# InterMine

## Object integration and warehousing software

# FlyMine

## An integrated database for *Drosophila* and *Anopheles* genomics

Andrew Varley,
Genetics Department,
Cambridge University

FlyMine

# Why yet another database?

- Currently lots of little databases
  - fine for "browsing", bad for "querying"
- Hard to query across them
  - lots of "cut and paste" on web pages
- Massive amounts of experimental data (microarray, proteomics) being produced
- Need to tie this information together

FlyMine

# Outline

External databases

Users

| | | |
|---|---|---|
| FlyBase | | Web user |
| Array Express | FlyMine integrated database | Web service |
| . . . | | |
| Proteomics | | Bulk access |

FlyMine

# InterMine - aims

- Integrate data from multiple sources

- Allow <u>arbitrary</u> queries from users

- Queries based on objects, not SQL

- Complex models (multiple inheritence)

- Multiple query interfaces (Java, OQL, graphical, HTML, etc)

- Different classes of user (web-based, SOAP)

- Open source!

FlyMine

# Query interfaces – OQL

*"Show gene expression data for genes which have GO term GO:0000278 applied"*

```
SELECT gene, exp
FROM Gene AS gene, GOTerm AS term,
     Experiment AS exp
WHERE gene.GOTerms CONTAINS term
AND term.code = "GO:0000278"
AND gene.experiments CONTAINS exp
AND exp.type = "Gene Expression"
```
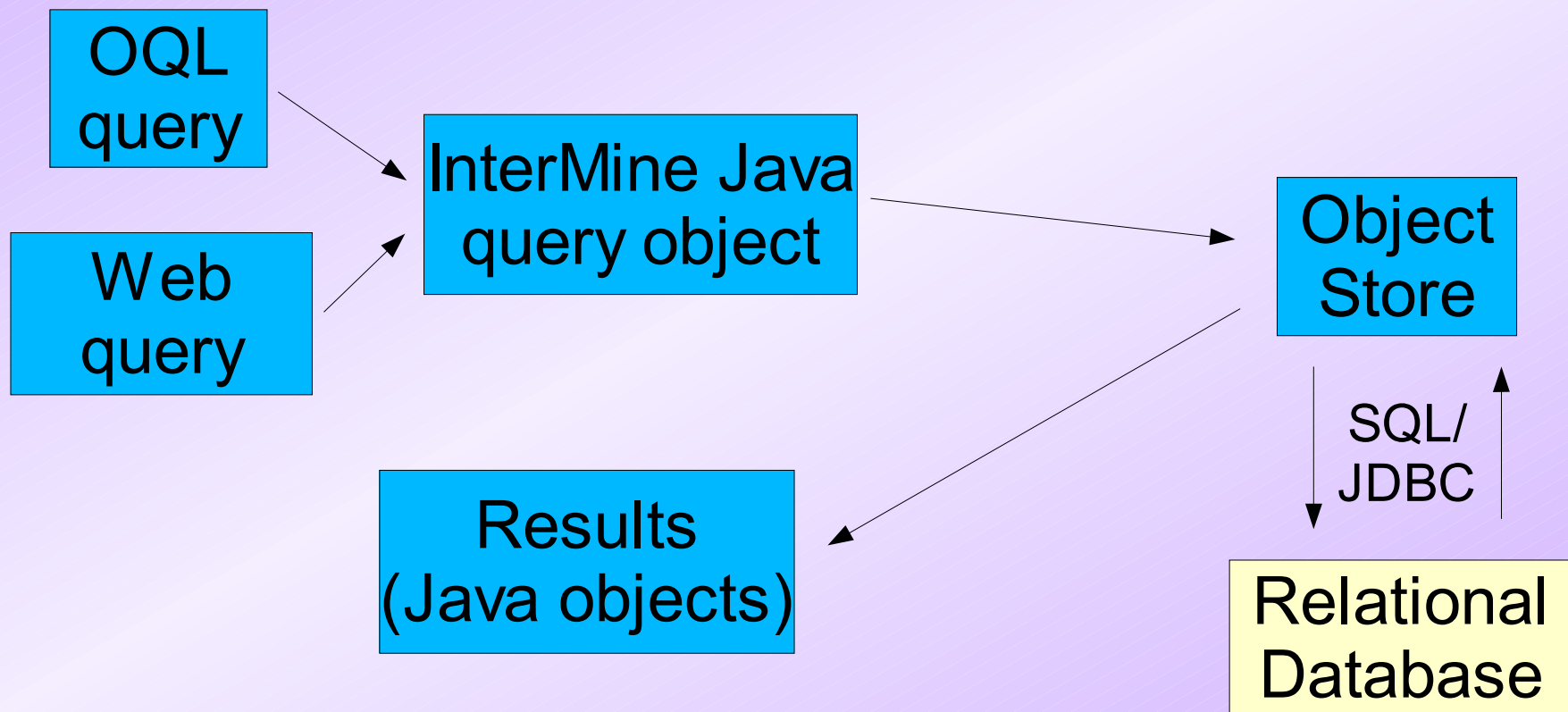
FlyMine

# Query interfaces – Java

*"Show gene expression data for genes which have GO term GO:0000278 applied"*

```
Query q = new Query();
QueryClass qcGene = new QueryClass(Gene.class);
QueryClass qcTerm = new QueryClass(GOTerm.class);
QueryField qfGeneTerm = new QueryField(qcGene, "GOTerms")
ConstraintSet cs = new ConstraintSet();
Constraint con1 = new ContainsConstraint(qfGeneTerm,
                    ContainsConstraint.CONTAINS, qcTerm);
cs.add(con1);
.
.
q.setConstraint(cs);
```

FlyMine

# InterMine ObjectStore

OQL query

Web query

InterMine Java query object

Object Store

Results (Java objects)

SQL/ JDBC

Relational Database

FlyMine

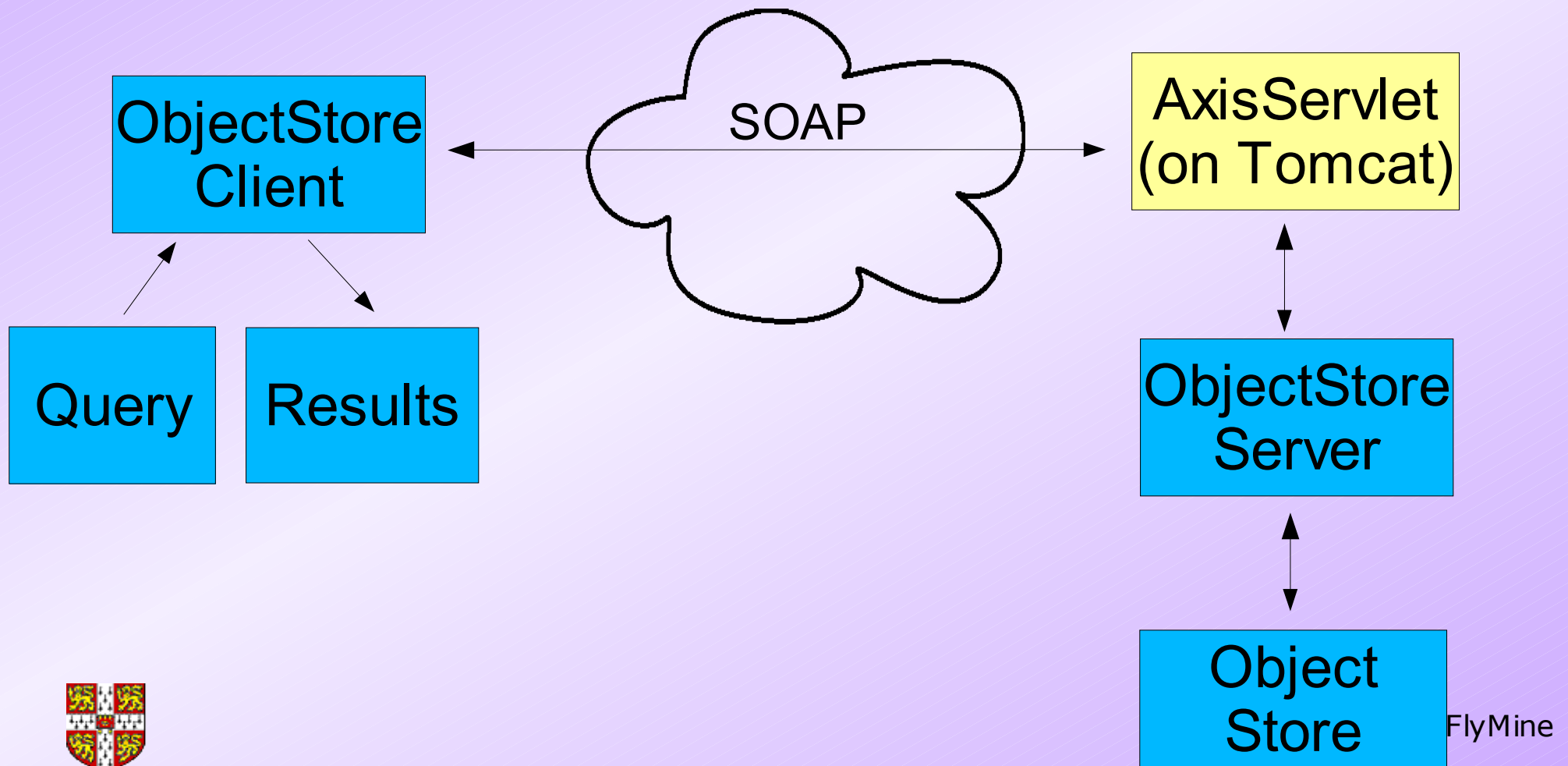# ObjectStore interface

```
public interface ObjectStore {

    Results execute(Query q);
    List execute(Query q, int start, int limit);
    int count(Query q);
    ResultsInfo estimate(Query q);
    Object getObjectById(Integer id);
    Model getModel();
}
```

# Web service

**Client**

**Server**

ObjectStore
Client

Query          Results

SOAP

AxisServlet
(on Tomcat)

ObjectStore
Server

Object
Store

FlyMine

# InterMine ObjectStore

- Results are tables of Java objects

  – All collections and references are proxied

- Results rows are fetched in batches from the underlying database server

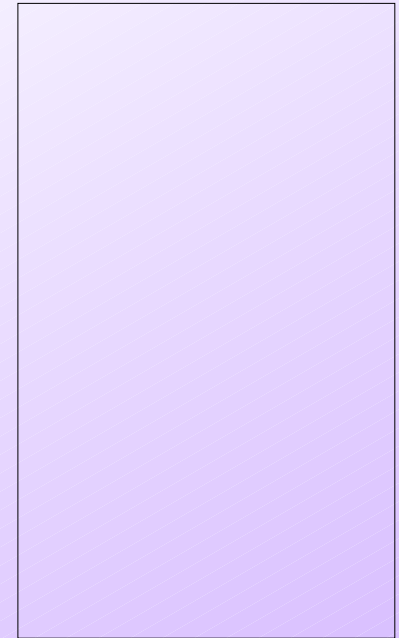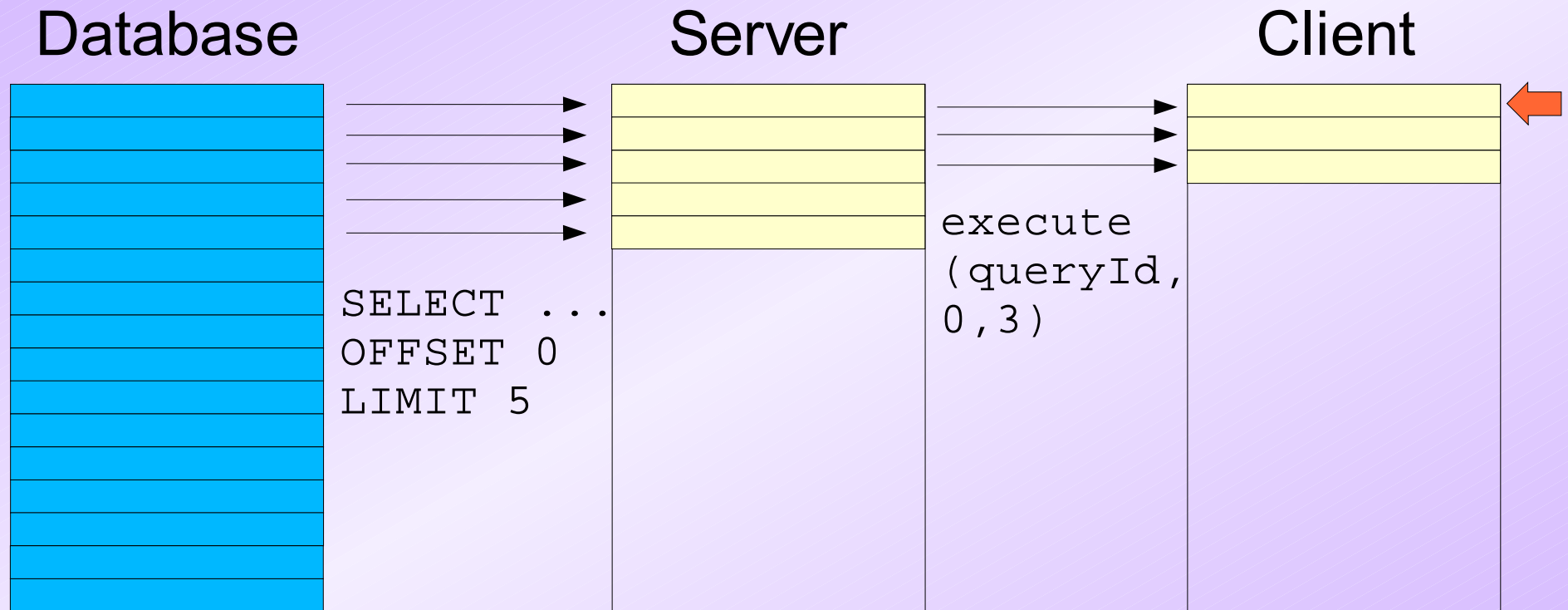- Results rows are pre-emptively fetched and cached (per JVM).
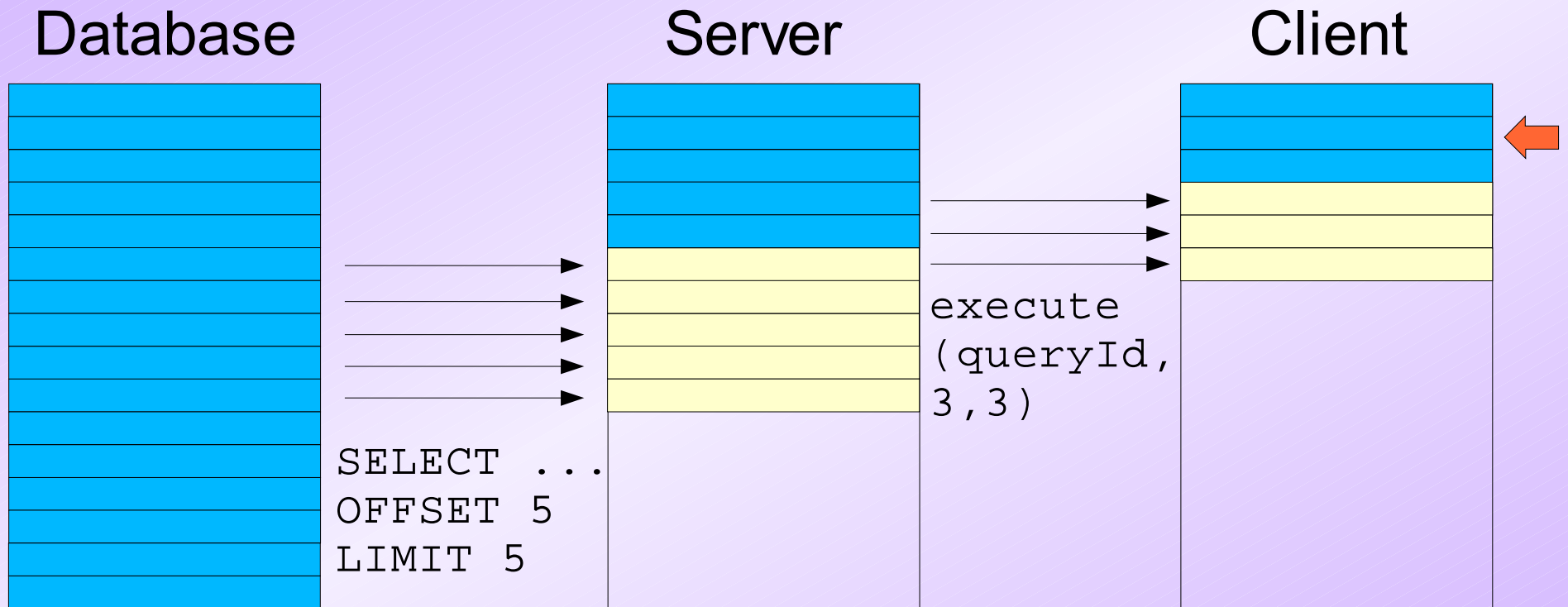
FlyMine

# Results

**Database**

**Server**

**Client**

# Results

**Database**                    **Server**                    **Client**



SELECT ...
OFFSET 0
LIMIT 5

execute
(queryId,
0,3)

FlyMine

# Results

**Database**

**Server**

**Client**

```
SELECT ...
OFFSET 5
LIMIT 5
```

```
execute
(queryId,
3,3)
```

FlyMine

# Results



Database          Server          Client

FlyMine

# Results

Database Server Client

FlyMine

# Results

Database                    Server                    Client

SELECT ...
OFFSET 10
LIMIT 5

execute
(queryId,
6,3)

FlyMine

# Results



Database

Server

Client

FlyMine

# Results

Database                    Server                      Client
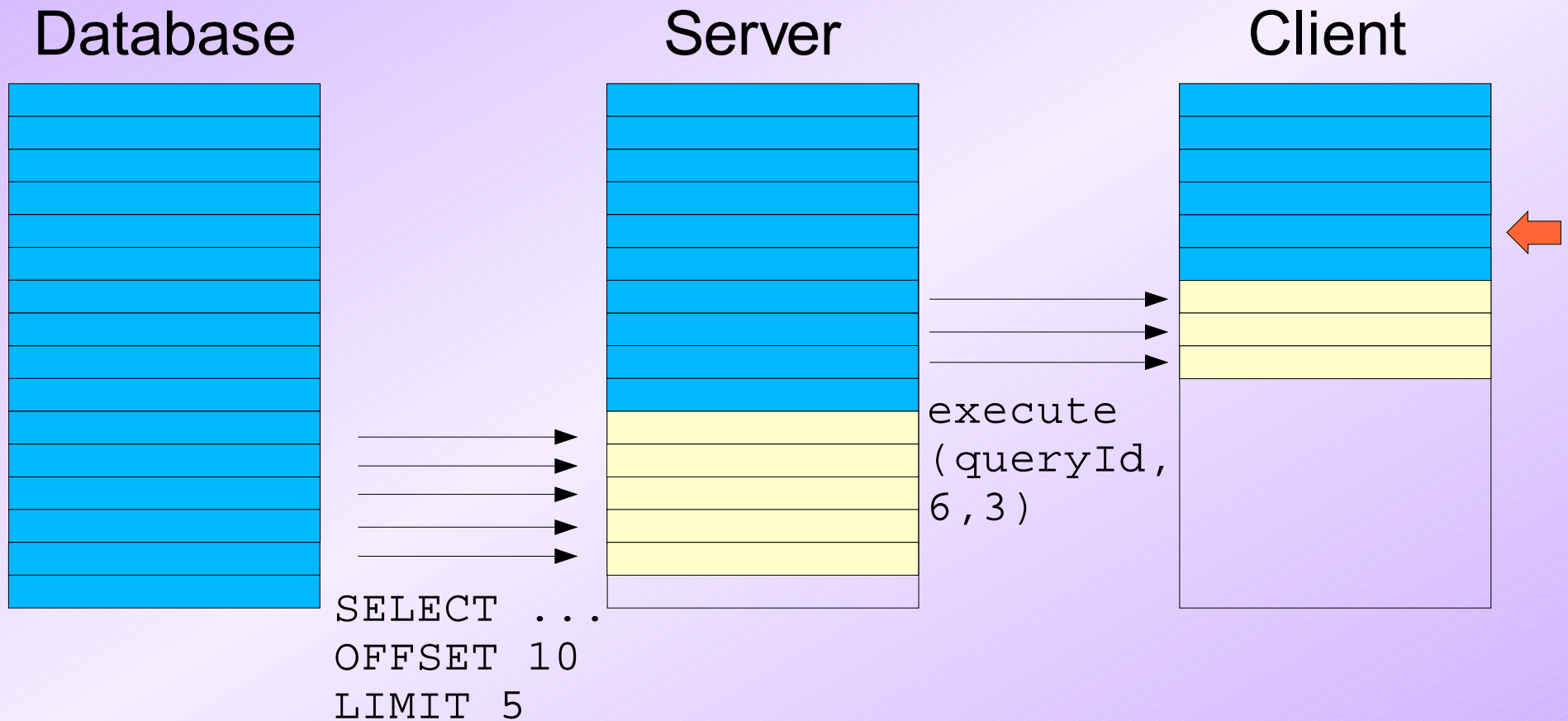
# Results

Database                    Server                      Client

execute
(queryId,
9,3)

FlyMine

# Results

Database        Server        Client

FlyMine

# Results

Database
Server
Client

# Results

Database              Server                Client

SELECT ...            execute
OFFSET 15             (queryId,
LIMIT 5               12,3)

FlyMine

# Data model

- Data model can be generated from one of a number of sources



UML (XMI) → InterMine model
XML Schema → InterMine model
OWL → InterMine model

InterMine model → Java classes
InterMine model → Web site
InterMine model → Torque schema → Database schema
InterMine model → XML binding

FlyMine

# Model Integration

SOFA (DAG) → SOFA OWL

ensembl (UML) → ensembl OWL

chado (UML) → chado OWL

MAGE (XMI) → MAGE OWL

orthologue OWL

PSI (XML Schema) → PSI OWL

SOFA OWL, ensembl OWL, chado OWL, MAGE OWL, orthologue OWL, PSI OWL → FlyMine OWL

OWL – Web Ontology Language
- standard for ontology description
- allows cardinality

merge description (OWL) → FlyMine OWL

biologist defined configuration for integrating models

FlyMine OWL → FlyMine model

FlyMine

# Integrate existing standards

- Link existing and emerging standards to define FlyMine model

  - e.g. SO, MAGE, PSI,GO, other ontologies
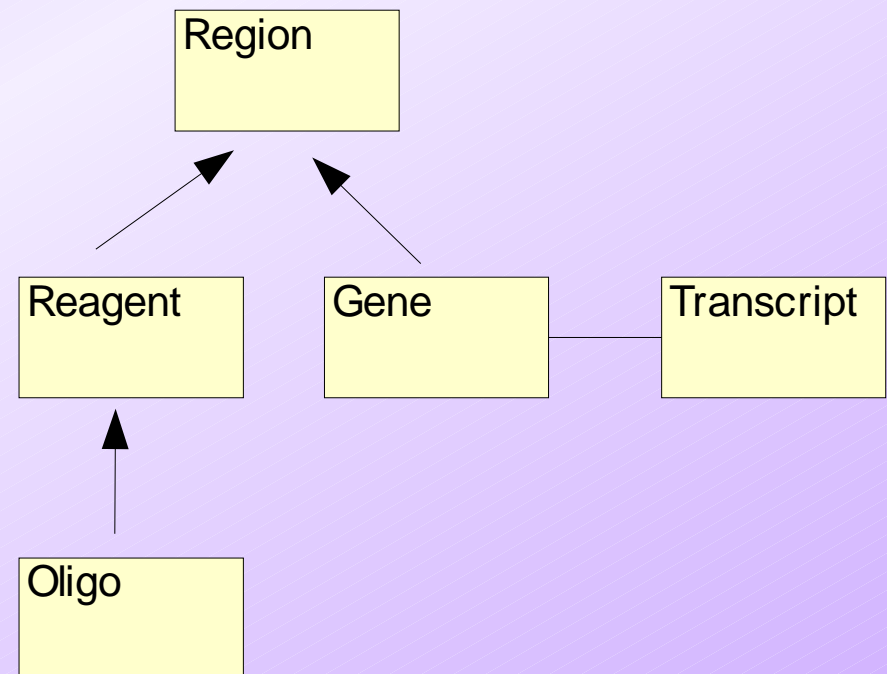
- Avoids creating a 'schema of everything'

- Easy to add new types of data

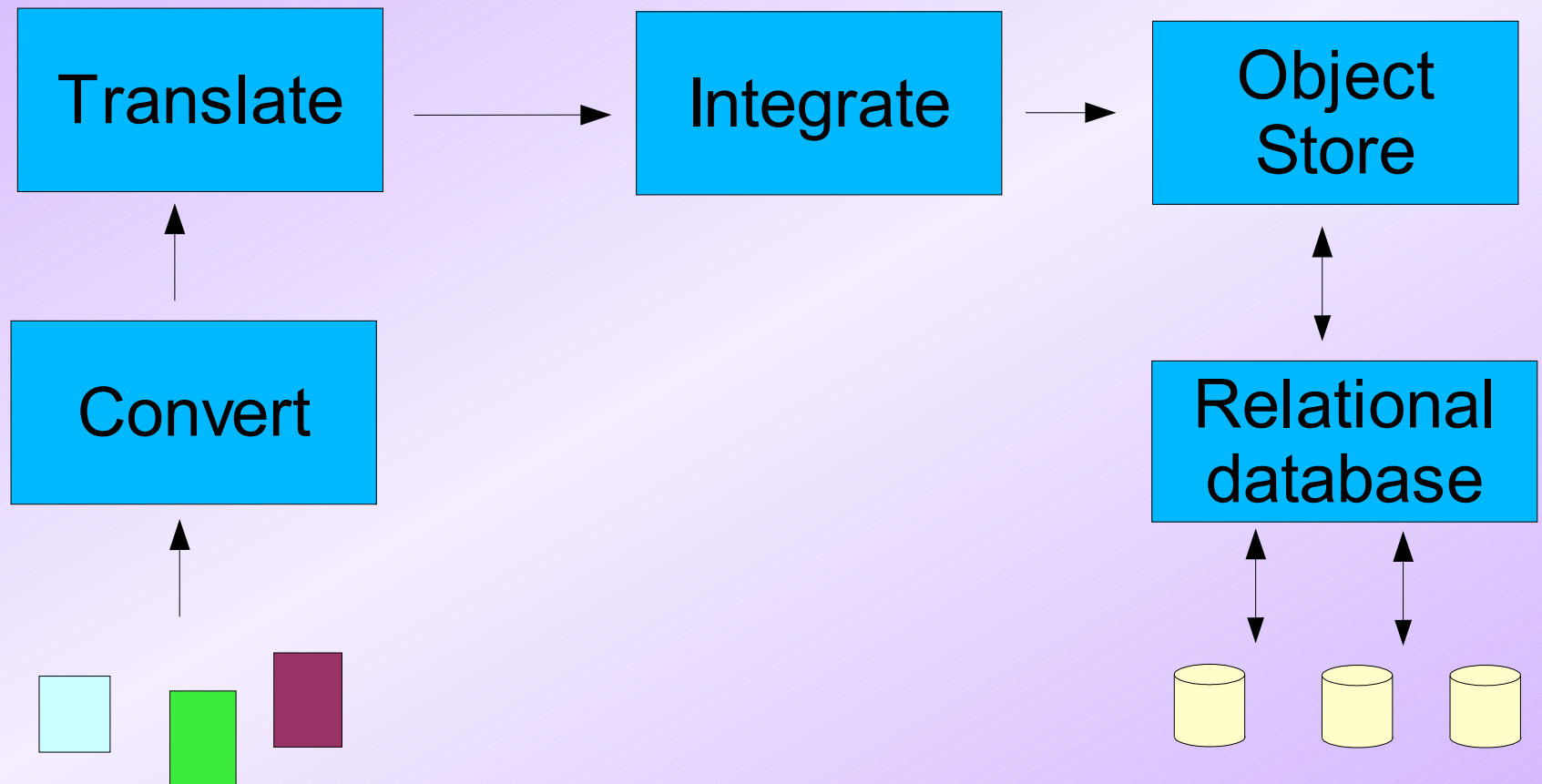- Evolve as standards change

FlyMine

# Example – Sequence Ontology

- All terms in ontology become classes in model

```
@is_a@region
   @is_a@gene
      @part_of@transcript
   @is_a@reagent
      @is_a@oligo
```

# Data loading pipeline



FlyMine

# Conversion from external sources

SQL database

DBConverter

"Items" XML

XMLConverter

XML

FileConverter/ custom converter

Flat files

```
<items>
  <item id="1" class="Feature">
    <attribute name="name"
               value="ABC1"/>
    <reference name="organism"
               ref_id="2"/>
    <collection name="synonyms">
      <reference ref_id="3"/>
      <reference ref_id="4"/>
    </collection>
  </item>
    .
    .
    .
</items>
```

FlyMine

# Translation

**Feature**

name=
    "ABC1"
type=
    "Gene"

**Transcript**

identifier=
    "ABC1_1"

→

**Gene**

name=
    "ABC1"

**Transcript**

name=
    "ABC1_1"

FlyMine

# Translation

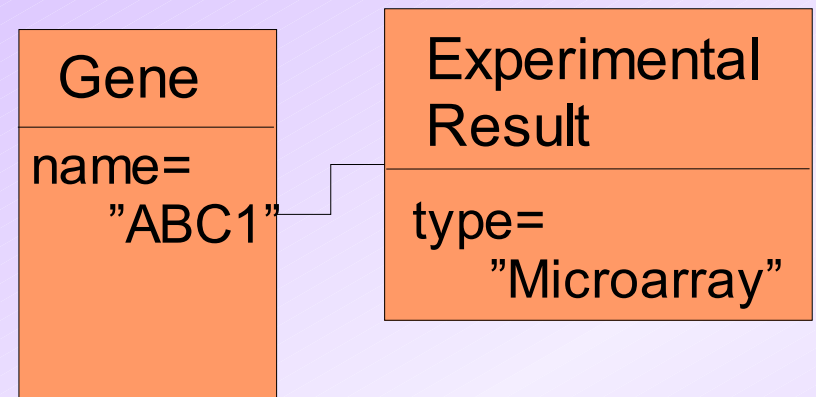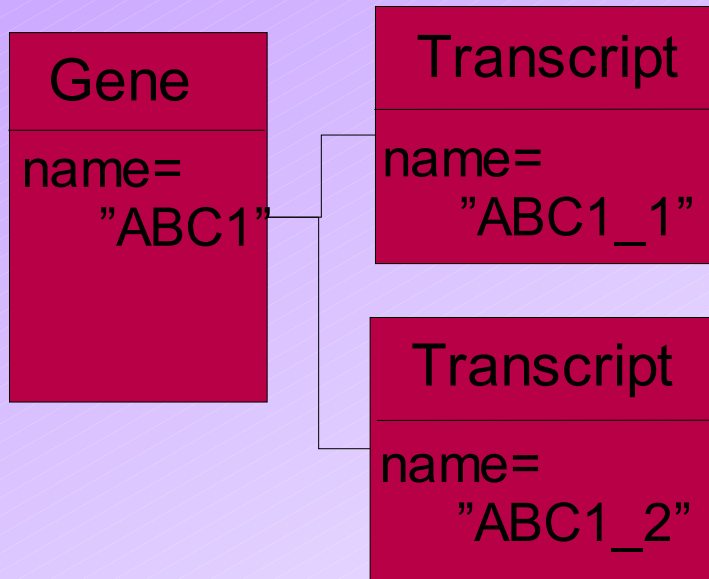- Controlled by OWL "merge description"

```
:Gene a owl:Class;
  rdfs:subClassOf exampleSrc:Feature ;
  rdfs:subClassOf
    [ a owl:Restriction;
      owl:onProperty exampleSrc:Feature__type
      owl:hasValue "Gene"
    ] .
:Transcript a owl:Class;
  owl:equivalentClass exampleSrc:Transcript .
:Transcript__name a owl:DatatypeProperty;
  owl:equivalentProperty exampleSrc:Transcript__identifier
```

FlyMine

# Integration

- Define "primary keys" for each object type
- Define "primary keys" that each source uses
- Define priorities for fields from different sources
- IntegrationWriter keeps track of originating sources for each field of each object
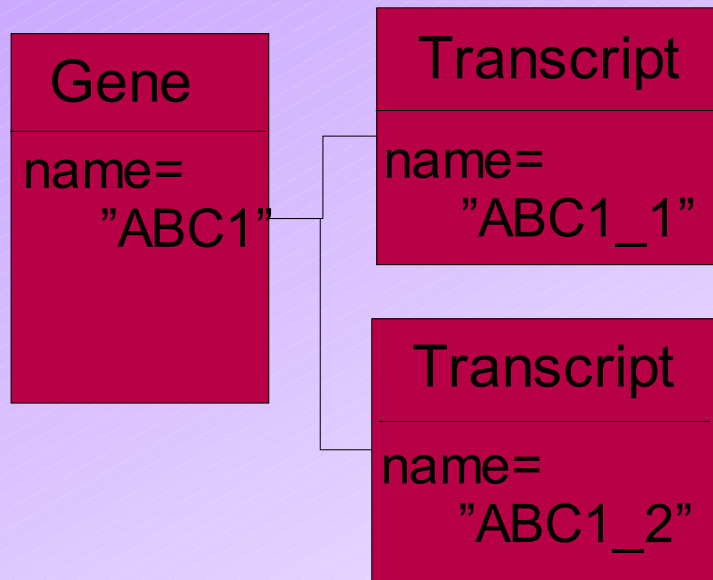
```
Primary keys:
   Gene: name
   Gene: flyBaseName
   Transcript: name,gene
```

New objects

Currently in database

**Gene**

name=
"ABC1"

**Transcript**

name=
"ABC1_1"

**Transcript**

name=
"ABC1_2"

**Gene**

name=
"ABC1"

**Experimental Result**

type=
"Microarray"

**Transcript**

name=
"ABC1_1"

**Transcript**

name=
"ABC1_2"

**Primary keys:**
    **Gene: name**
    **Gene: flyBaseName**
    **Transcript: name,gene**

New objects

Currently in database

FlyMine

## Gene

flyBaseName=
    "FBgn12345"

## Transcript

name=
    "ABC1_1"

## GOTerm

code=
    "GO56789"

## Gene

name=
    "ABC1"

## Experimental Result

type=
    "Microarray"

## Transcript

name=
    "ABC1_1"

## Transcript

name=
    "ABC1_2"

**Primary keys:**
  **Gene: name**
  **Gene: flyBaseName**
  **Transcript: name,gene**

New objects

Currently in database

FlyMine

**Gene**

name="ABC1"
flyBaseName=
    "FBgn12345"

**Gene**

name=
    "ABC1"

**Experimental
Result**

type=
    "Microarray"

**Transcript**

name=
    "ABC1_1"

**Transcript**

name=
    "ABC1_2"

**Primary keys:**
  **Gene: name**
  **Gene: flyBaseName**
  **Transcript: name,gene**

**Gene**

flyBaseName=
    "FBgn12345"

**Transcript**

name=
    "ABC1_1"

**GOTerm**

code=
    "GO56789"

## Gene

name="ABC1"
flyBaseName=
    "FBgn12345"

## Gene

name="ABC1"
flyBaseName=
    "FBgn12345"

## Experimental Result

type=
    "Microarray"

## Transcript

name=
    "ABC1_1"

## Transcript

name=
    "ABC1_2"

## GOTerm

code=
    "GO56789"

**Primary keys:**
  **Gene: name**
  **Gene: flyBaseName**
  **Transcript: name,gene**

New objects

Currently in database

FlyMine

# Lessons learned

- Need to batch database reads and writes to avoid round trip time

  - without: 6 months to load

  - with: 6 hours to load

- Cache objects to avoid database lookups

- Use multi-thread model
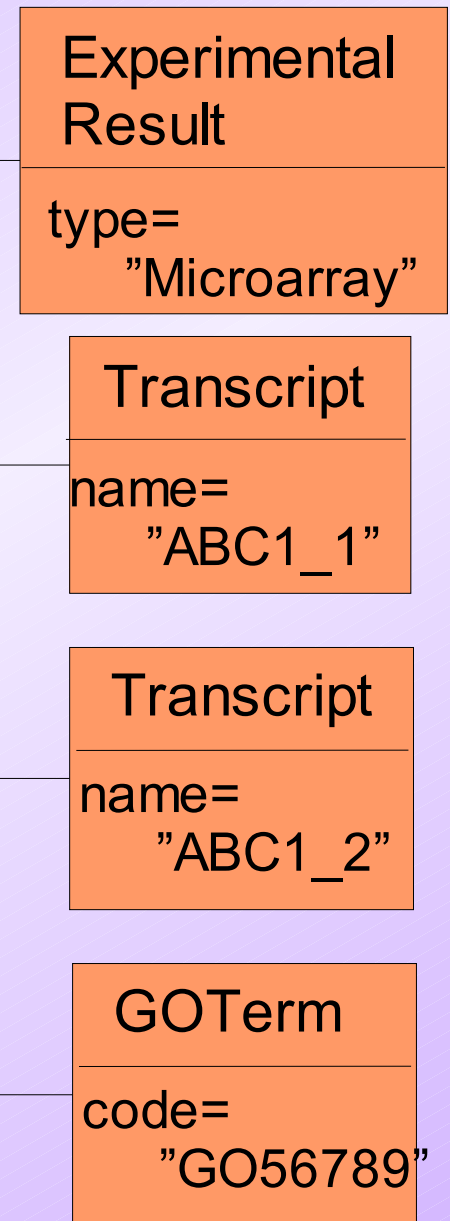
- GB ethernet link between computational and database machines useful

FlyMine

# Arbitrary queries – problems

- Badly formed queries may overload database server

- Difficult to optimise database for all queries
  - Which indexes to use?
  - Slow response to complex queries involving multi-table joins

# Arbitrary queries – solutions

- Close relationship with database server query planner

  - ask how long a query will take <u>before</u> attempting to run it (~3ms)

  - Disallow queries that will take longer than a certain threshold

- Store data massively redundantly in "precomputed tables" and rewrite incoming queries on-the-fly.

FlyMine

# Query optimisation – aims

- Provide a generic way of speeding up complex queries to any read-only SQL database

- Minimal parsing/computational overhead

- Transparent to users/applications

  - no new schema/model to learn

- Make available as standalone module

FlyMine

# Query optimisation architecture



ObjectStore  → SQL →  Query optimiser  → SQL →  Relational Database

JDBC ResultSet

Master tables

Precomputed tables

FlyMine

# Query optimisation

Precomputed tables

Incoming SQL query → Parse → ● → List of possible queries → EXPLAIN them all → Run the shortest

FlyMine

# Query optimisation – example

## Master tables

### Genes

| id | name |
|----|------|
|    |      |

### GOTerms

| geneid | code |
|--------|------|
|        |      |

### Experiments

| geneid | date | type |
|--------|------|------|
|        |      |      |

## Precomputed tables

### GenesGOTerms



```
CREATE TABLE genesgoterms AS
SELECT genes.id AS genes_id,
       genes.name AS genes_name,
       goterms.geneid AS goterms_geneid,
       goterms.code AS goterms_code
FROM genes, goterms
WHERE genes.id = goterms.geneid
```

### GenesGOTermsExperiments



```
CREATE TABLE ....
SELECT ....
FROM genes, goterms, experiments
WHERE genes.id = goterms.geneid
AND genes.id = experiments.geneid
```

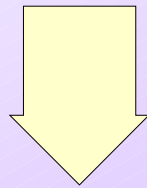FlyMine

# Query optimisation – example

*"Show dates that gene expression experiments were performed on genes which have GO term GO:0000278 applied"*

```
SELECT genes.name, experiments.date
FROM genes, goterms, experiments
WHERE genes.id = goterms.geneid
AND genes.id = experiments.geneid
AND goterms.code = "GO:0000278"
AND experiments.type = "Gene Expression"
```
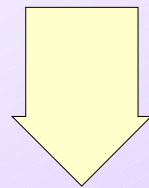
EXPLAIN

5 seconds

FlyMine

# Query optimisation – example

```
SELECT genes_name AS name,
       experiments.date
FROM genesgoterms  experiments
WHERE genes_id = experiments.geneid
AND goterms_code = "GO:0000278"
AND experiments.type = "Gene Expression"
```
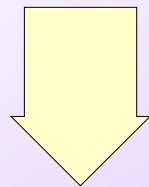
EXPLAIN

1 second

FlyMine

# Query optimisation – example

```
SELECT genes_name AS name,
       experiments_date AS date
FROM genesgotermsexperiments
WHERE goterms_code = "GO:0000278"
AND experiments_type = "Gene Expression"
```

EXPLAIN

200ms

FlyMine

# Query optimisation – example

- Without optimiser

  Time = 3ms  +  5 seconds

  (EXPLAIN)    (EXECUTE)

- With optimiser

  Time = 9ms   +  100ms  +  200ms

  (3 EXPLAINs)  (PARSE)    (EXECUTE)

# Query optimisation – summary

- For complex queries, the SQL optimisation module can produce large speed increases

- Optimisation is transparent to the user or application generating the SQL.

- Choosing which precomputed tables to store is important – may need to analyse incoming queries.

FlyMine

# Query optimiser – summary

- Optimiser trades off disk space for improved query performance.

- Schema independent.

- Can be used in conjunction with P6Spy to intercept JDBC calls from existing software to a database.

- Available from www.intermine.org !!

# Summary

- Current status
  - generic SQL query optimiser
  - powerful object data warehouse
  - 2 query interfaces (OQL + Java)
  - web front end
  - web service
  - framework for data loading/integration
  - no model-specific code
  - FlyMine database alpha release

FlyMine

# Summary

- Coming up
  - graphical query interface
  - tools for auto-generation of best set of precomputed tables
  - manual curation tools

FlyMine

# Acknowledgements

The FlyMine team:

Andrew Varley          François Guillier
Richard Smith          Rachel Lyne
Matthew Wakeling       Kim Rutherford
Mark Woodbridge        Gos Micklem

More information and download at www.intermine.org

FlyMine