300,000 rows in 36 sec (not via pre-computation)
- 



```

public class BakeOff {
    public static void main(String[] args) throws Exception {
        // code to get the "xfile" gene
        ObjectStore os = ObjectStoreFactory.getObjectStore("os.production");
        Query q = new Query();
        QueryClass qcObj = new QueryClass(Gene.class);
        q.addFrom(qcObj);
        QueryField qf = new QueryField(qcObj, "identifier");
        q.addToSelect(qf);
        SimpleConstraint sc = new SimpleConstraint(qf, ConstraintOp.EQUALS, new QueryValue("xfile"));
        q.setConstraint(sc);
        System.err.println("query: " + q);
        Results res = os.execute(q);

        // a Results object is a List of Lists
        List rr = (List) res.get(0);
        Gene gene = (Gene) rr.get(0);

        System.err.println ("symbol: " + gene.getIdentifier());

        // a BioEntity in FlyMine has a collection of Synonym objects -
        // we need Synonym.value for each Synonym
        System.err.print ("synonyms: ");
        Iterator synIter = gene.getSynonyms().iterator();
        while (synIter.hasNext()) {
            Synonym syn = (Synonym) synIter.next();
            System.err.print (syn.getValue() + ' ');
        }
    }
}

```





```

System.err.println ("description: " + gene.getDescription());

    // get the class name, but we already know that the gene is a Gene
    System.err.println ("type: " + gene.getClass().getName());

    // make a List from a the Set of exons for this Gene
    List exons = new ArrayList(gene.getExons());
    Exon exon1 = (Exon) exons.get(0);
    Exon exon2 = (Exon) exons.get(1);

    // get the start and end via the Location object
    System.err.println ("exon1 start: " + exon1.getChromosomeLocation().getStart());
    System.err.println ("exon1 end: " + exon1.getChromosomeLocation().getEnd());
    System.err.println ("exon2 start: " + exon2.getChromosomeLocation().getStart());
    System.err.println ("exon2 end: " + exon2.getChromosomeLocation().getEnd());

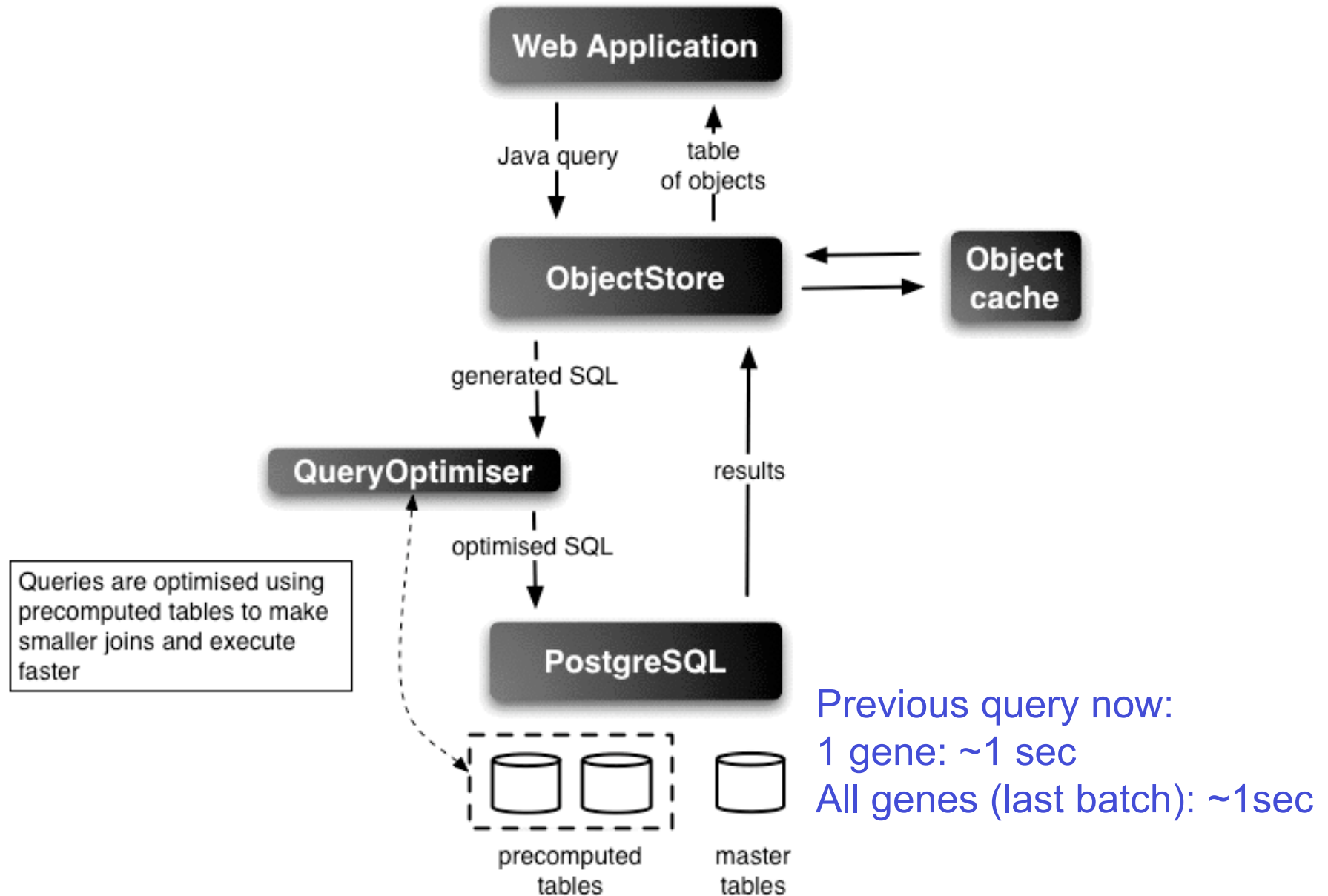
    // write out the first cds
    List cdss = new ArrayList(gene.getCDSs());
    FlyMineSequence flymineSequence = FlyMineSequenceFactory.make((CDS) cdss.get(0));

    // use BioJava to output the sequence
    Annotation annotation = flymineSequence.getAnnotation();
    annotation.setProperty(FastaFormat.PROPERTY_DESCRIPTIONLINE,
                           gene.getIdentifier() + " cds");
    SeqIOTools.writeFasta(System.err, flymineSequence);
}
}

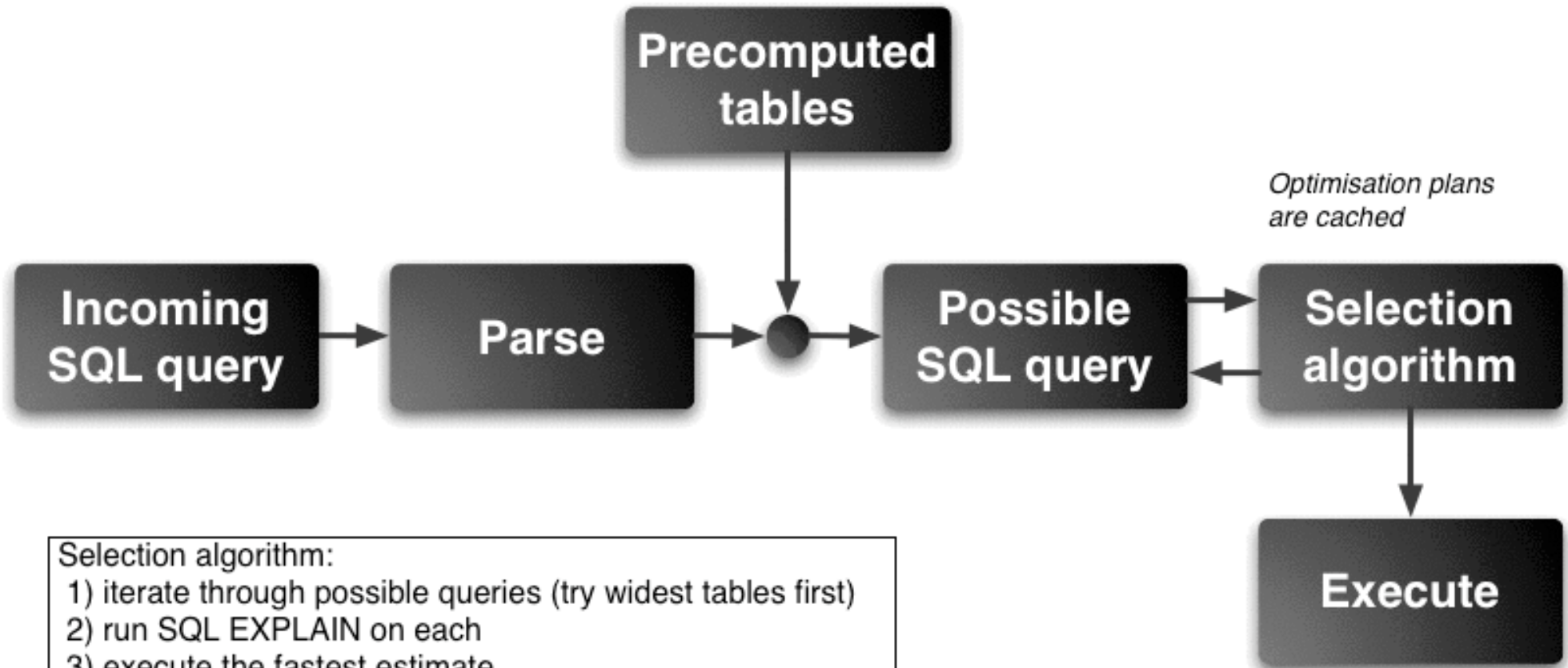
```



# Query Optimisation



# Query Optimisation



## Selection algorithm:

- 1) iterate through possible queries (try widest tables first)
- 2) run SQL EXPLAIN on each
- 3) execute the fastest estimate

## Heuristics to improve performance:

- a) avoid optimising very simple queries
- b) run immediately if a one table query found
- c) minimise number of EXPLAINS needed
- d) avoid EXPLAINS on very large joins
- e) bail out if taking too long



# Query Optimisation - implications

- \_ Performance optimisation not tied to schema design
- \_ Can adapt performance optimisation to usage
- \_ Template queries pre-computed (40 run per gene details page - renders in seconds)



# Future Directions

- Import of any Chado database
- Improvements to build process
- User groups for sharing of templates/ bags
- Public bags/ tagging
- Relative genome coordinate based queries
- REST-style web services from templates
- Elaboration of bag details page:
  - graphical widgets, bag comparison, protein interaction graph functions, gene expression, generalised GoSTAT, chromosome distribution, enriched publications...



# Cons

- slow build times (which are getting better)
- lots of disk space but it's a data warehouse and disks are cheap
- configuration is still a bit complicated, but improving fast
- we can't do much with locations yet except query starts and ends, although we do have the overlappingFeatures collections now...
- export is still limited to tab or comma delimited or FASTA
- we don't handle sequence very well:
  - e.g. each chromosome sequence is stored in one big text field in PostgreSQL
- we can't do queries involving sizes of collections of things (e.g. find the genes with only 1 transcript)



# Acknowledgements

Richard Smith  
Kim Rutherford  
Matthew Wakeling  
Xavier Watkins  
Julie Sullivan

Rachel Lyne  
Hilde Janssens  
François Guillier  
Philip North  
Gos Micklem

*Andrew Varley, Mark Woodbridge, Tom Riley,  
Peter McLaren, Debashis Rana, Wenyan Ji,  
Markus Brosch, Florian Reising*

[www.flymine.org](http://www.flymine.org)

[www.intermine.org](http://www.intermine.org)



FlyMine is funded by the Wellcome Trust (grant no. 067205), awarded to M. Ashburner, G. Micklem, S. Russell, K. Lilley and K. Mizuguchi.



UNIVERSITY OF  
CAMBRIDGE

# Bioinformatics Position in Berkeley

[www.berkeleybop.org/jobs](http://www.berkeleybop.org/jobs)

- \_ Work with Chris Mungall and Gos Micklem on Chado and InterMine
- \_ Gather requirements and design data models for novel, high-throughput, large-scale experiments
- \_ Build a workflow system, a la UCSC, for data submission
- \_ Integrate data from WormBase and FlyBase
- \_ Integrate database with GBrowse, BioMart, and other software
- \_ Applications and questions may be sent to Nicole Washington (nlwashington@lbl.gov)







UNIVERSITY OF  
CAMBRIDGE