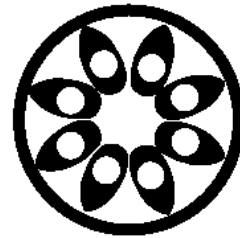# InterMine as a Science Gateway for Systems Biology
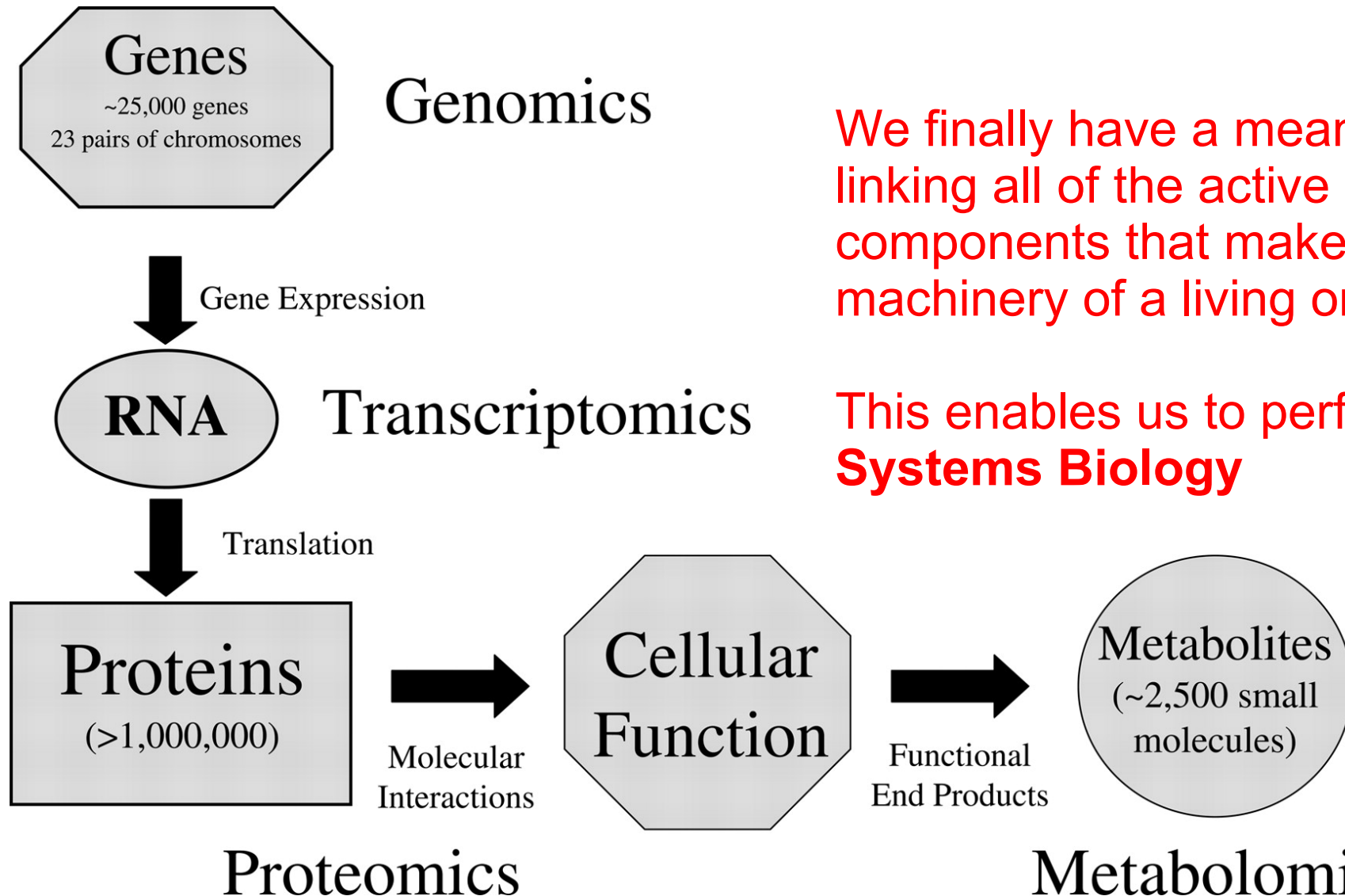
Aaron Golden
School of Mathematics, Statistics & Applied Mathematics
National University of Ireland, Galway

IWSG 2017 – Poznan, Poland

NUI Galway
OÉ Gaillimh

InterMine

# *Talk Roadmap*

(1) Desktop technologies & systems biology

(2) Problem – how to leverage resulting data space

(3) Sequence Ontology & Data Schemas
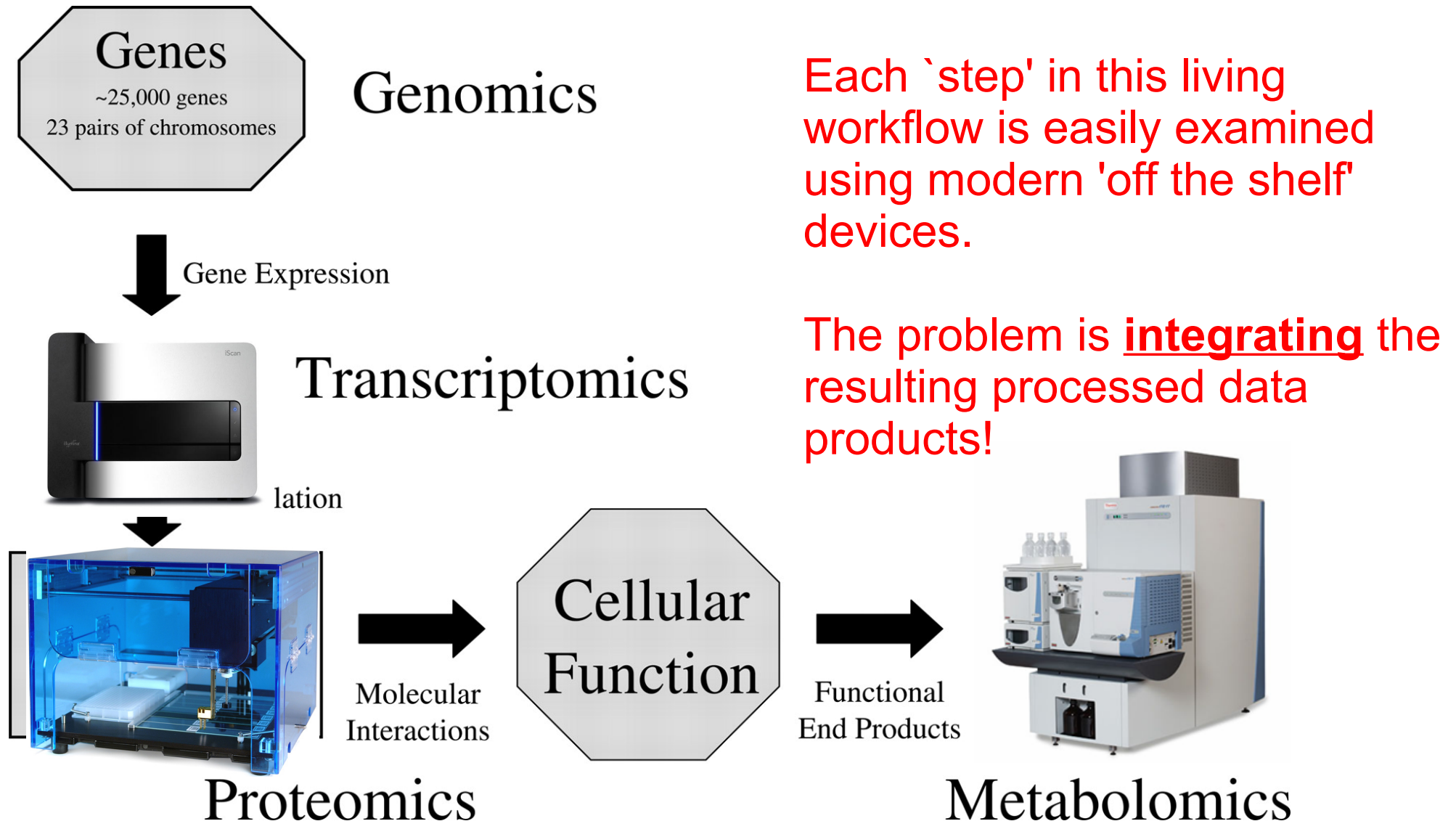
(4) InterMine – how it works as a Science Gateway

# The Good News...

Genes
~25,000 genes
23 pairs of chromosomes

Genomics

Gene Expression

RNA

Transcriptomics

Translation

Proteins
(>1,000,000)

Molecular Interactions

Proteomics

Cellular Function

Functional End Products

Metabolites
(~2,500 small molecules)

Metabolomics

We finally have a means of fully linking all of the active components that make up the machinery of a living organism.

This enables us to perform **Systems Biology**

NUI Galway
OÉ Gaillimh

InterMine

# *The Problem...*

Genes
~25,000 genes
23 pairs of chromosomes

**Genomics**

Gene Expression

**Transcriptomics**

lation

**Proteomics**

Molecular Interactions

**Cellular Function**

Functional End Products

**Metabolomics**

Each `step' in this living workflow is easily examined using modern 'off the shelf' devices.

The problem is **integrating** the resulting processed data products!

InterMine

# *Ontologies to the Rescue...*

Solution lies in the



Forms a schema based around the ultimate frame of reference
 - the genome

Genomic features described as parts of gene models, assembly components, experimental results generating annotations defined as being associated with a specific location...

Forms an excellent means ***to bind the various systems biology datasets together in the same data space***

# The basis for building a schema...



Genes
~25,000 genes
23 pairs of chromosomes

List of names, locations, properties
Genomics

Gene Expression

List of names, locations, quantity
RNA    Transcriptomics

Lab work yields different types of data files, with differing information content – yet common theme is a **name** or, failing that, an **originating location**

Translation

Proteins
(>1,000,000)

List of names, quantity
Proteomics

Molecular
Interactions

Cellular
Function

Functional
End Products

Metabolites
(~2,500 small molecules)

List of names, quantity
Metabolomics

# Sequence Ontology – e.g. mapping experimental results

# *Conceptual Solution...*

**Gene expression data**

**ChIP-seq data**

**Protein-protein data**

**Metabolomic data**

Student Data

Cafeteria Data

Financial Data

Transportation Data

Idealized system ultimately integrates data products – provides means for researchers to access & query data space

**Data Warehouse**

Data Mart
**Reporting Tool**

ETL Process
**Extract, Transform and Load**

InterMine

# *Introducing InterMine...*

Open Source Data Warehouse



github.com/intermine

InterMine has three parts:

1. Database

    a. Load data from different data sets into **a single database**

2. Webapp

    a. Mine the data

    b. Visualise!

3. Web services

# InterMine - Backend

Originally developed by Gos Micklem's group for FlyMine (2007)

*Uses read-only **ObjectStore** methodology to implement fast queries on linked dataspace using Seq. Ont. schema.*

Web-enabled access also via RESTful Perl, Python, Java, Ruby, R APIs

Java/PostgreSQL core code base, Tomcat for web services

# *Main step - Constructing & Loading a Data Model (using SO)*



**Data Model for ObjectStore based on**
 - experimental data space
 - organism's genome(s)
 - literature references
 - published/archived 'omics data
 - Data loaded from parsers to XML InterMine format

# The InterMine `Science Gateway' web portal

# The InterMine `Science Gateway' web portal

# Using the QueryBuilder tool to explore dataspace

# *Using Web Services remotely...*

List of Identifiers → List of database records → Homologues in X →
Their protein domains

```
use InterMine::Webservice;

my $flymine = Webservice::InterMine->get_service('www.flymine.org/query', 'TOKEN');
my $list = $flymine->new_list(type => 'Gene', content => "some/file.txt");

my $query = $list->to_query;
$query->add_views('homologues.homologue.symbol');
$query->add_constraint('homologues.homologue.organism.shortName', '=', 'R norvegicus');
my $rat_symbols = [map {$_->{homologues.homologue.symbol}} $query->results("hashrefs")];

my $ratmine = Webservice::InterMine->get_service('ratmine.mcw.edu/ratmine');
my $rat_query = $ratmine->new_query(class => 'Gene');
$rat_query->add_views(qw/symbol primaryIdentifier proteins.proteinDomains.name/);
$rat_query->add_constraint('symbol', 'ONE OF', $rat_symbols);
$rat_query->show;
```

NUI Galway
OÉ Gaillimh

InterMine

# InterMines, Intermines...

## EXISTING MINES

A number of different data warehouses powered by InterMine already exist. These include:

FlyMine - *Drosophila* genomics

modMine - fly and worm modENCODE data

MouseMine - at MGI

RatMine - at RGD

WormMine - at WormBase

YeastMine - at SGD

ZebrafishMine - at ZFIN

INDIGOmine - microbes

ThaleMine - Araport Project with data for Arabidopsis thaliana

ChickpeaMine - Desi & Kabul chickpea

TargetMine - drug target discovery

MitoMiner - proteomic data for mitochondria

HumanMine - human

FlyTF.org - *Drosophila* transcription factors

PhytoMine - plants

MedicMine - *Medicago truncatula*

BovineMine - *Bos Taurus*

HymenopteraMine - Bees, Ants & Wasps

SoyMine - Soybase soy bean data

CHOMine - *Cricetulus griseus* and CHO cells

BeanMine - LegFed chado bean data

LegumeMine - String bean, Soy, and Peanut

PeanutMine - Peanut chado/GFF data

Shaare - Gene candidate prioritisation

PlanMine - Planarian flatworms

Wheat3BMine - Wheat chromosome 3B

GrapeMine - Grapevine

RepetDB - repetitive DNA elements

XenMine - Xenopus

TetraMine - *Tetrahymena thermophila*

NUI Galway
OÉ Gaillimh

InterMine

# *What about the DCI thing?*

InterMines' tend to be standalone `silos' based one one model organism/biological process.
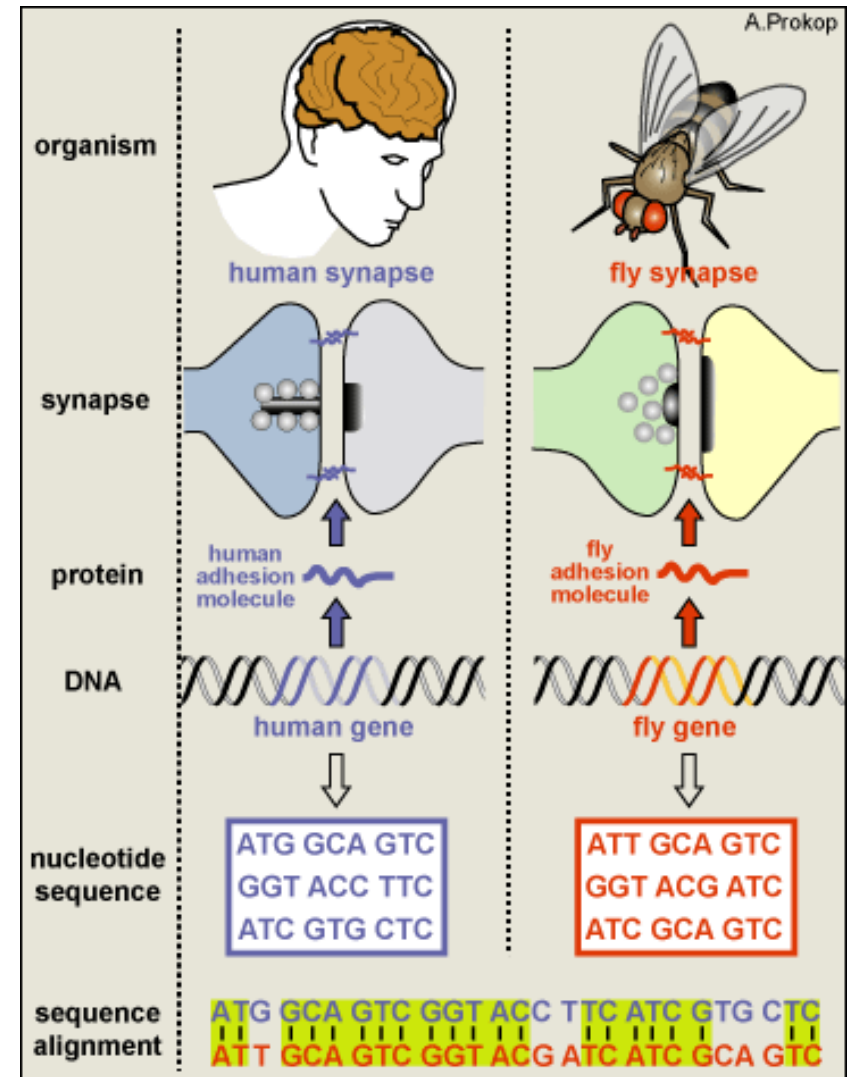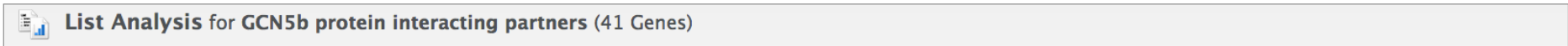
ToxoMine:
  DB server (8 core, 32GB RAM, 2 TB)
  Web servers (4 core, 16GB, 256 GB)

However – *fundamental homology between all living organisms enables interoperability 'for free' through the Sequence Ontology*.

Offers a Virtual Observatory like experience across registered Mines.



NUI Galway
OÉ Gaillimh

InterMine

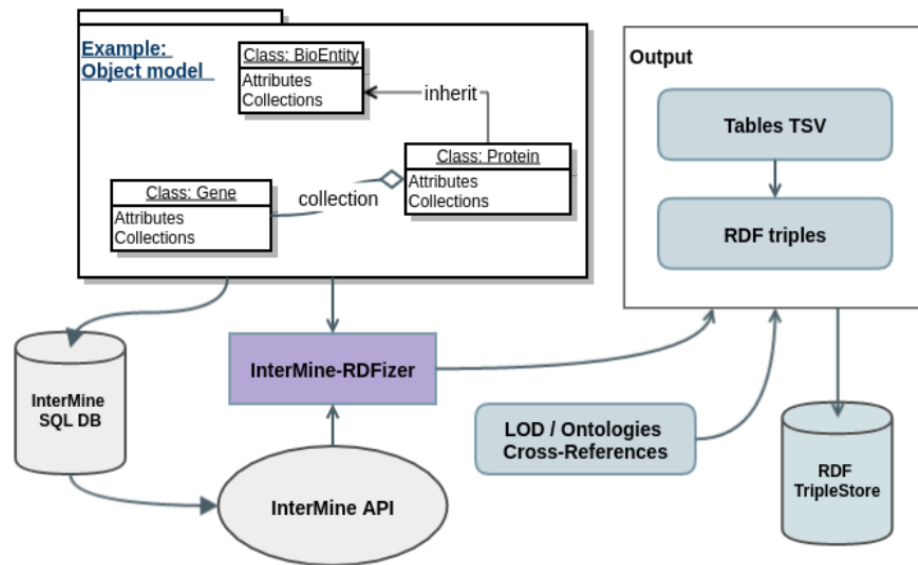# Cross-species query can be performed using the InterMOD infrastructure

# *Ongoing InterMine Community Activity*

## RDFization of Intermine

Create MO-LD (Model Organism Linked Data)

- improve interoperability
- enable federated queries across larger RDF space
- better packaging for local/cloud deployment

## Enchanced Webservice functionality

Evolve infrastructure to current state-of-the-art

- move away from Java/Tomcat paradigm
- modularity to integrate third-party web apps
- Clojure, Angular, PostgreSQL, Web Services

# Acknowledgements

NUI Galway
OÉ Gaillimh

InterMine