

The Taverna Workbench: Integrating and analysing biological and clinical data with computerised workflows

Dr Katy Wolstencroft

myGrid

University of Manchester

Vrije Universiteit, Amsterdam



- Why workflows are important
- WSDL, REST and other Workflow Services
- Getting started with Taverna
- Taverna in Use
- Sharing and reusing workflows
- Workflows on servers, grids and clouds
- Taverna Future Plans





Taverna



[Introduction](#) [Documentation](#) [Download](#) [Developers](#) [News](#) [Publications](#) [About](#)

Taverna Workflow Management System

Powerful, scalable, open source & domain independent tools for designing and executing workflows. Access to 3500+ resources.

RECENT NEWS

- October 8, 2012 Taverna Server 2.4.1 release
- September 18, 2012 Software Sustainability Institute Fellowships
- May 30, 2012 Taverna 2.4.1 patch

Get

Download for Windows,
Mac OS X or Linux

Use

Learn about the features &
functionality

Extend

Learn about the internals &
how to develop plugins

COMMUNITY

Next generation sequencing on Amazon cloud • Taverna-Galaxy integration • CDK plugin for cheminformatics
• Taverna 3 OSGi • SCUFL2 workflow bundle language • Taverna infrastructure VMs



Download, unpack and run



- 21st century is the century of information
- eGovernment
- World bank data
- Climate change data
- Large scale physics
 - Large Hadron collider
 - Astronomy
- 'Omics data
- Next Gen Sequencing



Where is the data?

- In repositories run by major service providers (e.g. NCBI, EBI)
- Group/Institute web sites
- On ftp servers
- In local project stores

- Few defined formats
- Inconsistent metadata



National Center for
Biotechnology Information (USA)



Tokyo, Japan



European Bioinformatics Institute

Cambridge, UK



PathPort
The Pathogen Portal Web Project



SRS



SeqHound

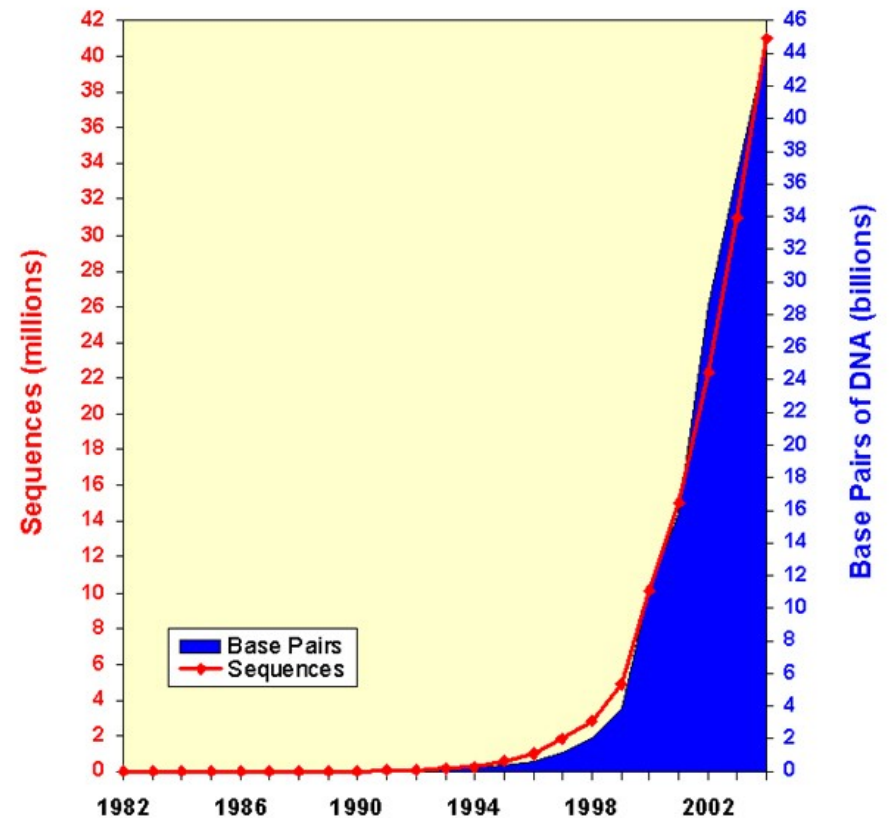


Lots of Resources

NAR 2012 – 1500 databases



Growth of GenBank
(1982 - 2004)



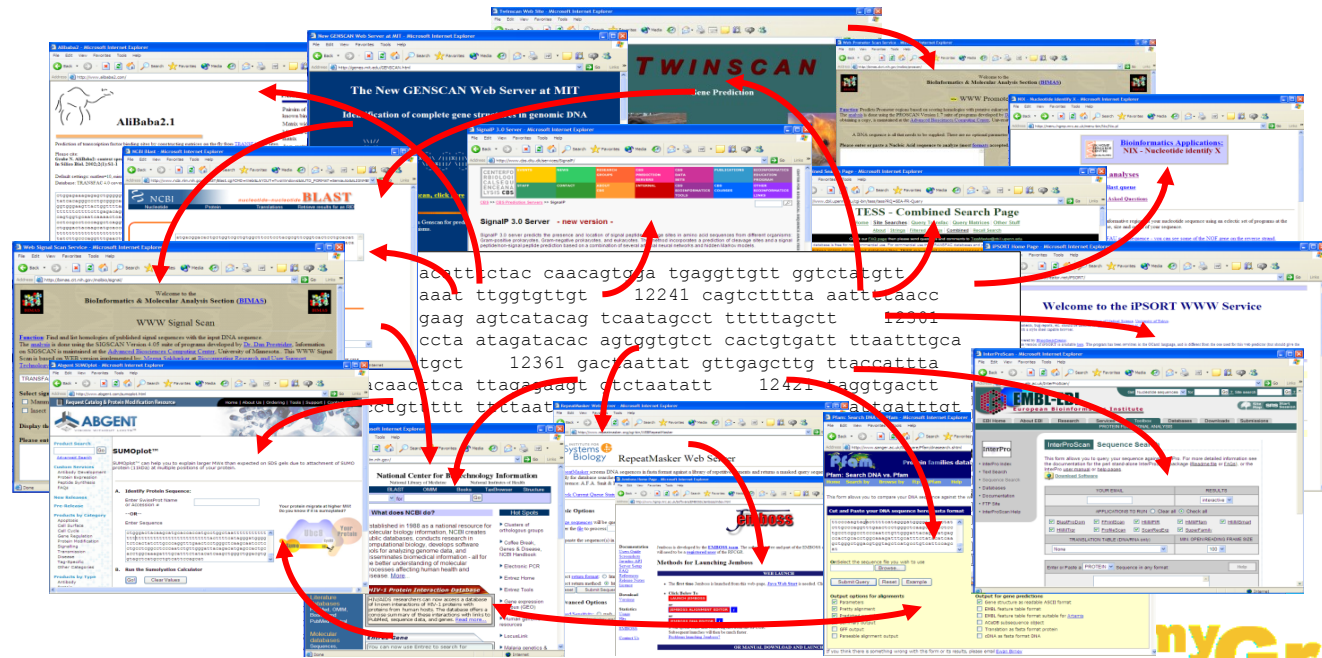
Distribution

- Data resources – databases, analysis tools
- Computational power – servers, clusters, cloud/grid
- Researchers and collaborators – skills and expertise need to be shared and exchanged
Analysis scripts need to be shared and exchanged



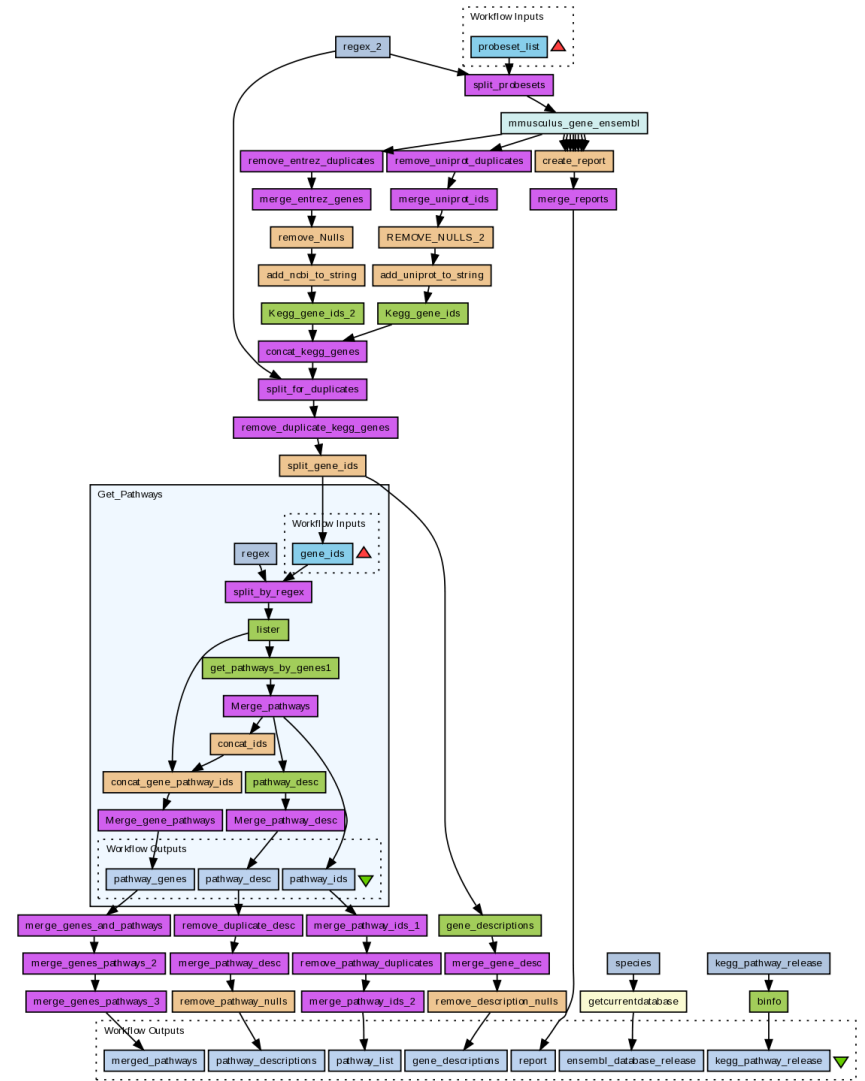
What that means for Bioinformatics

- Sequential use of distributed tools
- Incompatible input and output formats
- Analysis of large data sets by multiple researchers
- Difficult to record parameter selections
- Difficult to reproduce analyses



Workflow as a Solution

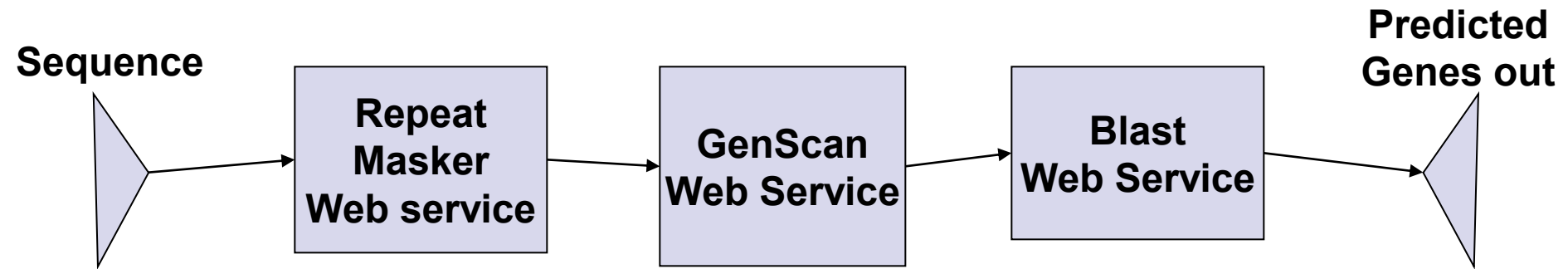
- Automating the process
- Sophisticated analysis pipelines
- A set of **services** to analyse or manage data (either local or remote)
- Data flow through services
- Control of service invocation
- Iteration



What is a Workflow?

Describes *what* you want to do, rather than *how* you want to do it

Simple language specifies how processes fit together



Workflows are ideal for...

- High throughput analysis
 - Transcriptomics, proteomics, next gen sequencing
- Data integration, data interoperation
- Data management
 - Model construction
 - Data format manipulation
 - Database population
 - Semantic integration
 - Visualisation



Promoting Reproducible Research

Informatics involves

- Complex, multi-step analyses
- Lots of data as inputs
- Lots of data generated
- Workflows encapsulate the methods and parameters
- Workflows allow you to visualise the methods



Preventing Irreproducible Research

- An array of errors

<http://www.economist.com/node/21528593>

- Duke University, 2006 -Prediction of the course of a patient's lung cancer using expression arrays and recommendations on different chemotherapies from cell cultures
- 3 different groups could not reproduce the results and uncovered mistakes in the original work



If the Analyses were done using Workflows.....

- Reviewers could re-run experiments and see results for themselves
- Methods could be properly examined and criticised
- Mistakes could be pinpointed



Workflows are ...

- ... records and protocols (i.e. your *in silico* experimental method)
- ... know-how and intellectual property
- ... hard work to develop and get right
-re-usable methods (i.e. you can build on the work of others)

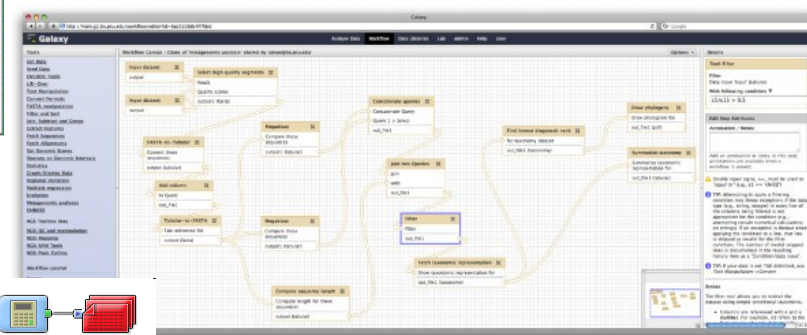
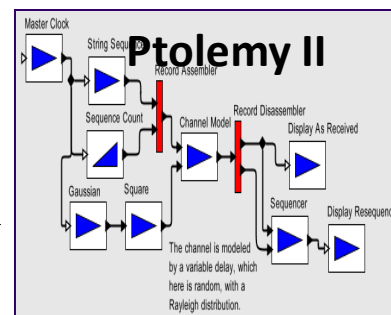
So why not share and re-use them

my experiment



WORKFLOW SYSTEMS



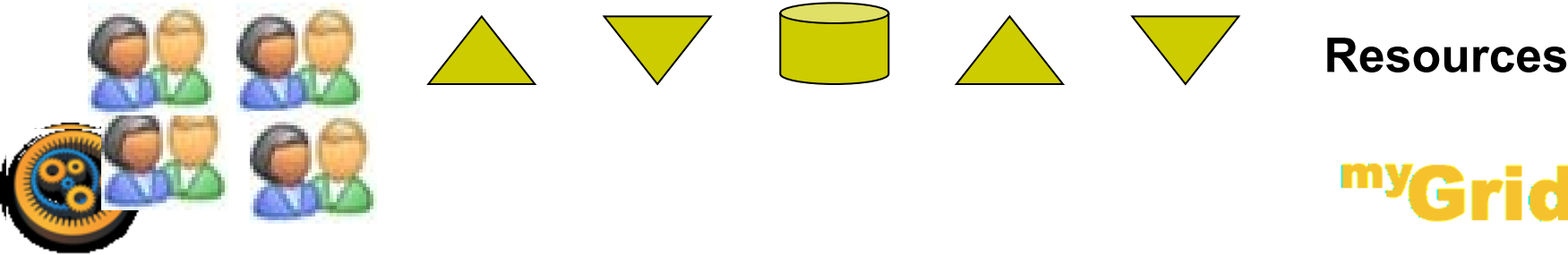
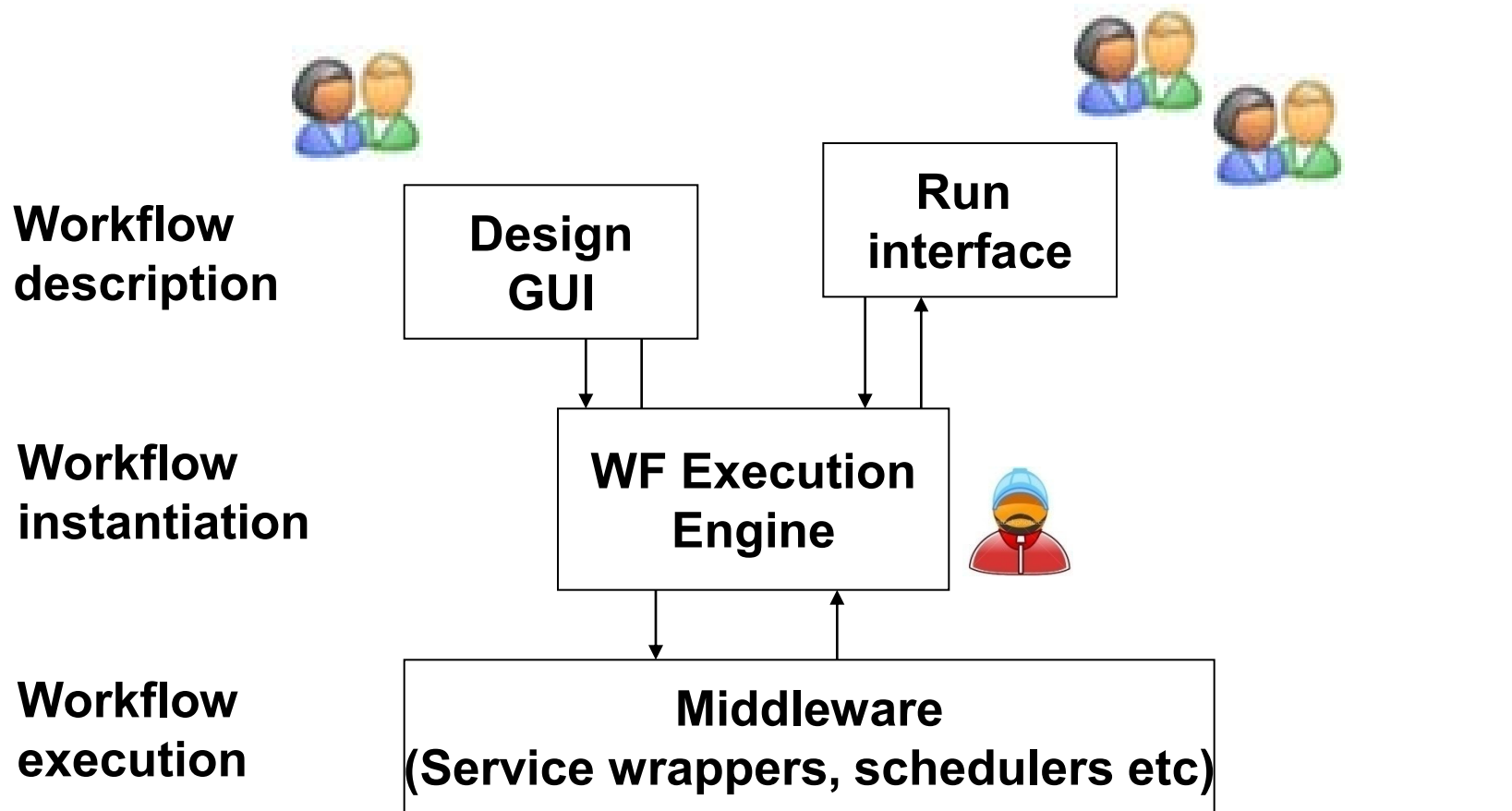


Galaxy

myGrid



All Workflow Systems at 50,000 feet



Different Types of Workflows

- Sequences of concatenated steps
- Two types of workflows:

- Data workflows

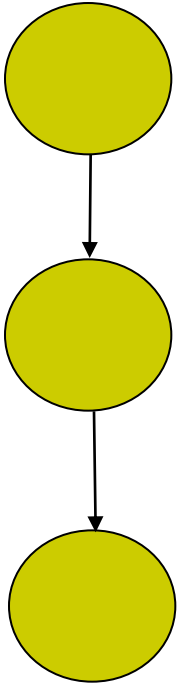
A task is invoked once its **expected data** has been received.
When complete, it passes any resulting data downstream

- Control workflows

A task is invoked once its **dependant tasks** have been completed

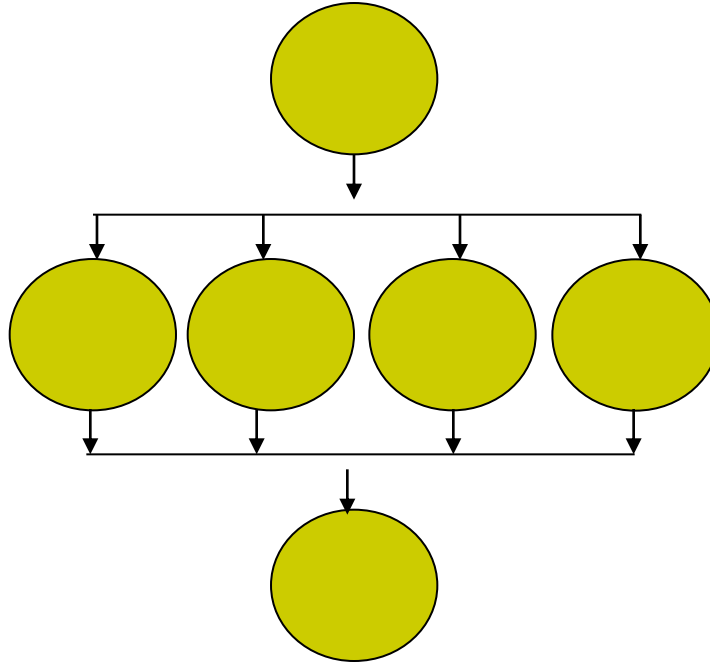


Possible Workflow Structures



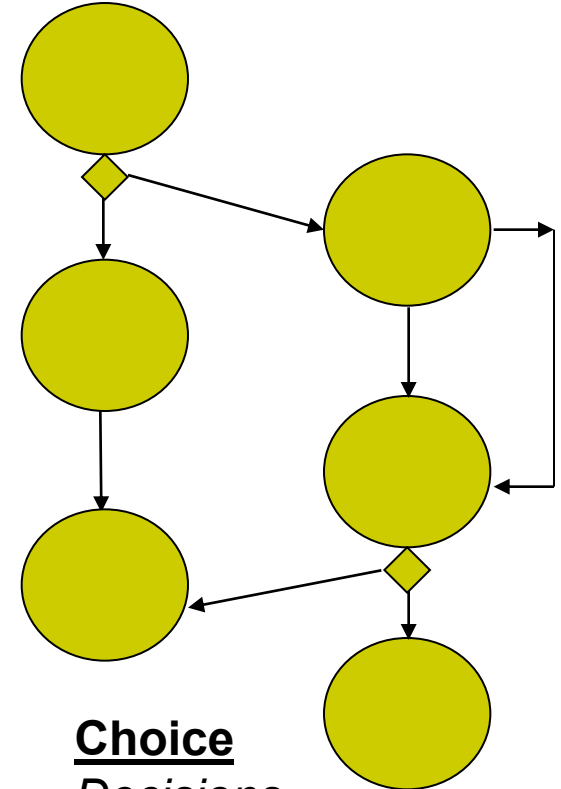
Sequence

Store intermediate results



Parallel

Apply multiple components to a set of data



Choice

Decisions at runtime

Iteration

Loop through datasets



Taverna Workbench

<http://www.taverna.org.uk/>

Freely available
open source
Current Version 2.4

80,000+ downloads
across version

Part of the myGrid Toolkit

Windows/Mac OS X/
Linux/unix



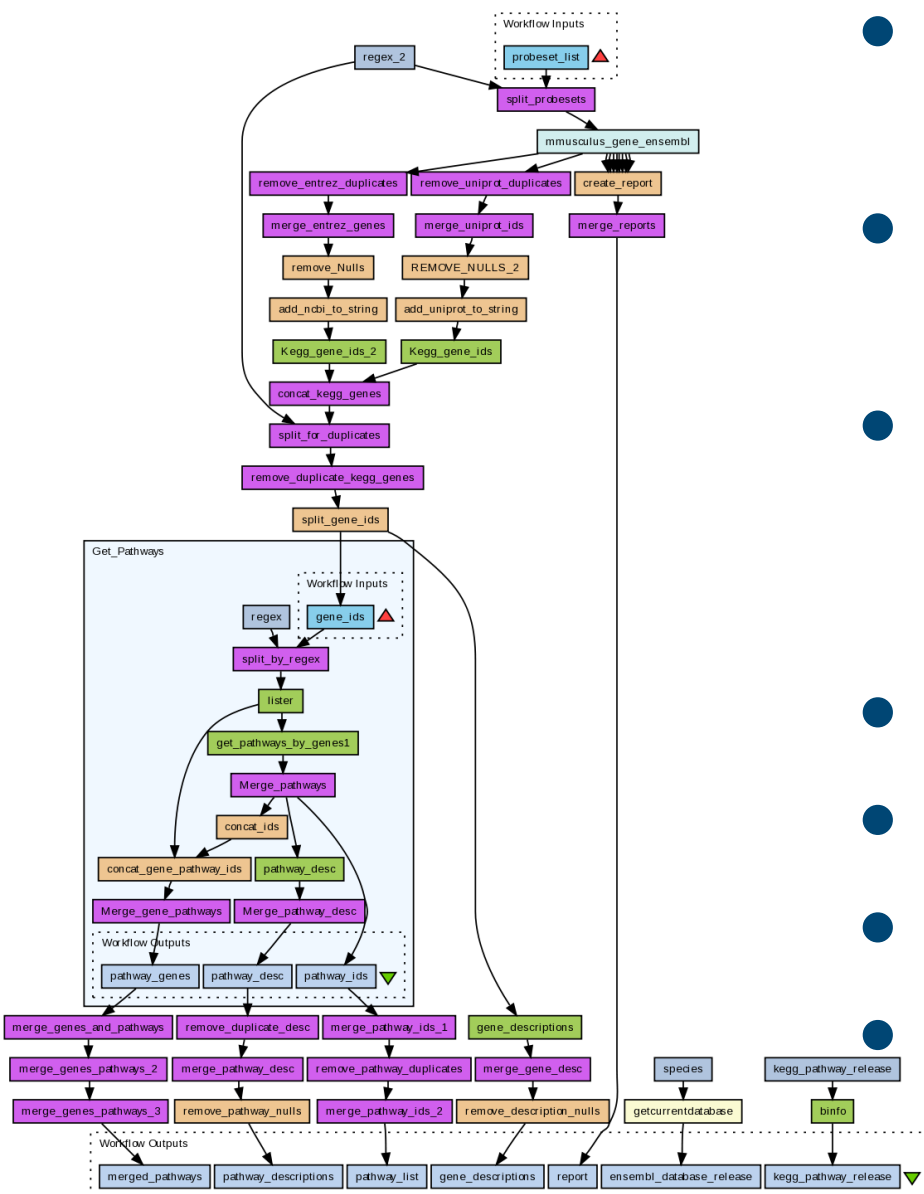
The screenshot shows the Taverna website homepage. At the top is the Taverna logo (a gear with a sun-like center) and the word "Taverna". To the right is the myGrid logo and a search bar. Below the header is a navigation menu with links: Introduction, Documentation, Download, Developers, News, Publications, and About. The main content area features a large banner with the text "Taverna Workflow Management System" and a description: "Powerful, scalable, open source & domain independent tools for designing and executing workflows. Access to 3500+ resources." Below this banner are three buttons: "Get" (Download for Windows, Mac OS X or Linux), "Use" (Learn about the features & functionality), and "Extend" (Learn about the internals & how to develop plugins). To the right of the banner is a "RECENT NEWS" section with three items: "February 15, 2011 Opal plugin for Taverna 2.2", "January 27, 2011 BitesizeBio Webinar on Taverna, myExperiment and BioCatalogue", and "January 11, 2011 PDF and HTML". Below the banner is an "IN PRESS" section with four items: "Taverna 2.3", "Taverna 3 Next Generation", "SCUFL2 workflow bundle language", and "Taverna infrastructure VMs". The bottom section is titled "Taverna" and contains three paragraphs of text. The first paragraph describes Taverna as an open source and domain independent Workflow Management System. The second paragraph mentions that Taverna has been created by the myGrid team and funded through the OMII-UK. The third paragraph describes the Taverna suite, including the Taverna Engine, Taverna Workbench, Taverna Server, and Command Line Tool. To the right of the text is a video player titled "See Taverna 2.2 in action" showing a workflow diagram and a play button.

Nucleic Acids Res. 2006 Jul 1;34(Web Server issue):W729-32.
Taverna: a tool for building and running workflows of services.
Hull D, Wolstencroft K, Stevens R, Goble C, Pocock MR, Li P, Oinn T.



Taverna Workflows

- Part of UK E-Science myGrid project
- Started in 2001, collaboration across UK
- Now: Manchester (Goble), Oxford/Southampton (DeRoure)
- <http://www.taverna.org.uk>
- Local Taverna desktop
- Taverna Server
- Taverna on the cloud



Open source, open development

- Taverna suite of tools are all **open source**, free to use and **customise**
- Large user **community**, active mailing lists
- Lead developers: **myGrid** in Manchester UK
- **Contributors** from across the world
- **Plugins** developed and shared by contributors
 - XPath, REST, R, BioCatalogue, PBS, SADI, External Tools (UseCase), UNICORE, CDK, Opal, caGrid, XWS, gLite



Taverna Workbench

List of services

Workflow engine to run workflows

Construct and visualise workflows

Web Services

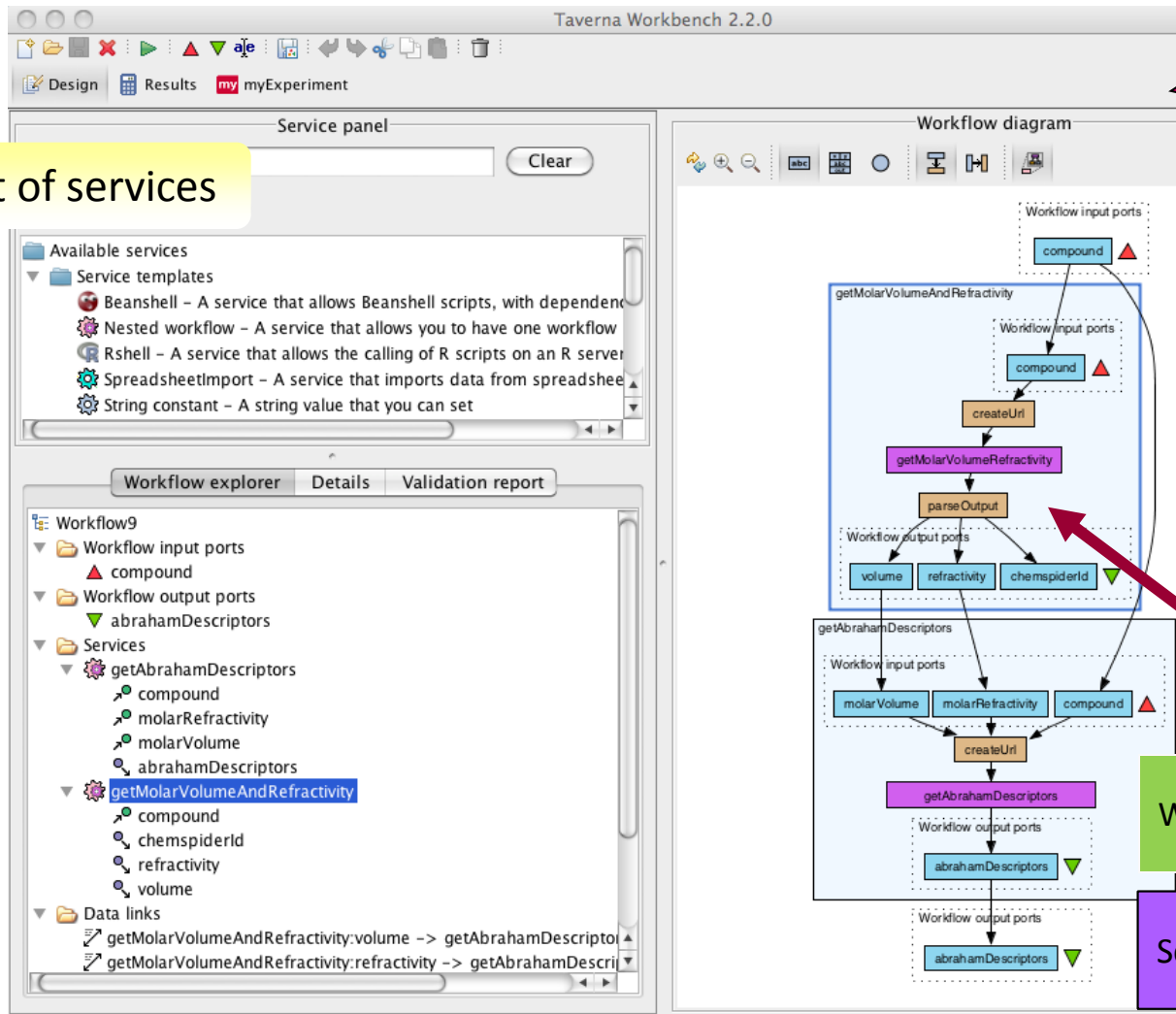
e.g. KEGG

Scripts

e.g. beanshell, R

Programming libraries

e.g. libSBML



Workflows and the in Silico Life Cycle

Create and run workflows



Workflows and the in Silico Life Cycle

BioCatalogue 

Discover,
understand and
assess services

Create and run workflows



Workflows and the in Silico Life Cycle

BioCatalogue 

Discover,
understand and
assess services

Discover, reuse and
share workflows

 **myexperiment**

Create and run workflows

 **Taverna**



Workflows and the in Silico Life Cycle

BioCatalogue 

Discover,
understand and
assess services

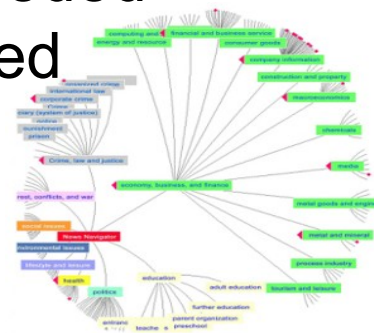
Discover, reuse and
share workflows

 **myexperiment**

Create and run workflows

 **Taverna**

Manage the
metadata needed RDF, OWL
and generated



 **myGrid**



SERVICES IN WORKFLOWS



What are Web Services?

NOT the same as services on the web (i.e. web forms)

Web services support machine-to-machine interaction over a network

Therefore, you can automatically connect to and use remote services from your computer in an automated way



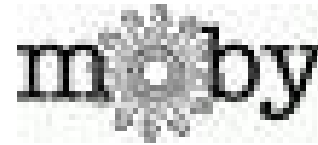
Web Services – Brief Glossary

- WSDL (Web Service Definition Language)
 - A machine-readable description of the operations supported
- SOAP (Simple Object Access Protocol)
 - An xml protocol for passing messages
- REST (Representational State Transfer)
 - An alternative interface to SOAP



Using Remote Tools and Services with Taverna

- Web Services
 - WSDL
 - REST
- Grid Services
- Local services
- Beanshell (small, local scripts)
- Secure Services
- Workflows
- BioMart
- R-processor
- And more.....



Specialist services

BioMart Queries

- Federated database system that provides unified access to distributed data sources
- Ensembl, Pride.....

Dataset: Rattus norvegicus genes (RGSC3.4)

Filters: Chromosome: X

Attributes: Ensembl Gene ID, Ensembl Transcript ID

Export all results to: File (Go) | TSV | Unique results only

View: 10 rows as HTML | Unique results only

Ensembl Gene ID	Ensembl Transcript ID	RGD ID	SO term
ENSRNOG00000000007	ENSRNOT00000000000	628617	alpha-Thalassemia
ENSRNOG00000000003	ENSRNOT00000000004	628618	Mental Retardation
ENSRNOG00000000000	ENSRNOT00000000001	628619	Dysagammaglobulinemia
ENSRNOG00000000001	ENSRNOT00000000002	628620	
ENSRNOG00000000004	ENSRNOT00000000005	628621	
ENSRNOG00000000000	ENSRNOT00000000003	628622	

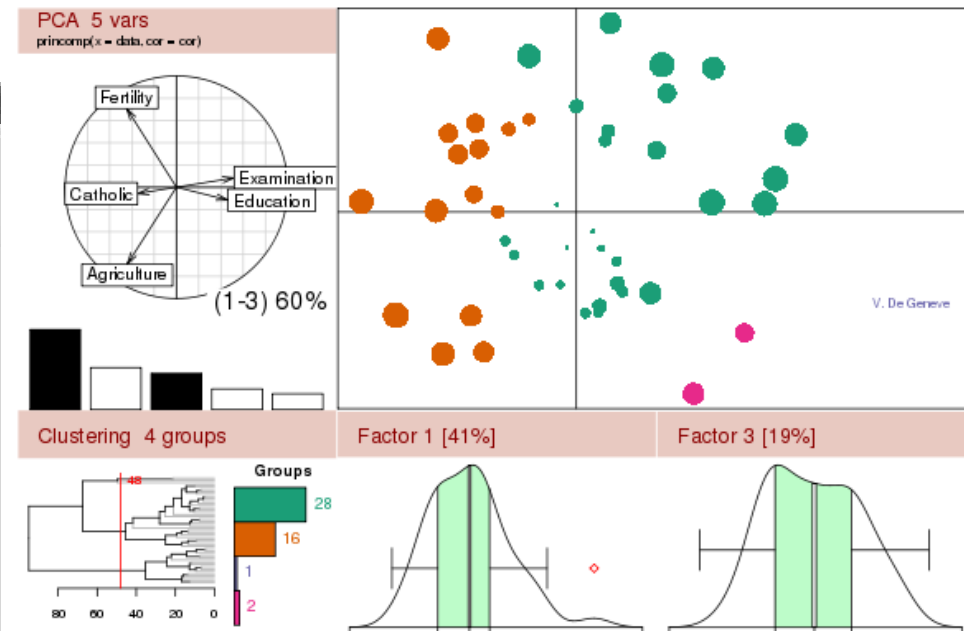
Dataset 623 / 36367 Entries

Filters: Biological Process: response to stress

Attributes: RGD ID, SO term

R-scripts

- R is a free software environment for statistical computing and graphics



Different Approaches to Service Connections

- Open – connect to ANY service regardless of type and structure
 - More services, but more heterogeneity
 - Easy to add new services
 - Taverna, Kepler
- Closed – connect to services designed specifically to work together,
 - Less heterogeneity, but fewer services
 - Harder to add new services
 - Galaxy server, Knime



Who Provides the Services?

Open domain services and resources

- Taverna accesses thousands of services
- Third party – we don't own them – we didn't build them
- All the major providers
 - NCBI, DDBJ, EBI ...
- Enforce NO common data model.



National Center for
Biotechnology Information (USA)



Tokyo, Japan



EMBL-EBI

European Bioinformatics Institute



Cambridge, UK



SRS

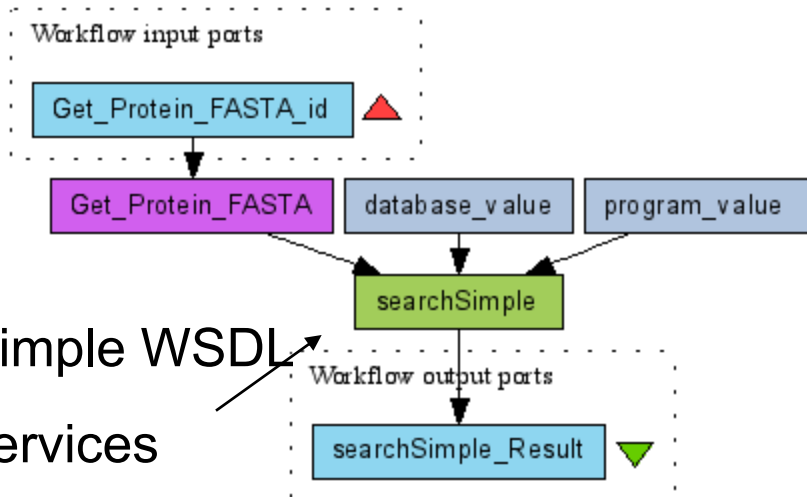


SeqHound

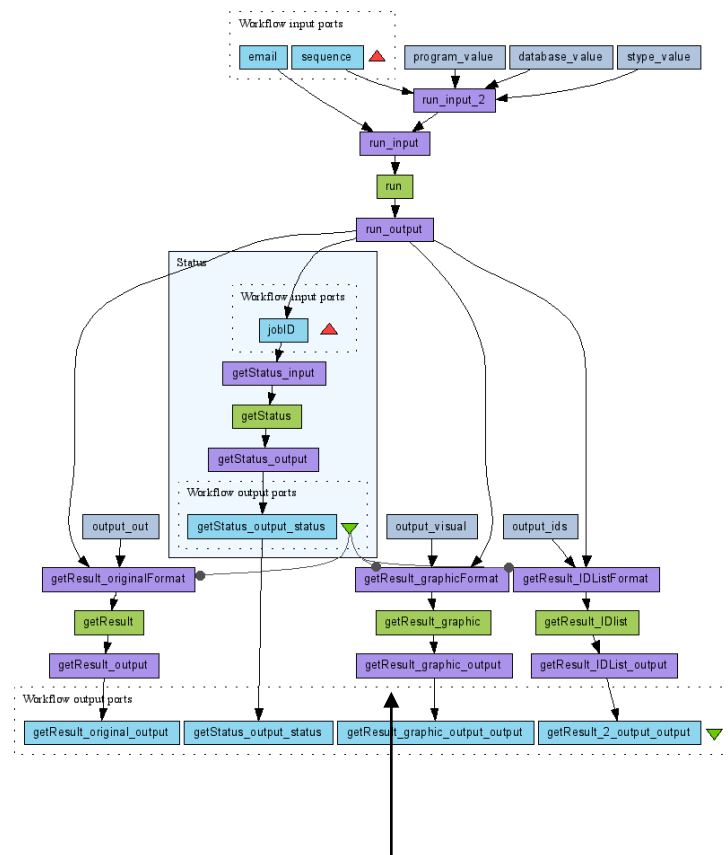
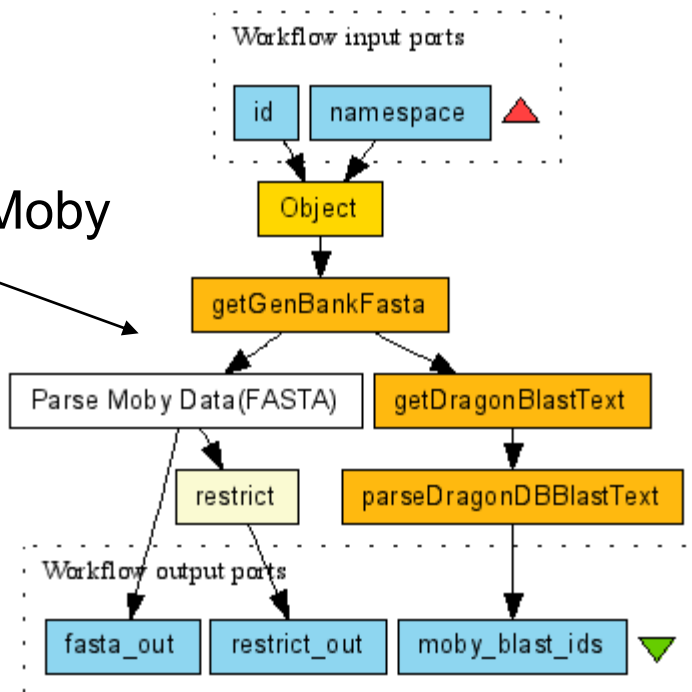


How do you use the services?

Simple WSDL
services



SADI / BioMoby
'Semantic'
Services



Asynchronous services



Managing Heterogeneities

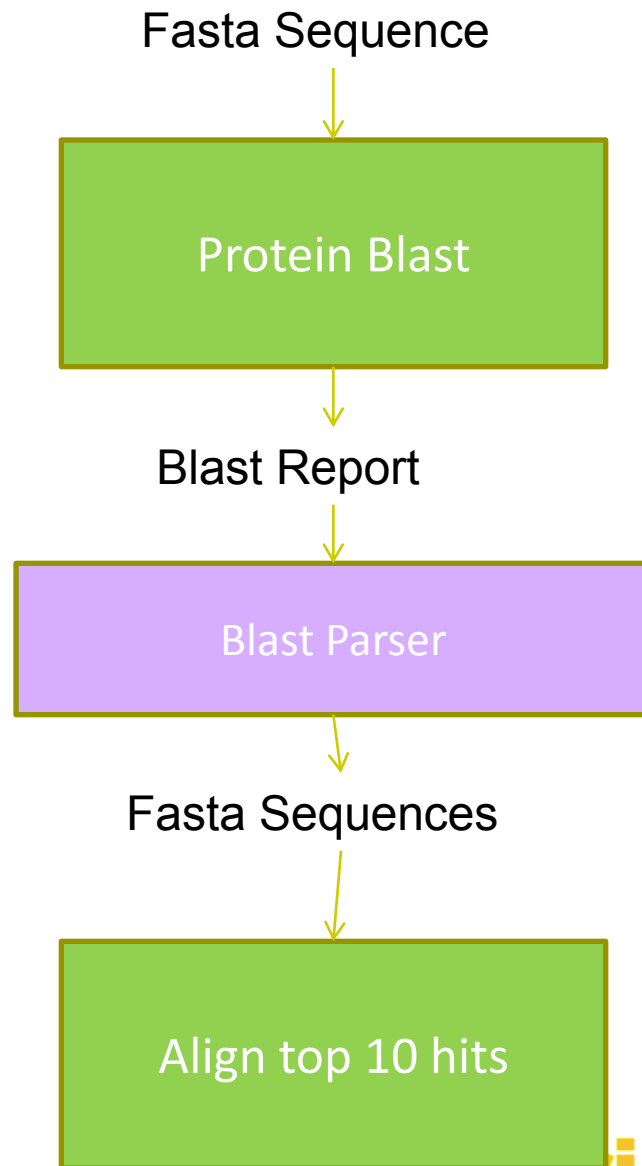
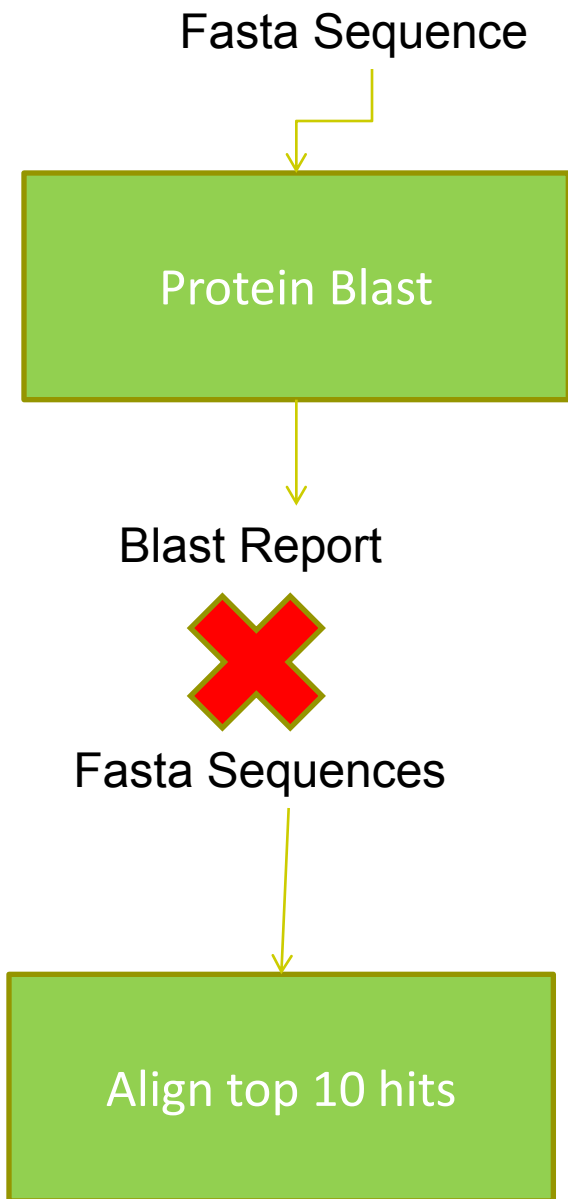
1. Understand how services work – inputs, outputs, dependencies → service descriptions and documentation
2. Find and use SHIM (or helper) services to combat incompatibilities

A Shim Service is a service that:

- doesn't perform an experimental function, but acts as a connector, or glue, when 2 experimental services have incompatible outputs and inputs



Shim Example












Understanding how services work

The BioCatalogue: providing a curated catalogue of Life Science Web Services

The BioCatalogue currently has **1730 services**, **130 service providers** and **445 members**

Latest Activity

Last 7 days

-  **vasun joined** the BioCatalogue
-  **Peter Taschner added** a publication annotation to the Soap Service of Service: [MutalyzerService](#)
-  **Peter Taschner added** a contact annotation to the Service Deployment of Service: [MutalyzerService](#)
-  **Peter Taschner added** a tag annotation to Service: [MutalyzerService](#)
-  **Peter Taschner added** a tag annotation to Service: [MutalyzerService](#)
-  **Peter Taschner added** a tag annotation to Service: [MutalyzerService](#)
-  **Peter Taschner added** an alternative name annotation to Service: [MutalyzerService](#)
-  **Peter Taschner added** a documentation url annotation to the Soap Service of Service: [MutalyzerService](#)
-  **Peter Taschner added** a description annotation to the Soap Service of Service:

"Web Services are hard to find"

DISCOVER

- Find the right Web Service
- Powerful search and filtering
- Information from providers and community

[More info](#)

"My Web Services are not visible"

REGISTER

- Easily register Web Services
- Instantly available to everyone
- Providers can advertise, describe and monitor their Services

[More info](#)

"Web Services are poorly described"

ANNOTATE

- Anyone can describe and annotate
- Ongoing expert curation
- Social curation by the community

[More info](#)

"Web Services are volatile"

MONITOR

- Services change and get outdated
- BioCatalogue monitors Services
- Monitors availability and reliability

[More info](#)

Site Announcements

Have your say about BioCatalogue by taking part in the BioCatalogue users's survey

By [Franck Tanoh](#) (4 days ago)

BioCatalogue Maintenance - 7 December 2010 @ 9:30 am (GMT)

By [Eric E. Nzuobontane](#) (6 days ago)

BioCatalogue iPhone and iPad app now available to download for free

By [Franck Tanoh](#) (2 months ago)

The BioCatalogue Functional Unit paper presented at IEEE 2010 Fourth International Workshop on Scientific Workflows

By [Franck Tanoh](#) (2 months ago)

The National Cancer Research Institute (NCRI) joins forces with the BioCatalogue

By [Franck Tanoh](#) (4 months ago)

[More](#)

Our Partners



The EMBRACE Registry and the BioCatalogue have now been merged

Latest Services

[MutalyzerService](#)

[dbfetch](#)

[graphtools](#)

[PRANK \(REST\)](#)

[FASTM \(REST\)](#)



SHARE

REST

99 0

aka **InterProScan**

Categories: Function Prediction

Monitoring

REST Endpoints (6)

Monitoring

News

European Bioinformatics Institute (EBI)

Provider

UNITED KINGDOM 

Submitter

 Hamish McWilliam (about 1 month ago)

<http://www.ebi.ac.uk/Tools/services/rest/iprscan>

Service Description

Tags

Tags (19)

bioinformatics	ebi	embl-ebi	Gene3D
HAMAP	interpro	interproscan	Panther
Pfam	PIRSF	PRINTS	ProDom
ProSite	protein domain	protein family	
protein function	SMART	SUPERFAMILY	
	TIGRFAMs		

 [Login to add tags](#)

 Favourited By (0)

Managing Changes to Services

- Monitoring detects changes, but the community site can notify users about changes → advanced warning
- EBI – Soaplab EMBOSS tools discontinued Feb 13
 - Redirect to alternative services (also from EBI)
 - KEGG – SOAP services discontinued December 12
 - Replacing with equivalent REST services
 - Help identify equivalent or similar services



GETTING STARTED WITH TAVERNA: DEMO



Enrichment Analysis

- Many experiments result in a list of genes (e.g. microarray analysis, Chip-Seq, SNP identification etc)
- Today, we will use Taverna to perform enrichment analyses on a list of genes
 - We will enrich our dataset by discovering:
 1. Which pathways our genes are involved in and visualising those pathways
 2. The functions of the genes using Gene Ontology annotations



TAVERNA IN USE



What do Scientists use Taverna for?

Systems biology model building

Sequence analysis Protein structure prediction

Gene/protein annotation Microarray data analysis

Phylogeny Model simulations sweeps **Astronomy**

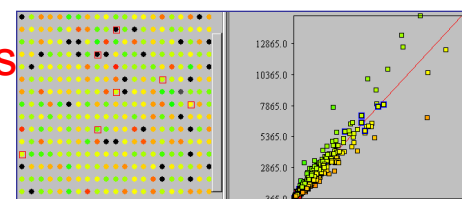
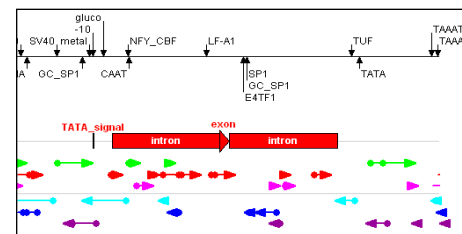
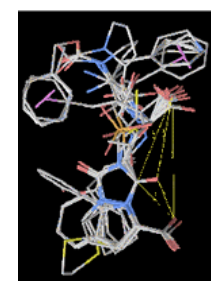
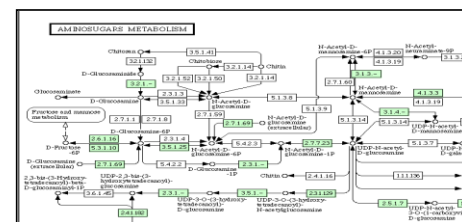
High throughput screening Proteomics **Music**

Phenotypical studies Text mining **Meteorology**

Public Health care epidemiology **Social Science**

Medical image analysis QTL studies **Cheminformatics**

QSAR studies Genome Wide Association Studies



Taverna for Omics

Functional Genomics

<http://www.myexperiment.org/workflows/126>

Publication: Solutions for data integration in functional genomics: a critical assessment and case study.

Smedley, Swertz and Wolstencroft, et al Briefings in Bioinformatics. 2008 Nov;9(6):532-44.

Genotype to Phenotype

<http://www.myexperiment.org/workflows/16>

Publication: A systematic strategy for large-scale analysis of genotype phenotype correlations: identification of candidate genes involved in African trypanosomiasis. Fisher et al Nucleic Acids Res. 2007;35(16):5625-33

Next Generation Sequencing

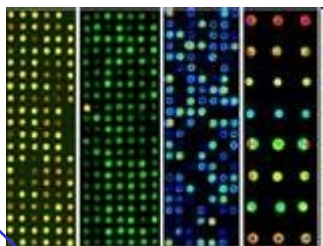
- Whole Genome SNP analysis of different cattle species in response to trypanosomiasis infection (sleeping sickness)
- Large data processing strategies
- Taverna in the cloud – deploying and running large data processes using cloud computing services



Lymphoma Prediction Workflow

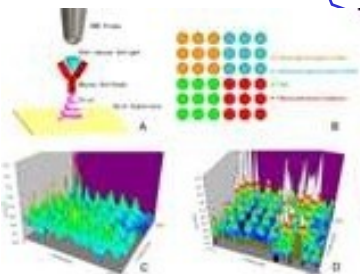


MicroArray from
tumor tissue



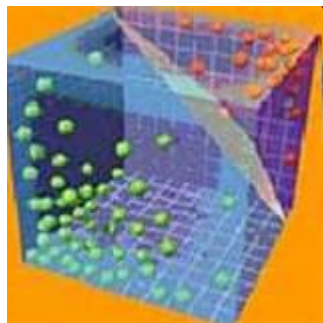
caArray

Microarray
preprocessing

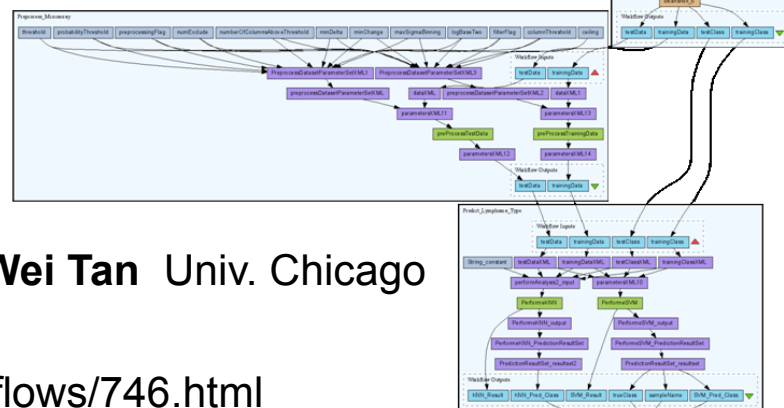


Lymphoma
prediction

GenePattern



Use **gene-expression** patterns associated with two lymphoma types to predict the type of an unknown sample.



Wei Tan Univ. Chicago

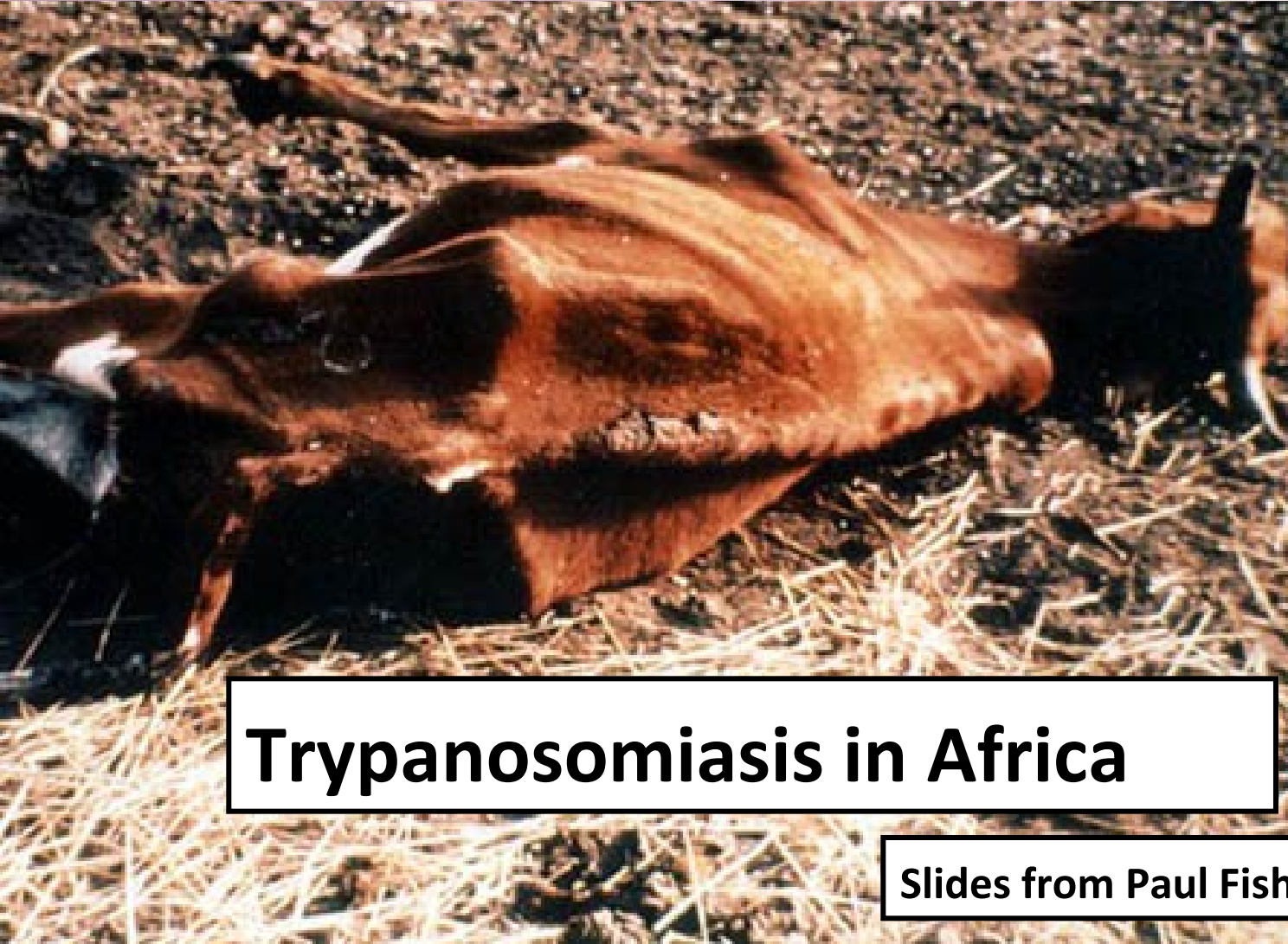
Wei Tan: <http://www.myexperiment.org/workflows/746.html>

Ack. Juli Klemm, Xiaopeng Bian, Rashmi Srinivasa (NCI)

Jared Nedzel (MIT)



The Wellcome Trust Funded Host-Pathogen Project



Trypanosomiasis in Africa

Slides from Paul Fisher



Steve Kemp



Andy Brass



Paul Fisher

<http://www.genomics.liv.ac.uk/tryps/trypsindex.html>

myGrid



Cattle Disease Research

\$4 billion US

Different breeds of African Cattle

- Some resistant
- Some susceptible

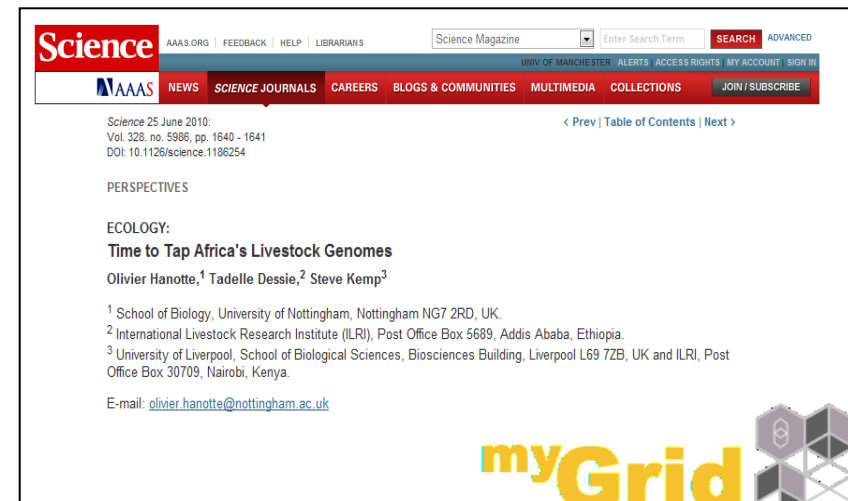
African Livestock adaptations:

- More productive
- Increases disease resistance
- Selection of traits

Potential outcomes:

- Food security
- Understanding resistance
- Understanding environmental
- Understanding diversity

<http://www.bbc.co.uk/news/10403254>



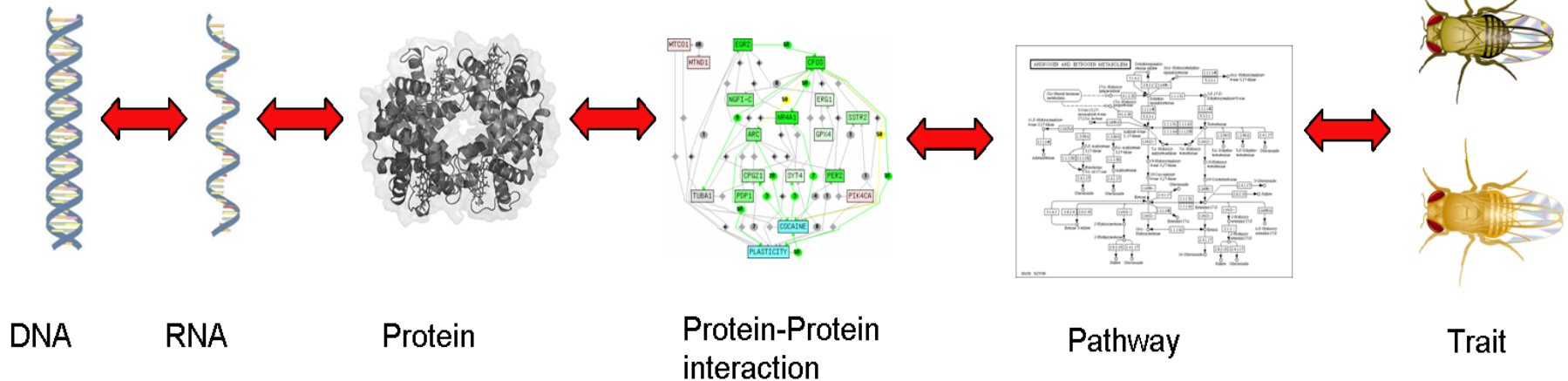
myGrid



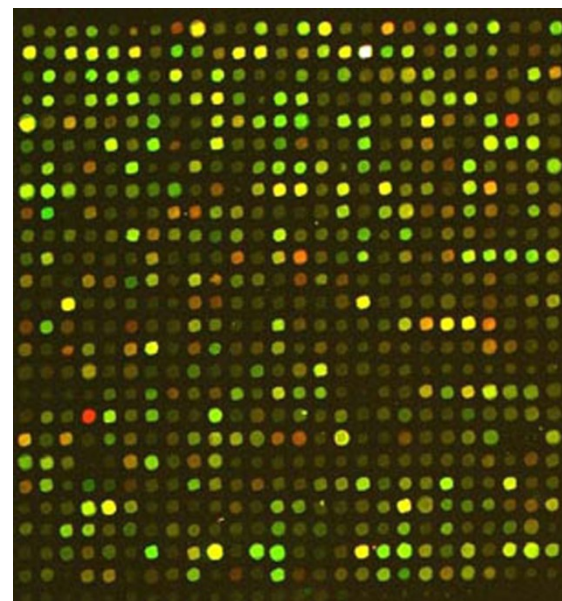
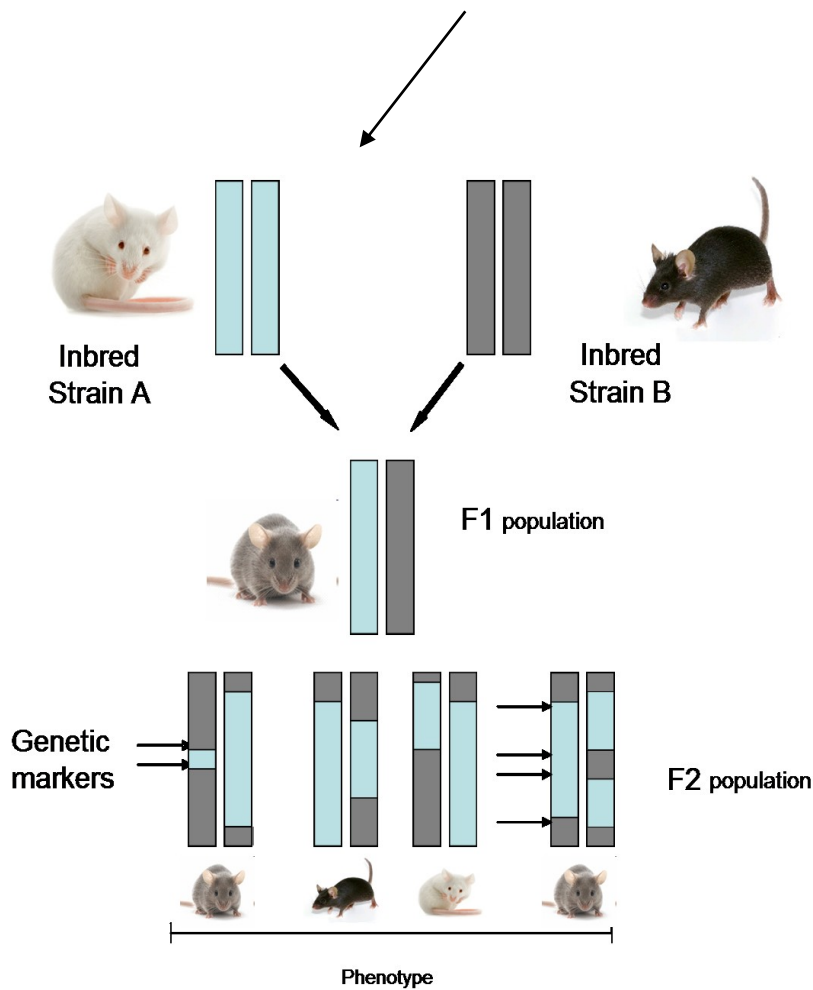
Understanding the process: Genotype - Phenotype

Genotype

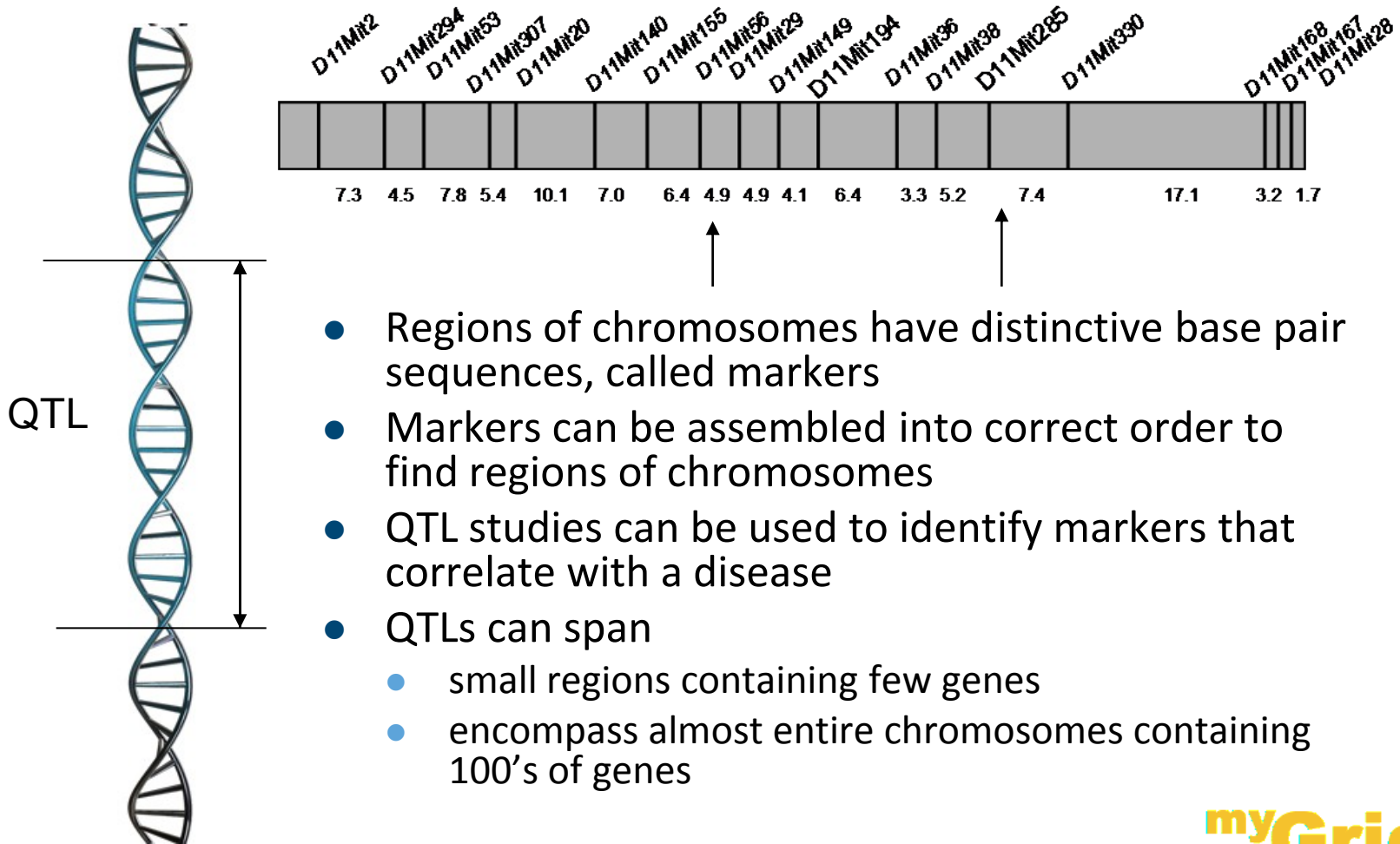
Phenotype



QTL + Microarrays

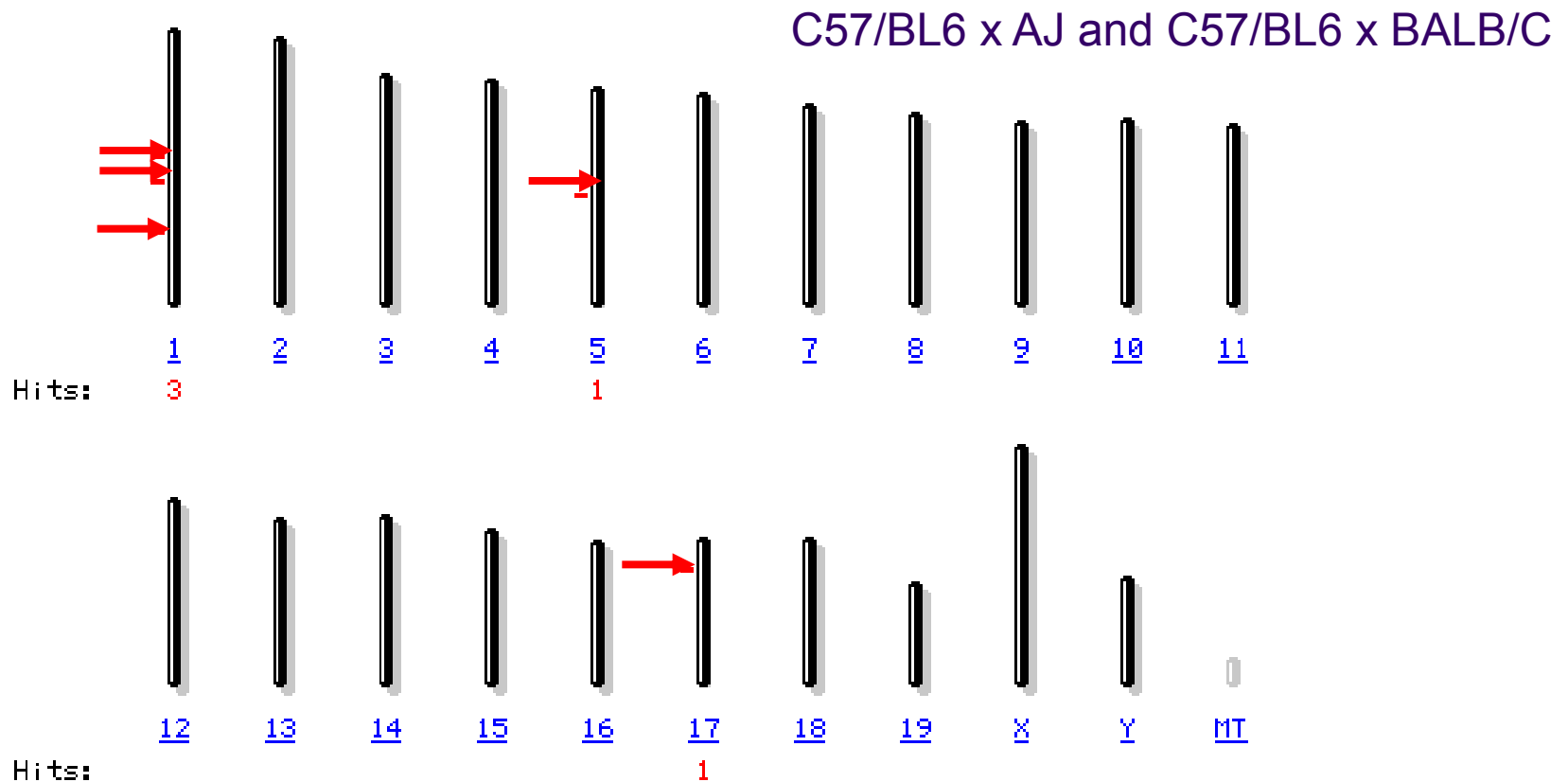


Quantitative Trait Loci (QTL)



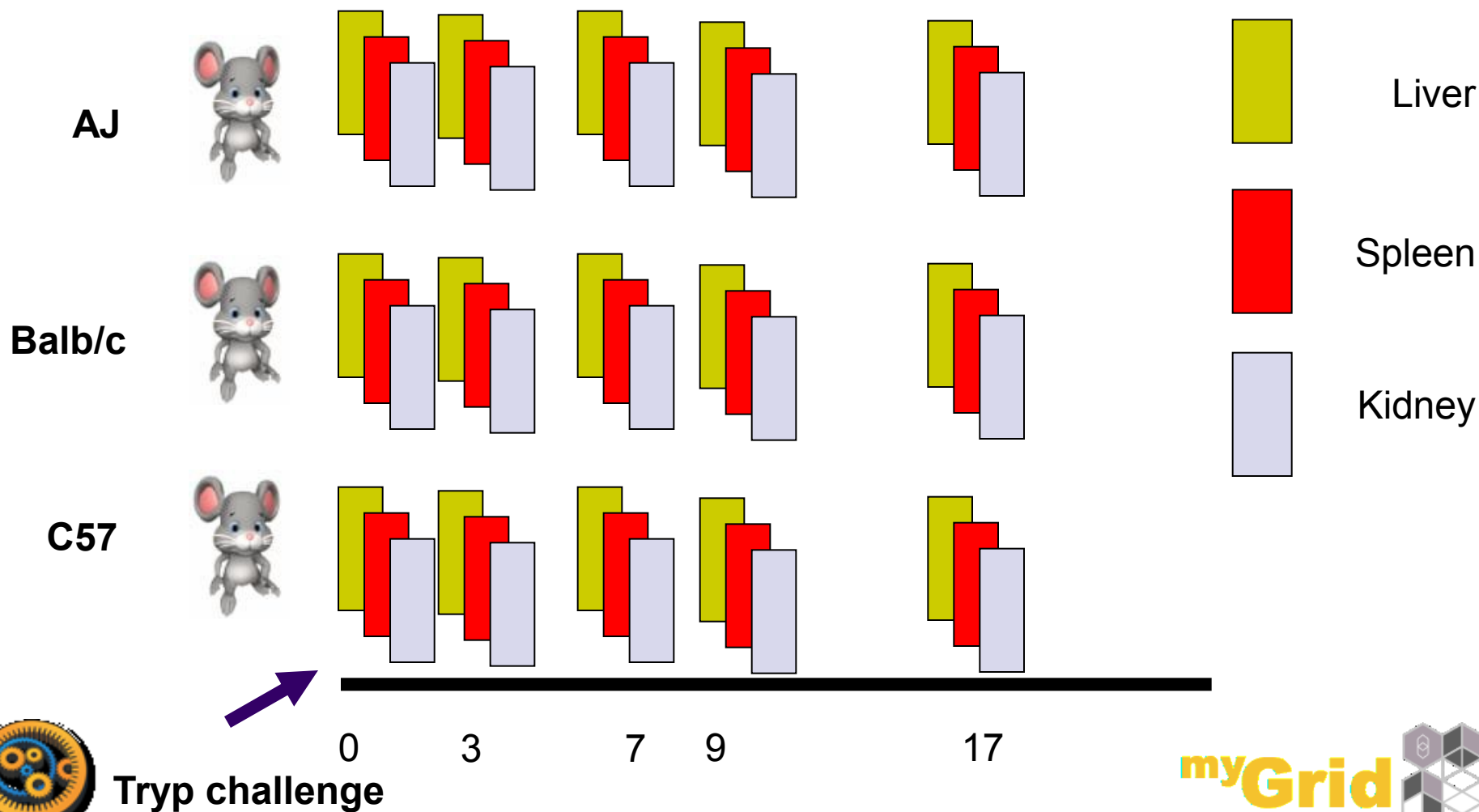
- Regions of chromosomes have distinctive base pair sequences, called markers
- Markers can be assembled into correct order to find regions of chromosomes
- QTL studies can be used to identify markers that correlate with a disease
- QTLs can span
 - small regions containing few genes
 - encompass almost entire chromosomes containing 100's of genes

Trypanosoma infection response (Tir) QTL

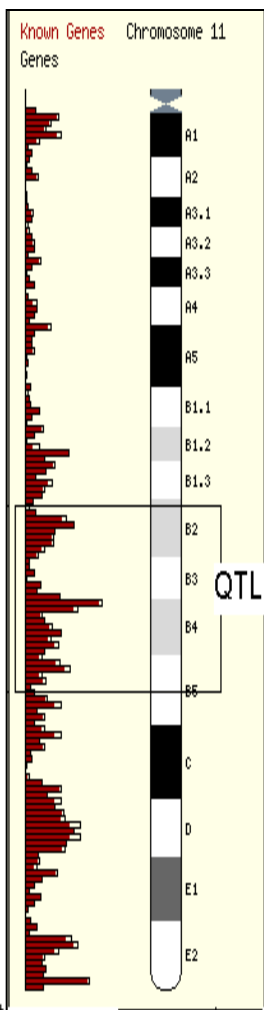


The experiment

A total of 225 microarrays

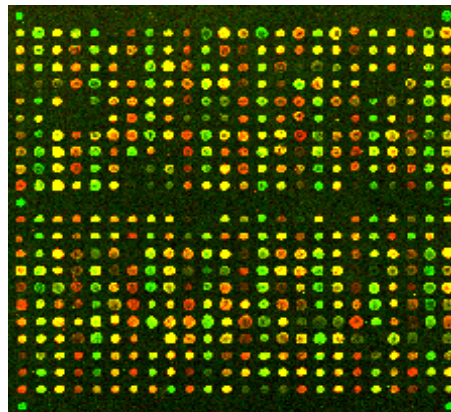


Huge amounts of data



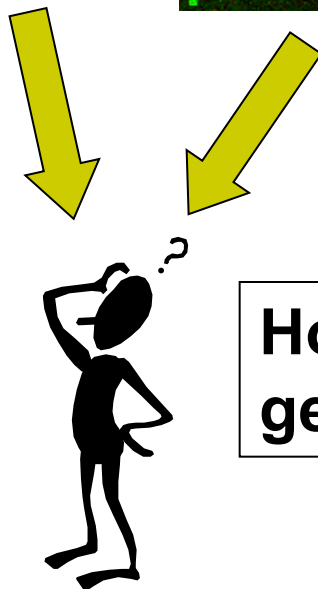
QTL region on
chromosome

200+ Genes



Microarray

1000+ Genes



**How do I look at ALL the
genes systematically?**



Phenotype

myGrid 

Microarray + QTL

Data analysis

- Identify pathways that have differentially expressed genes (from microarray studies)
- Identify pathways from Quantitative Trait genes (QTg)
- Track genes through pathways that are suspected of being involved in resistance/susceptibility



Trypanosomiasis Resistance Results

- DAXX gene identified in the workflows
- **Daxx** gene not found using manual investigation methods
- Sequencing of the Daxx gene in **Wet Lab** (at Liverpool) showed mutations that are thought to change the structure of the protein
- These mutations were also published in scientific literature, noting its effect on the **binding of Daxx protein to p53 protein**
- **p53 plays direct role in cell death and apoptosis, one of the Trypanosomiasis phenotypes**



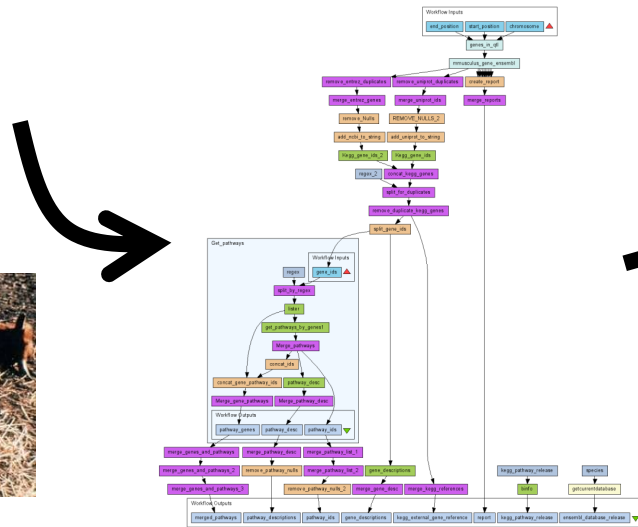
Reuse, Recycle, Repurpose Workflows



Dr Paul Fisher

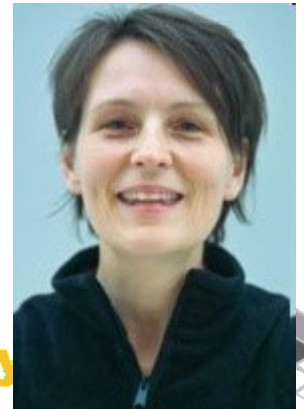


Identify QTg and pathways implicated in resistance to Trypanosomiasis in the mouse model



Dr Jo Pennock

Identify the QTg and pathways of colitis and helminth infections in the mouse model





Same Host, another Parasite...but the SAME Method

- Mouse whipworm infection - parasite model of the human parasite - *Trichuris trichuria*

Understanding Phenotype

- Comparing resistant vs susceptible strains – Microarrays

Understanding Genotype

- Mapping quantitative traits – Classical genetics QTL

Joanne Pennock, Richard Grencis
University of Manchester





Workflow Results

- Identified the biological pathways involved in sex dependence in the mouse model, previously believed to be involved in the ability of mice to expel the parasite.
- Manual experimentation: **Two year study** of candidate genes, processes unidentified
- Workflow experimentation: **Two weeks study** – identified candidate genes

Joanne Pennock, Richard Grecis
University of Manchester



“Traditional” Hypothesis-Driven Analyses

‘Cherry Pick’
genes

200 genes



Pick the genes involved in
immunological process

40 genes



Pick the genes that I am most
familiar with

2 genes



What about the other 198
genes? What do they do?

Biased view



Workflow Success

- Workflow analysed each piece of data *systematically*
 - Eliminated user bias and premature filtering of datasets
- The size of the QTL and amount of the microarray data made a manual approach impractical
- Workflows capture exactly where data came from and how it was analysed
- Workflow output produced a manageable amount of data for the biologists to interpret and verify
 - “make sense of this data” -> “does this make sense?”



Sharing and Reusing Workflows

myexperiment



[Home](#)[Users](#)[Groups](#)[Workflows](#)[Files](#)[Packs](#)[Services](#)[Topics](#)

Workflows ▾

[Home](#) > [Workflows](#)

Workflows

Search filter terms

Filter by type

- | | |
|--|------|
| <input type="checkbox"/> Taverna 2 | 1152 |
| <input type="checkbox"/> Taverna 1 | 645 |
| <input type="checkbox"/> RapidMiner | 223 |
| <input type="checkbox"/> Kepler | 43 |
| <input type="checkbox"/> Biodipse Scrip... | 34 |
| <input type="checkbox"/> LONI Pipeline | 26 |
| <input type="checkbox"/> GWorkflowDL | 24 |
| <input type="checkbox"/> BioExtract Server | 17 |
| <input type="checkbox"/> Tesla | 11 |
| <input type="checkbox"/> Galaxy | 10 |

Filter by tag

- | | |
|---|-----|
| <input type="checkbox"/> example | 230 |
| <input type="checkbox"/> mygrid | 104 |
| <input type="checkbox"/> localworker | 103 |
| <input type="checkbox"/> bioinformatics | 102 |
| <input type="checkbox"/> graph | 91 |

[« previous](#)[1](#)[2](#)[3](#)

...

[230](#)[next »](#)Sort by: [Rank](#) ▾

Showing 2293 results. Use the filters on the left and the search box below to refine the results.

Taverna 2

Pathways and Gene annotations for QTL region
(v7) [View](#) [Download \(v7\)](#)

Created: 19/11/09 @ 18:18:52 | Last updated: 07/09/12 @ 18:23:36

Credits: Paul Fisher

License: [Creative Commons Attribution-Share Alike 3.0 Unported License](#)Original
Uploader Paul
Fisher

This workflow searches for genes which reside in a QTL (Quantitative Trait Loci) region in the mouse, *Mus musculus*. The workflow requires an input of: a chromosome name or number; a QTL start base pair position; QTL end base pair position. Data is then extracted from BioMart to annotate each of the genes found in this region. The Entrez and UniProt identifiers are then sent to KEGG to obtain KEGG gene identifiers. The KEGG gene identifiers are then used to search for pathways in the KEGG path...

Rating: 4.6 / 5 (10 ratings) | Versions: 7 | Reviews: 1 | Comments: 7 |

New/Upload

Workflow ▾

 Katy
Wolstencroft

- [My Profile](#) [[edit](#)]
- [My Messages](#)
- [My Memberships \(4\)](#)
- [My History](#)
- [My News](#)

3 new friendship requests

- [mihaionita_me](#)
- [Pankaj chauhan](#)
- [Hanny](#)

4 new group requests

- [From Rolando Milian](#)
(for Group: Tutorial)
- [From Charlyb](#)
(for Group: Msc Tutorial)
- [From Mkh](#)
(for Group: MIB_Tutorial)
- [From Mateenraj](#)

Just Enough Sharing....

myExperiment can provide a central location for workflows from one community/group

- You specify:
 - Who can look at your workflow
 - Who can download and run your workflow
 - Who can modify your workflow
- Ownership and attribution

Share with my Groups:

<input type="checkbox"/> UsefulChem	View and Download only	▼
<input type="checkbox"/> Taverna 2 beta tester programme	View and Download only	▼
<input type="checkbox"/> Social Scientific Land	View and Download only	▼
<input type="checkbox"/> Mark's Project	View and Download only	▼
<input type="checkbox"/> GNU	View and Download only	▼
<input type="checkbox"/> myExperiment	View and Download only	▼
<input type="checkbox"/> Music workflows	View and Download only	▼



[Home](#) [Users](#) [Groups](#) [Workflows](#) [Files](#) [Packs](#) [Services](#) [Topics](#)[Home](#) > [Workflows](#)

Workflows

New/Upload

Workflow

GO

Search filter term

Filter by type

- ☐ Taverna 2
- ☐ Taverna 1
- ☐ RapidMiner
- ☐ Kepler
- ☐ Biodipse Scrip...
- ☐ LONI Pipeline
- ☐ GWorkflowDL
- ☐ BioExtract Server
- ☐ Tesla
- ☐ KNIME

Filter by tag

- ☐ example
- ☐ mygrid
- ☐ localworker
- ☐ ...

**BioVeL**
Biodiversity Virtual e-Laboratory
on myexperiment[Log in](#) | [Register](#) | [Give us Feedback](#) | [Invite](#)[Home](#) [Users](#) [Groups](#) [Workflows](#) [Files](#) [Packs](#) [Services](#) [Topics](#)[Home](#) > [Groups](#) > BioVeL

| Members (46) |

Group for sharing workflows
FP7-283359 BioVeLFor more information visit [http://...](#)

Created at: Saturday 06 Augus

Unique name: biovel

Search filter terms

Filter by type

☐ Taverna 2 59

Filter by tag

- ☐ panacea 19
- ☐ basicxces 17
- ☐ freeing 17
- ☐ example 15
- ☐ taqqing 15
- ☐ crawled 11
- ☐ bilingual 7
- ☐ cqp 5

[Home](#) [Users](#) [Groups](#) [Workflows](#) [Files](#) [Packs](#) [Services](#) [Topics](#)[Home](#) > [Workflows](#)

Workflows

[« previous](#) [1](#) [2](#) [3](#) ... [6](#) [next »](#)

Sort by: Rank

Showing 59 results. Use the filters on the left and the search box below to refine the results.

Search

Taverna 2

bilingual word aligner for crawled data (v3)

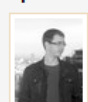
View

Created: 19/04/11 @ 15:18:06 | Last updated: 03/06/11 @ 10:18:35

Credits: Marcpoch

Download (v3)

License: Creative Commons Attribution-Share Alike 3.0 Unported License

Original
Uploader

Marcpoch



This is a word alignment workflow using hunalign and giza++.

New/Upload

Workflow

GO

Log in / Register

Username or Email:

Password:

Remember me: ☐

Log in

Need an account?
[Click here to register](#)[Forgot Password?](#)Popular Tags
25 tags

Workflow Entry: Microarray CEL file to candidate pathways

All versions of this workflow are licensed under a [Creative Commons Attribution-NoDerivs 3.0 License](#).

Version: 2 (latest)

Change to: **2 (latest)**

Title: Microarray CEL file to candidate pathways

Version created on: Wednesday 03 October 2007 @ 18:35:55 (GMT Daylight Time)

Diagram: (click to expand)



Description:

This workflow takes in a <http://www.bioinf.manchester.ac.uk/MADAT/index.html>. Also retruned by this workflow are a list of the top differentially expressed genes (size dependant on the number specified as input – geneNumber), which are then used to find the candidate pathways which may be influencing the observed changes in the microarray data. By identifying the candidate pathways, more detailed insights into the gene expression data can be obtained.

CEL

Uploader



Paul Fisher

Credits (People)

None

Attribution (Workflows)

None

Tags

All Tags

Uploader's Tags

insertional mutation | shim | shotgun method | similarity | simplifier

Note: the size of the tags show how popular they are on myExperiment.

Add Tags

Ratings

Hover and click to rate



Current: 4.50/5 (2 ratings)

You rated it:

Breakdown

Paul Fisher - 5/5
Jits - 4/5

Statistics

New/Upload

Workflow



Jits

My Stuff

11 friends | 1 group | 2 files | 4 workflows

cmuzys1_2 readme.txt
rssbanditlist.opml

Workflows

Example of an alter...
ExcelFileReader and...
Show Gene Ontology ...
Fetch today's Dilbe...

My Tags

20 tags

a | add | another | are | boy! | fg | i | j | is a | its | jh | jh fg | kool | oh boy | really? | tag | test | testing mania | this | we

Popular Tags

25 tags

another | bioinformatics | i | j | insertional mutation | nucleic acid | pluripotency | scaffold | SEG | sequencing | shim | similarity | simplifier |

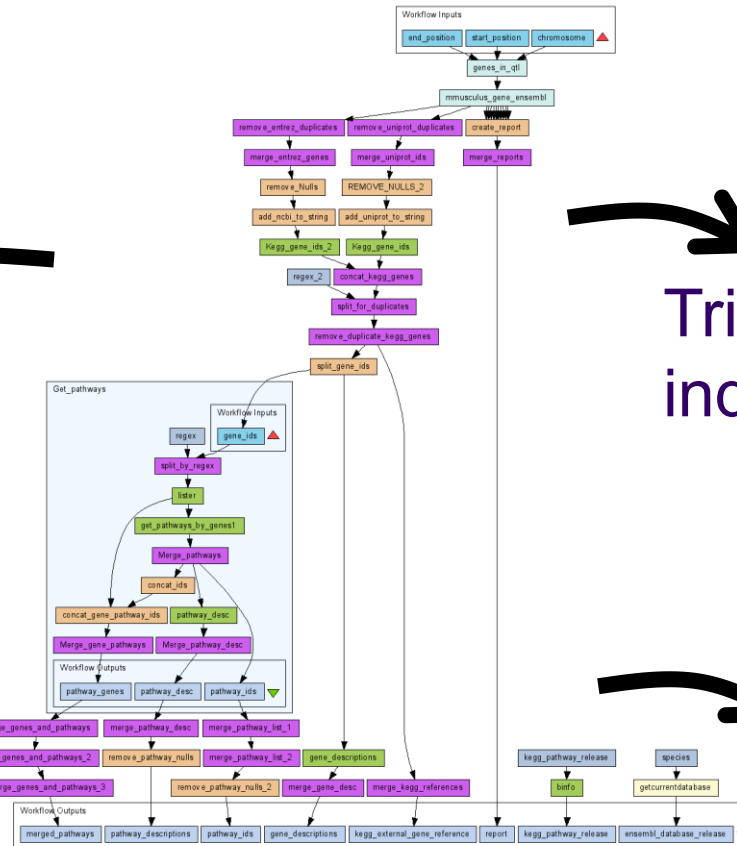
Reuse, Reuse, Reuse

Atopic Dermatitis

Trichuriasis induced Colitis

Blood Pressure

Epilepsy



FINDING AND USING A MYEXPERIMENT WORKFLOW: DEMO



Workflow engine features

- Implicit iterations
 - With customisable list handling
- Parallelisation
 - Run as soon as data is available
- Streaming
 - Process partial iteration results early
- Retries, failover, looping
 - For stability and conditional testing



Data and Provenance

- Workflows can generate vast amount of data - how can we manage and track it?
- We need to manage data **AND** metadata **AND** experimental provenance
- Scientists need to check back over past results, compare workflow runs and share workflow runs with colleagues
- Scientists need to look at intermediate results when designing and debugging



Data and Provenance Handling

- **Provenance** captured for workflow runs
 - **Trace** execution steps, view **intermediate values** while running
 - Export as Open Provenance Model (OPM) / RDF
 - Proof and **origin** of produced outputs
 - Extensible **annotations**
- Wf4Ever: reproducible **research objects**
 - Workflow/data as a scientific publication → preservation
 - Need to capture more service data and metadata



Advanced users design and build workflows (informaticians)

Spectrum of Users

myexperiment beta

Users Groups Workflows Files Blogs Forums

Home » Workflows » View: Microarray CEL file to candidate pathways

Workflow Entry: Microarray CEL file to candidate pathways

All versions of this workflow are licensed under a [Creative Commons Attribution-NoDerivs 3.0 License](#).

Version: 2 (latest) Change to: 2 (latest) GO

Title: Microarray CEL file to candidate pathways
Version created on: Wednesday 03 October 2007 @ 18:35:55 (GMT Daylight Time)
Diagram: (click to expand)

Description:
This workflow takes in a <http://www.bioinf.manchester.ac.uk/MADAT/index.html>. Also returned by this workflow are a list of the top differentially expressed genes (size dependant on the number specified as input – geneNumber), which are then used to find the candidate pathways which may be influencing the observed changes in the microarray data. By identifying the candidate pathways, more detailed insights into the gene expression data can be obtained.

Uploader: Paul Fisher

Credits: (People)
None

Attribution: (Workflows)
None

Tags:
All Tags Uploader's Tags
insertional mutation | shim | shotgun method | similarity | simplifier
Note: the size of the tags show how popular they are on myExperiment.

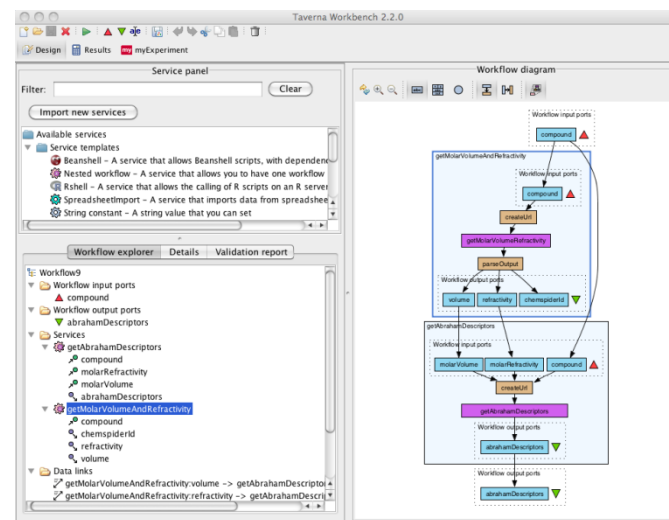
Ratings:
Current: 4.50/5 (2 ratings)
You rated it:

Popular Tags:
25 tags
bioinformatics | insertional mutation | nucleic acid | pluripotency | scaffold | sequencing | shim | similarity | simplifier

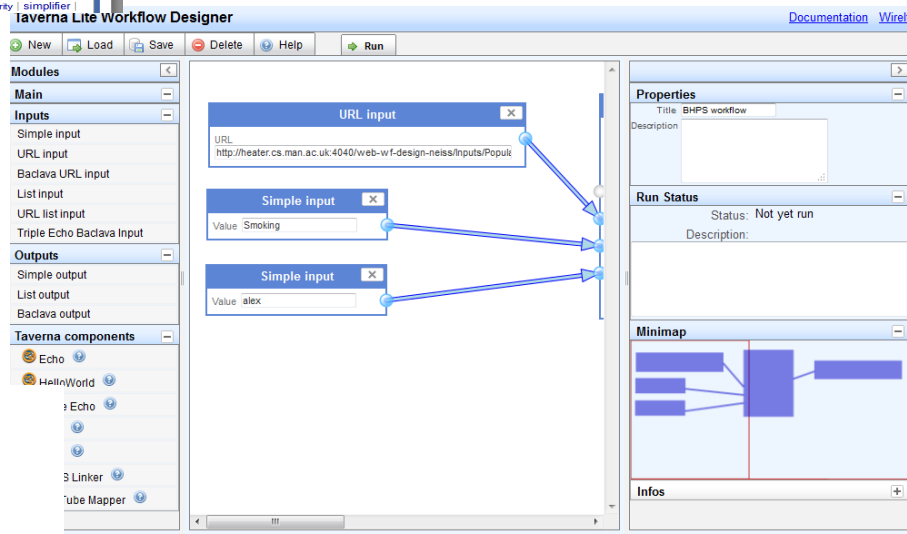
<http://www.myexperiment.org>

Load Data:

Run Workflow



Intermediate users reuse and modify existing workflows



Others “replay” workflows through a web interface or Taverna Lite

TAVERNA SERVER



Taverna Server

- Running workflows remotely
 - Through other client software
 - Via a web interface
- Tapping into remote computing resources
 - Execution on servers, grids or clouds



Limitations of the Desktop workbench

- You have to install it and learn how to use it
- Although computation could happen at remote service locations, data and computation can also happen locally
- High throughput experiments take a lot of compute and a lot of time
- Long running workflows need uninterrupted execution

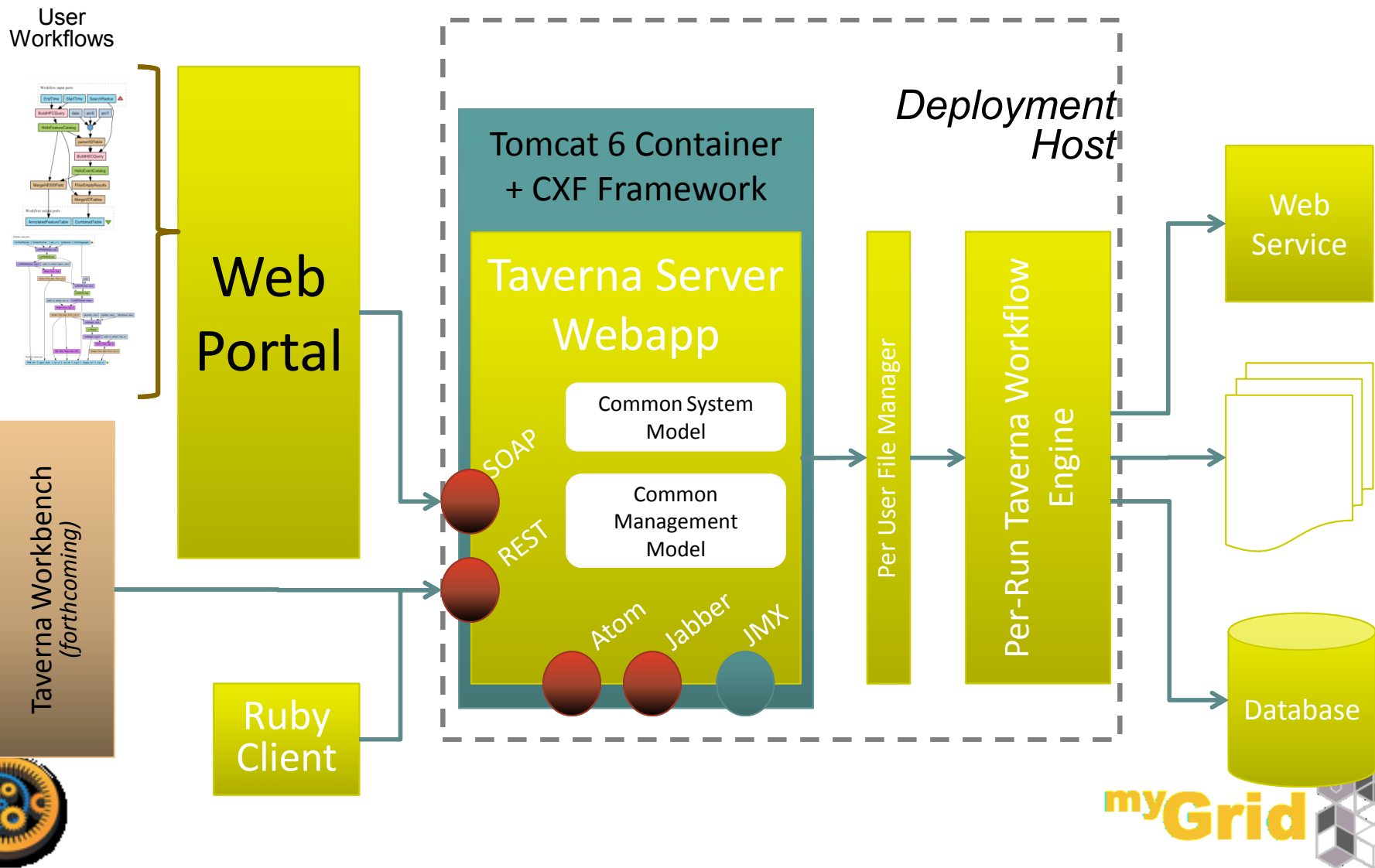


Data Limitations with the Desktop Workbench

- Running the Workbench is limited by:
 - Local disk space for storing data
 - Network speeds for up/download
 - Firewall access



Taverna Server



Taverna Server in Use

- T2Web, running myExperiment workflows through web interface
- HELIO - Heliophysics Integrated Observatory
- SCAPE - SCalable Preservation Environment (digital archives)
- BioVel – Biodiversity Virtual e-laboratory
- Cloud analytics for the life sciences – Taverna on the cloud
- Running Taverna through Galaxy



Workflow: BioAID_ProteinDiscovery - Mozilla Firefox

Workflow: BioAID_ProteinDis...

www.mybiobank.org/t2web/workflow/74

Getting Started Getting Started BioBank Bookmarklett Latest Headlines Remember The Mil...

nbtic netherlands bioinformatics centre

BioAID_ProteinDiscovery

workflow by [Marco Roos](#) Leiden University Medical Centre

Configure Workflow Inputs

Description: Fill in a search query, similar to pubmed. For advanced queries look up the Lucene syntax (http://lucene.apache.org/java/2_9_1/queryparsersyntax.html).

Enter Query:
"transmembrane"

Upload file? ☐

Enter maxHits parameter:
5

Upload file? ☐

Execute

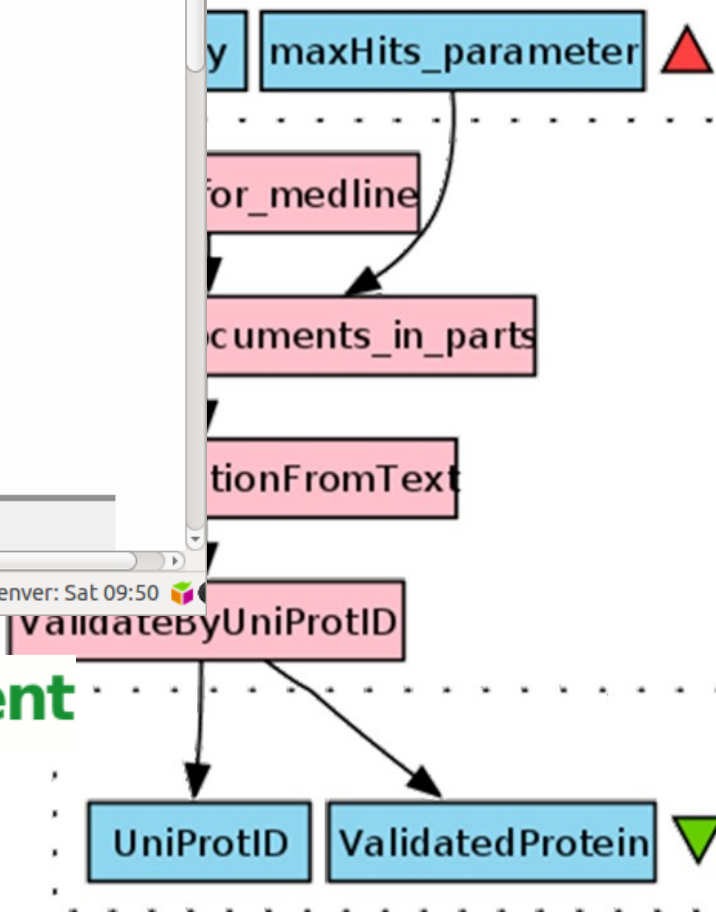
Workflow Description

zotero UK: Sat 16:50 US Pacific: Sat 08:50 Hong Kong: Sat 23:50 GMT/UTC: Sat 15:50 Denver: Sat 09:50

T2 Web

Marco Roos

Kostas Karasavvas



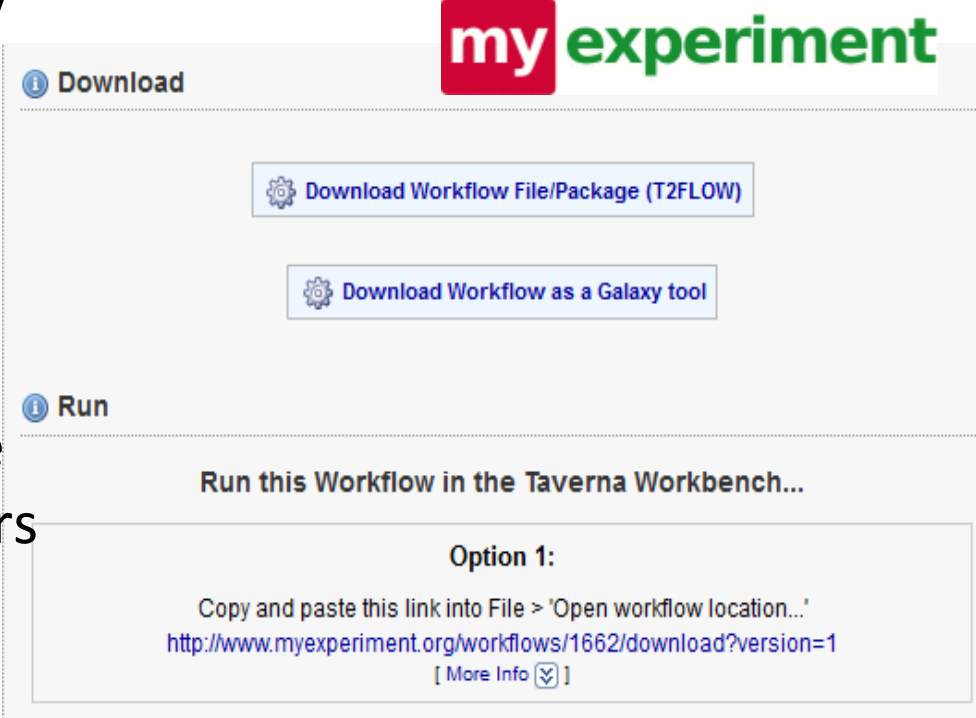
my experiment

myExperiment workflow ID



Running Taverna Through Galaxy

- Workflow interoperability
 - The methods are more important than the platform
 - Workflows in Galaxy and Taverna already exist
 - Any Taverna workflow can be made available to Galaxy users
 - Discover and import from myExperiment

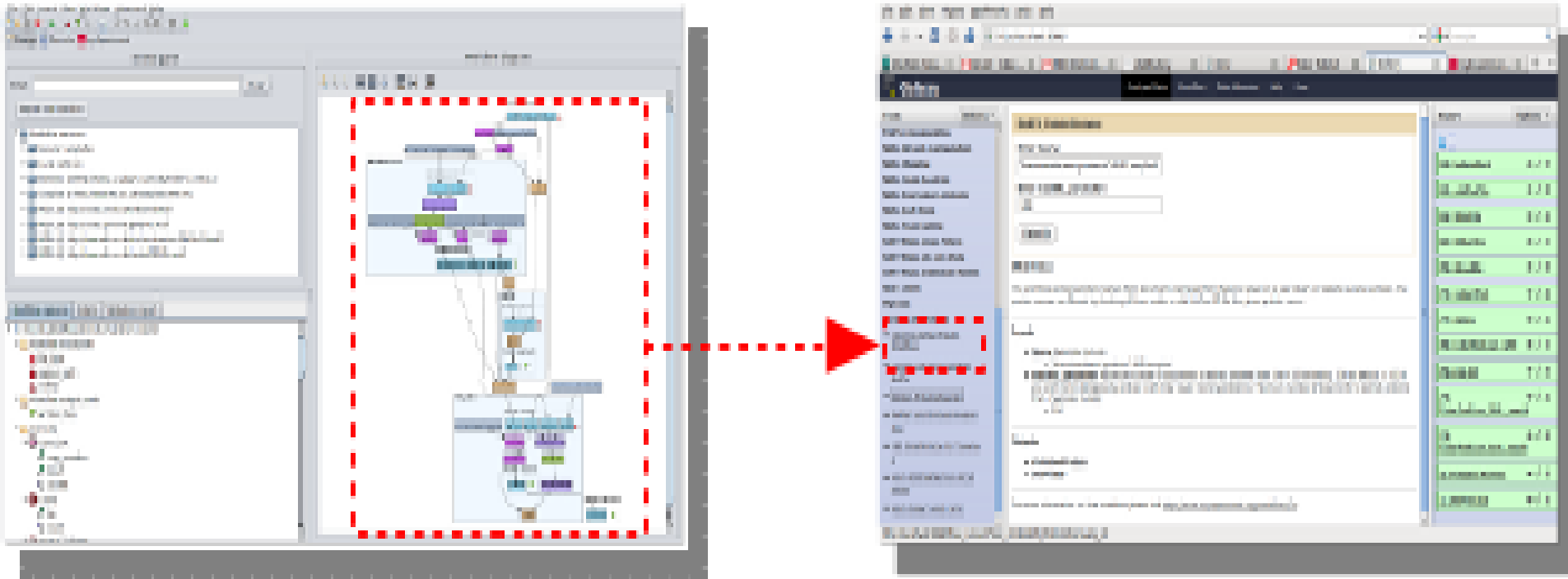


The screenshot shows the 'myexperiment' interface. Under the 'Download' section, there are two buttons: 'Download Workflow File/Package (T2FLOW)' and 'Download Workflow as a Galaxy tool'. Under the 'Run' section, it says 'Run this Workflow in the Taverna Workbench...' and provides 'Option 1: Copy and paste this link into File > 'Open workflow location...'' with the URL <http://www.myexperiment.org/workflows/1662/download?version=1> and a '[More Info]' link.



Running Taverna through Galaxy

Kostas Karasavvas, NBIC



- Connect the Taverna and Galaxy communities
- Galaxy specialises in genomics, next gen sequencing etc
- Taverna can access more 'downstream' analysis services – e.g. pathway analyses, literature, GO enrichment etc



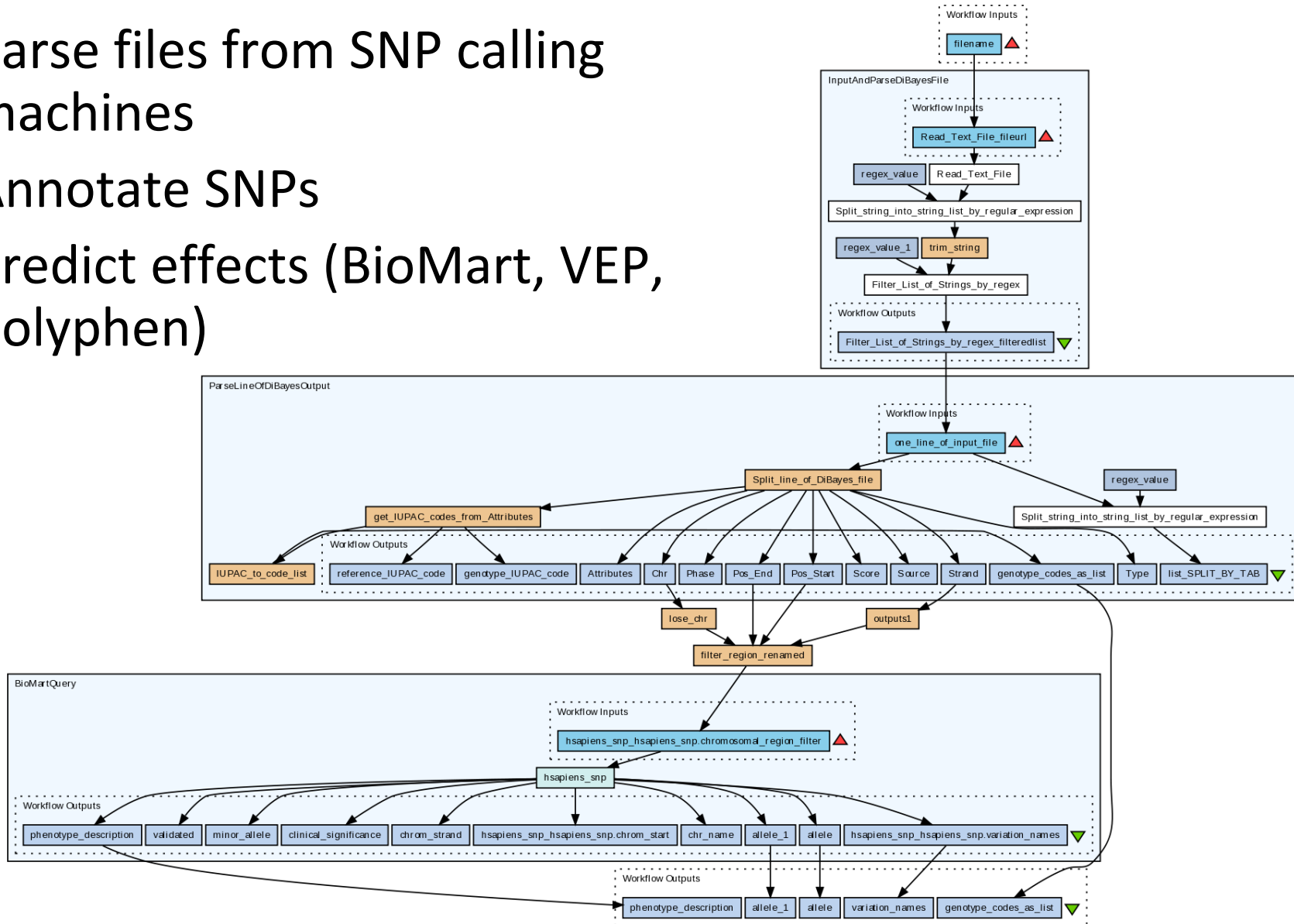
Cloud Analytics for the Life Sciences

- Workflows for genetic diagnostics (for the NHS)
 - Exome and whole genome
 - SNP analysis and annotation
- Execution on the cloud
 - Secure execution and results handling
 - Elastic to cope with demand
 - Pay-as-you-go – cheap at the point of use



A Typical Workflow

- Parse files from SNP calling machines
- Annotate SNPs
- Predict effects (BioMart, VEP, polyphen)



A Typical Workflow

ElasticView

(v 0.1 Alpha) Powered by Eagle Genomics Ltd.

madhuLogout »

My Workspace

Data+↻?

Workflows-↻?

Hello, World!📁

Simple pass thr...📁

Simple SNP work...📁

madhuruns🗑️

My Activity

Getting Startedmadhuruns

Inputs+↻?

Progress-↻?

Overall Run Progress

InputAndParseDiBayesFile1/1 done

ParseLineOfDiBayesOutput1/1 done

BioMartQuery1/1 done

Results+↻?

Advantages

- Workflows are reusable
- Cloud computing infrastructure manages large data and processes – no need for big local resources
- Genomic analyses easy to run in parallel
- Simple submission through web interface for researchers
- Selecting ready-made workflows
- Simple and limited configuration of workflows
- Collaboration with industry – commercialisation of the services

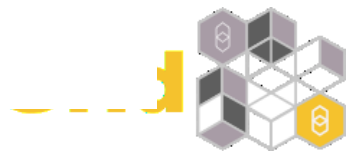
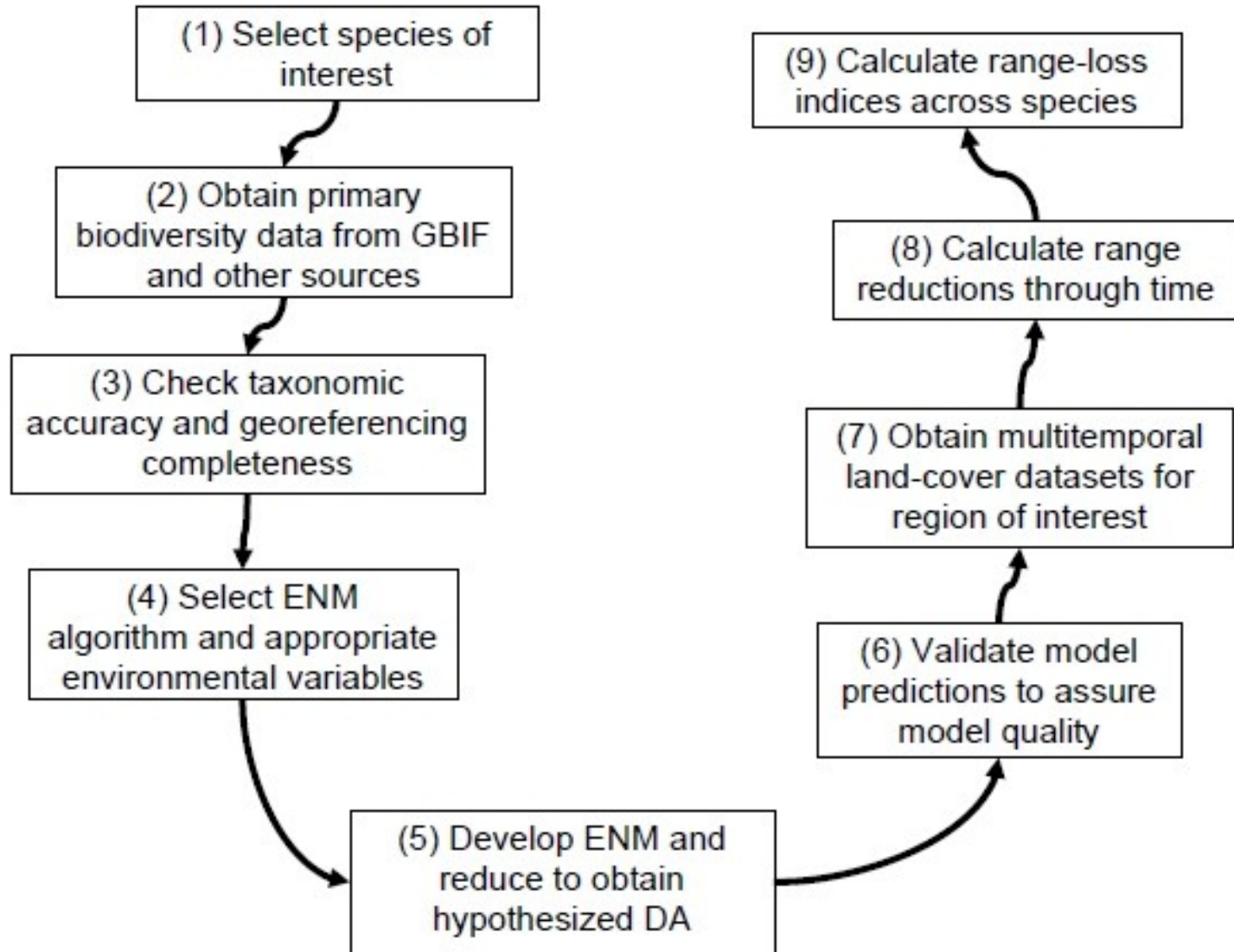


Biodiversity Virtual e-Laboratory

- A network of expert scientists who develop, support, and use workflows and services in biodiversity
- Workflows, including:
 - Phylogenetics
 - Metagenomics
 - Ecological niche modelling
 - Species distribution modelling
 - Models how environmental niches of a species shift due to the changing climate.



Case Study: Ecological Niche Modelling



Interaction Service: Communicating with your Remote Workflow

- Service suspends workflow execution to wait for further input from the user
- Interaction through the web interface
- Messages between workflow engine and web page via ATOM feeds, using Javascript



TAVERNA SERVER DEMO



A RECAP ON TAVERNA WORKFLOWS



Taverna Advantages

- Allows complex analysis pipelines
- Access to local and remote services (>8000 in biology)
- New services 'added' instantly
- Workflows can be shared and run in any Taverna instance
- Can be used for any areas of bio or non-bio research



Issues and Problems

- Transferring large data over networks
 - Take services to data (like in the cloud example)
 - Pass by reference, rather than by value
 - Transfer only what you need for analysis
- Service incompatibility
 - shims – sharing and reusing
 - Creating integrated sets of services → components
- Services changing and vanishing
 - Use BioCatalogue and myExperiment to identify alternatives and find similar methods

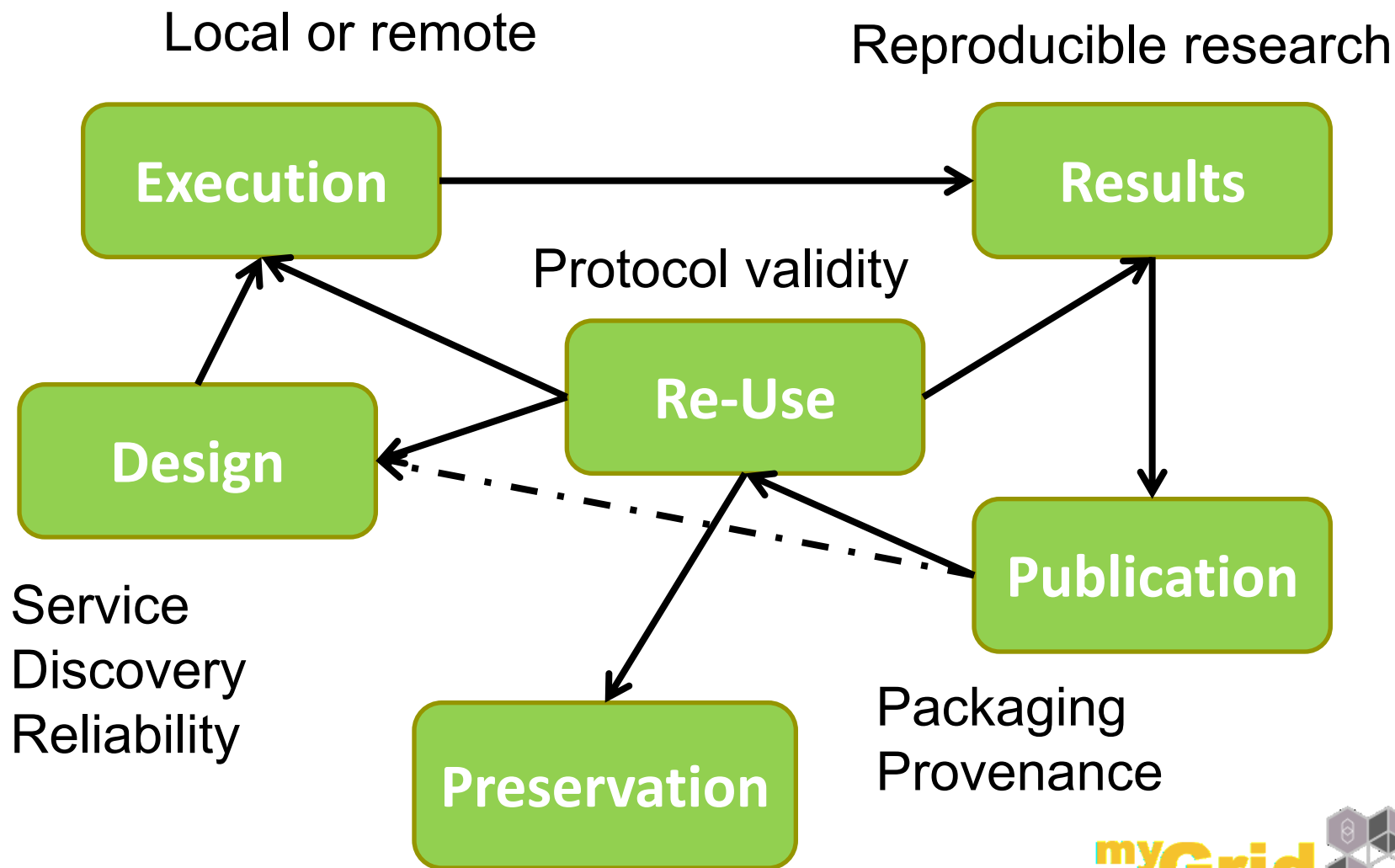


Components

- A set of services designed to be compatible by
 - Consistent annotation to help understand how they work
 - Combining with shims to provide uniform (or predictable) input and output formats
- Hiding the complexity of public web services



Taverna Workflows Supporting *in silico* Science



Taverna 3 roadmap

- OSGi plugin system
- Workflow language: ScufI2
 - Making programmatic interaction easier
 - Compound format; embedding metadata, dependencies, independent API for creating/inspecting workflows
- Components
 - Finding/sharing command line tool descriptions
 - Richer way of finding compatible services



Summary – Workflow Advantages

- Informatics often relies on data integration and large-scale data analysis
- Workflows are a mechanism for linking together resources and analyses
- Automation
- Large data manipulation
- Promote reproducible research
- myExperiment allows you to reuse workflows and benefit from others work
- Easy to find and use successful analysis methods



More Information

- Taverna
 - <http://www.taverna.org.uk>
- myExperiment
 - <http://www.myexperiment.org>
- BioCatalogue
 - <http://www.biocatalogue.org>



Acknowledgements

- myGrid consortium, in particular
 - Paul Fisher
 - Carole Goble
 - Alan Williams
 - Stian Soiland
 - Khalid Belhajjame
 - Rob Haines
 - Donal Fellows
 - Helen Hulme
- Trypanosomiasis project
 - Andy Brass
 - Paul Fisher
 - Harry Noyes



