

Introduction to Workflows with Taverna and myExperiment

Stian Soiland-Reyes and Christian Brenninkmeijer
University of Manchester

materials by Dr Katy Wolstencroft and Dr Aleksandra Pawlik

<http://orcid.org/0000-0001-9842-9718>

<http://orcid.org/0000-0002-2937-7819>

<http://orcid.org/0000-0002-1279-5133>

<http://orcid.org/0000-0001-8418-6735>



*This work is licensed under a
[Creative Commons Attribution 3.0 Unported License](http://creativecommons.org/licenses/by/3.0/)*

Bonn University, 2014-09-01

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



- ‘Omics data
- Next Gen Sequencing
- eGovernment
- World bank data
- Climate change data
 - Large Hadron collider
 - Astronomy

Data Deluge

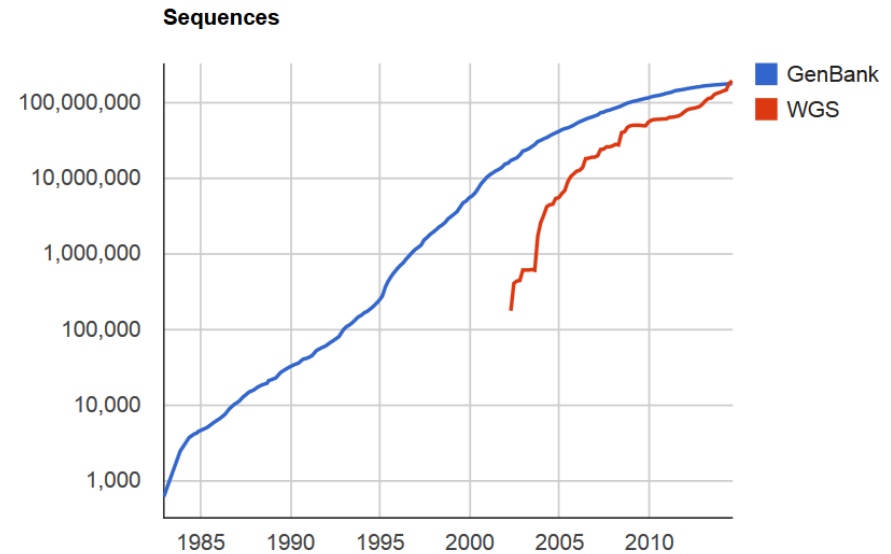


Lots of Resources

NAR 2014: 1552 databases

Genbank 2014-04: 172 million sequences,
162 billion basepairs

WGS 2014-04: 774 billion basepairs



Next Generation Sequencing

- 2008-2012: 1000 Genome Project
 - A Deep Catalog of Human Genetic Variation
- 2009-: Genome 10k project
 - A genomic zoo—DNA sequences of 10,000 vertebrate species, approximately one for every vertebrate genus.
- 2012-: Human Microbiome Project
 - Characterise the microbial communities found at several different sites on the human body



Where is the data?

- Repositories run by major service providers (e.g. NCBI, EBI)
- Local project stores
- Static web pages
- Dynamic web applications
- FTP servers (!)
- Inside PDFs ☹️
- Web Services 😊



National Center for Biotechnology Information (USA)



Tokyo, Japan



Cambridge, UK



SRS



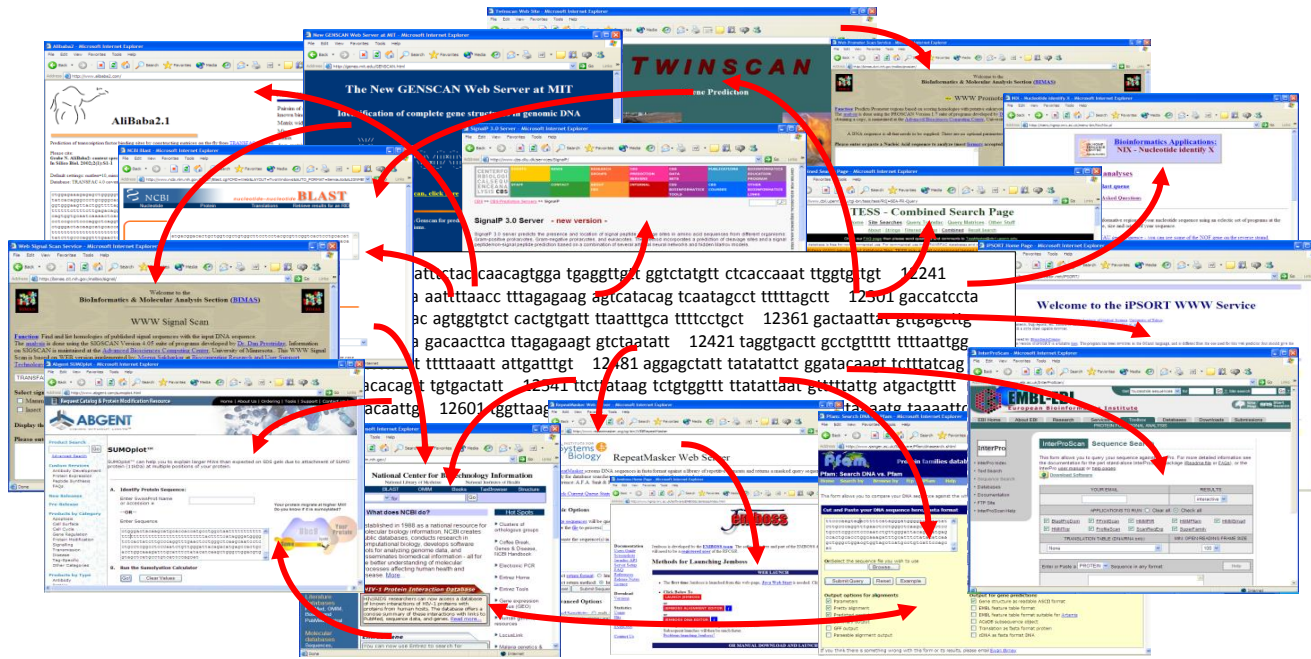
SeqHound



The implicit workflow

Bioinformatics research combines:

- Data resources (public and private)
- Computational power (standard and custom)
- Researchers and collaborators



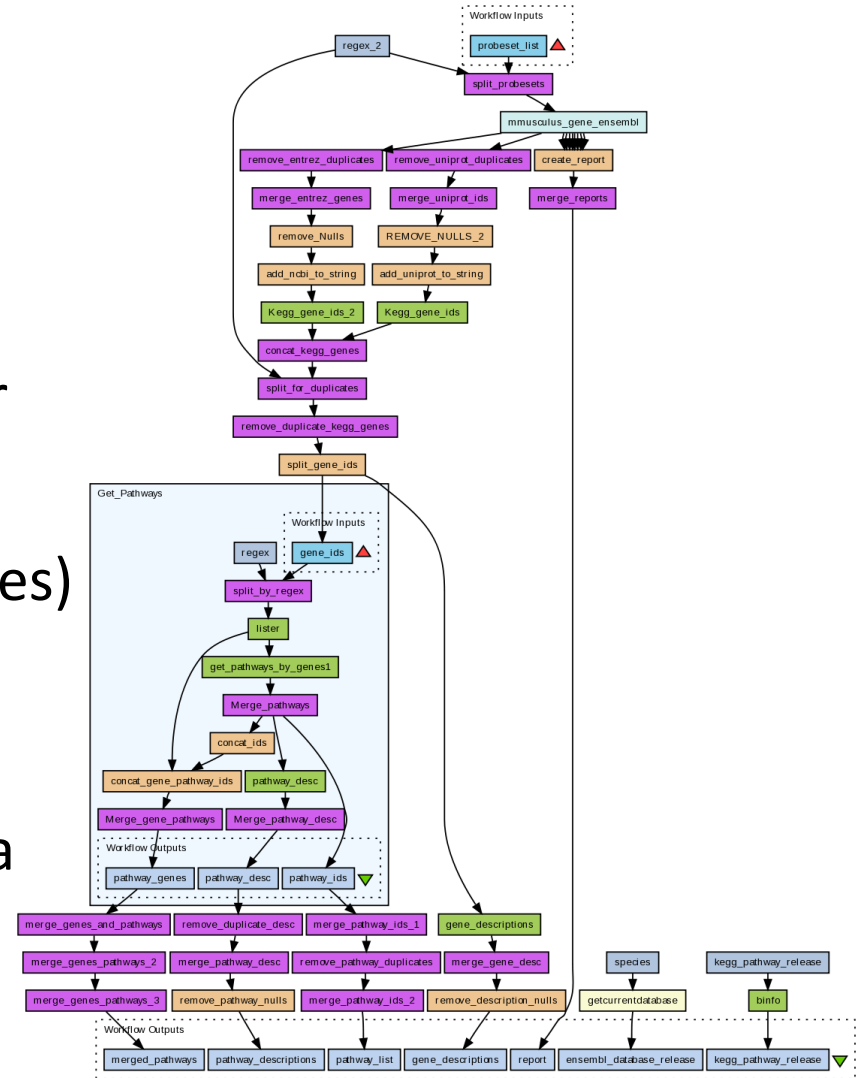
What that means for Bioinformatics

- Sequential use of distributed tools
- Incompatible input and output formats
- Challenging to record/reproduce/tweak
 - parameter selections
 - service selection
 - results of each step
- OK for one gene or one protein, but what about 10,000?
 - Analysing large data sets requires programmatic help

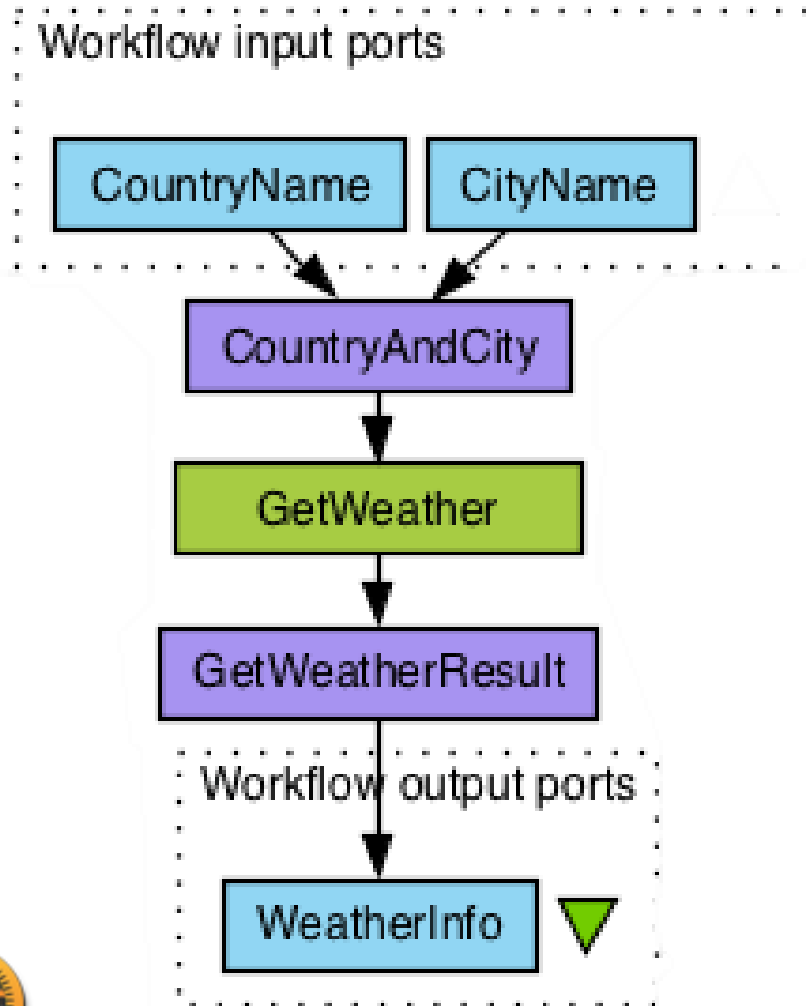


Workflow as a Solution

- Sophisticated analysis **pipeline**
- **Graphical** representation of executable analysis
- Combine a set of **services** to analyse or manage data (local or remote)
- Data **flow** from one service (boxes) to the next (connected with arrows)
- **Iteration** – process multiple data items
- Automation – **rerun** workflow



Example Taverna Workflow



Workflow: Get the weather forecast of the day given the city and the country

Green box is a **Web Service**

Purple boxes are local **XML** services to assemble/ extract XML

Blue boxes are workflow **input** and **output** ports

Arrows define the direction of **data flow**



Workflows as a solution

- **Flow of data** from one tool to the next is automatic – just connect inputs and outputs
- Incompatibilities overcome in the workflow with helper services (*shims*)
 - Allowing new tool combinations
- Workflow engine records parameter values and algorithms – **provenance**
- Workflows can include data **integration** and **visualization**
- **Iteration** over large data sets automatic – ideal for high throughput analysis (e.g. omics)



Reproducible Research

Preventing non-reproducible 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 – reported in *Nature Medicine*
- 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 the *in-silico* experiments and see results for themselves
- Methods could be properly examined and criticized by inspecting the workflow
- Mistakes and opportunities could be pinpointed earlier



Taverna Workbench

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

Freely available,
open source

80,000+ downloads
across versions

Installers for Windows,
Mac OS X, Linux

Current version: 2.5.0



The screenshot shows the Taverna website homepage. At the top left is the Taverna logo, a gear with a blue and yellow color scheme. To its right is the word "Taverna" in a large, bold font. Further right is the "myGrid" logo and a search bar. Below the logo and name is a navigation menu with links for "Introduction", "Documentation", "Download", "Developers", "News", "Publications", and "About". The main content area features a large banner for "Taverna Workflow Management System" with a description: "Powerful, scalable, open source & domain independent tools for designing and executing workflows. Access to 3500+ resources." To the right of the banner is a "RECENT NEWS" section with three entries: "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 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). Below these buttons is an "IN PRESS" section with four items: "Taverna 2.3", "Taverna 3 Next Generation", "SCUFL2 workflow bundle language", and "Taverna infrastructure VMs". The main text area below the banner describes Taverna as an open source and domain independent Workflow Management System, used for designing and executing scientific workflows and *in silico* experimentation. It mentions that Taverna was created by the myGrid team and funded through the OMII-UK, with guaranteed funding until 2014. It also states that the Taverna suite is written in Java and includes the Taverna Engine, Taverna Workbench, Taverna Server, and a 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. The video player has a YouTube logo and a caption: "As the workflow runs, you can check its progress and...".

Wolstencroft et al. (2013): **The Taverna workflow suite: designing and executing workflows of Web Services on the desktop, web or in the cloud**, *Nucleic Acids Research*, **41**(W1): W557-W561. doi:[10.1093/nar/gkt328](https://doi.org/10.1093/nar/gkt328)



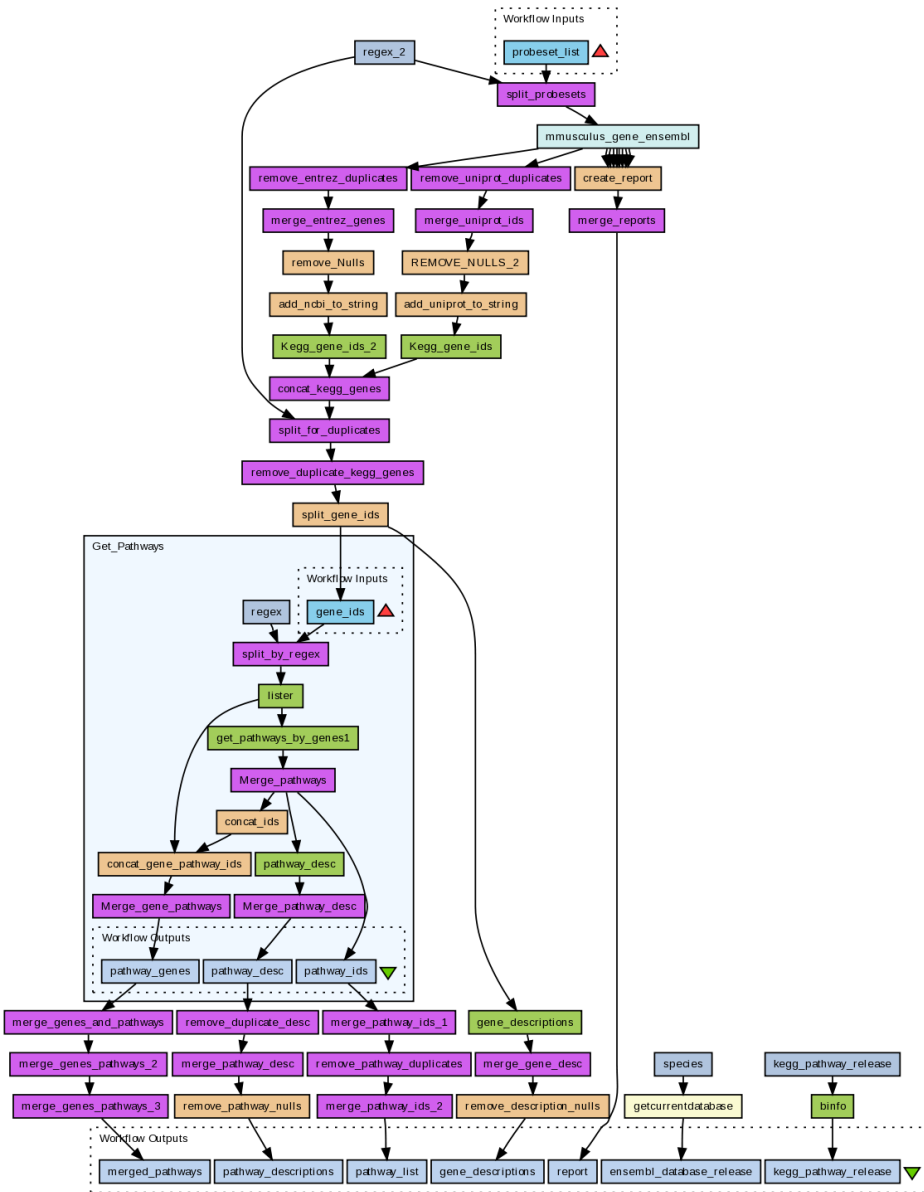
Taverna Workflow System

History:

- 2003: Taverna 0.1 (300 downloads)
- 2014: Taverna 2.5.0 (5100 downloads)

Products:

- Taverna Workbench
- Taverna Server
- Taverna Command line
- Taverna Online
- Taverna Player
- Plugins and integrations



Taverna editions and extensibility

Taverna is a generic workflow system that can be extended by **plugins** and customized for use in different domains.

The Taverna **editions** are pre-built downloads of Taverna with plugins for the most popular domains.

- Core
- Astronomy
- **Bioinformatics**
- Biodiversity
- Digital Preservation
- Enterprise



<http://www.taverna.org.uk/download/workbench/2-5/>



Taverna Workbench

Workflow engine to run workflows

List of services

Construct and visualise workflows

Web Services

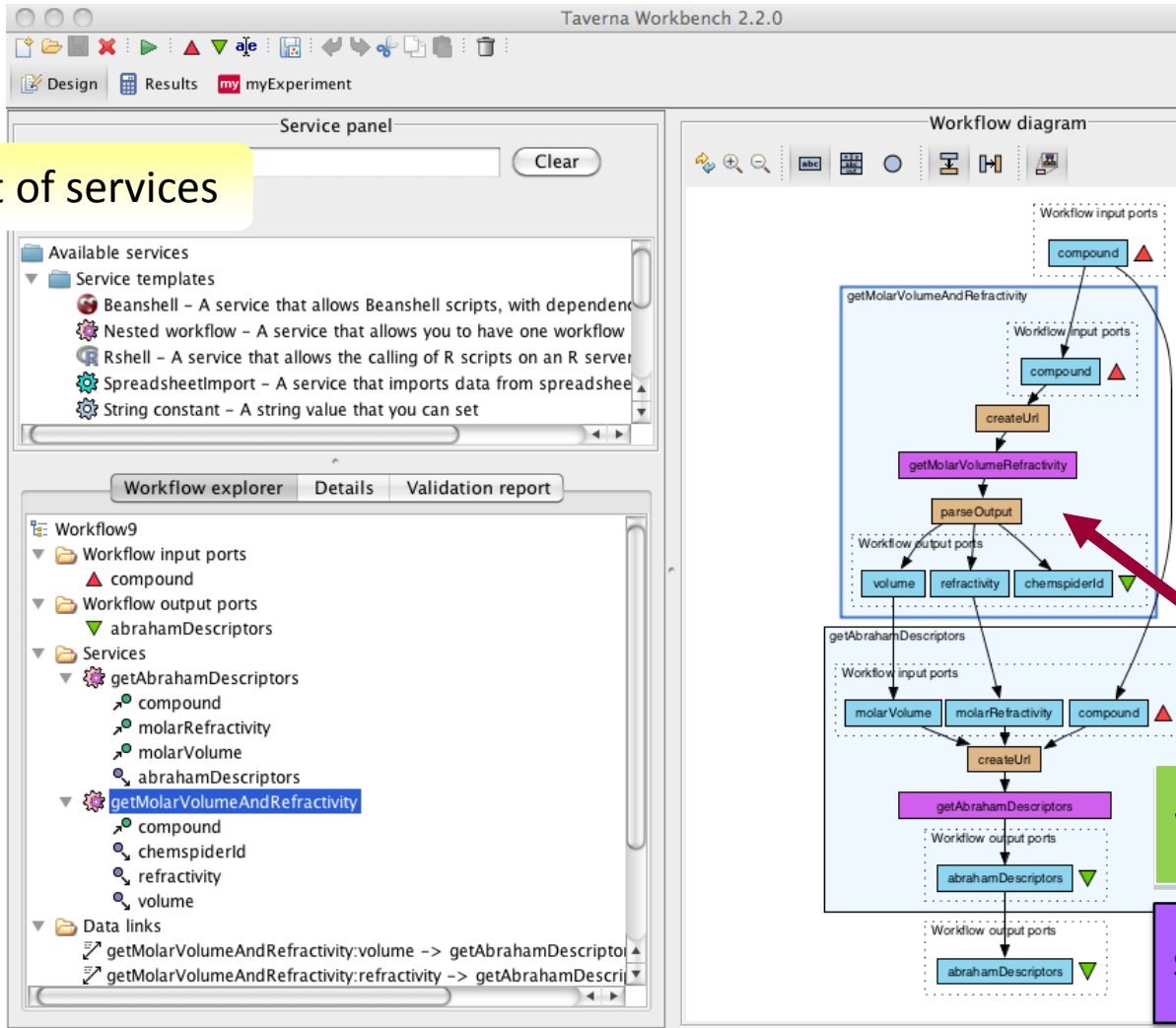
e.g. KEGG

Scripts

e.g. beanshell, R

Programming libraries

e.g. libSBML



Using Tools and Services from Taverna workflows

- Web Services
 - WSDL
 - REST
- Data services
 - BioMart
- Local scripts:
 - R
 - Beanshell
 - Command line (e.g. Python, Perl)
- Other workflows
- And more..... Add your own!



What are Web Services?

Web Services: HTTP-based programmatic access (API).

Instead of “*GET me the web page*
http://example.com/cat-pics”,

Web Services allow “*GET me a genome sequence*
http://example.com/gene/WAP_RAT”

Connect to and use remote services from your
computer in an automated way

NOT the same as services on the web (i.e. forms that
shows results as a web page)



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



Cambridge, UK



SRS

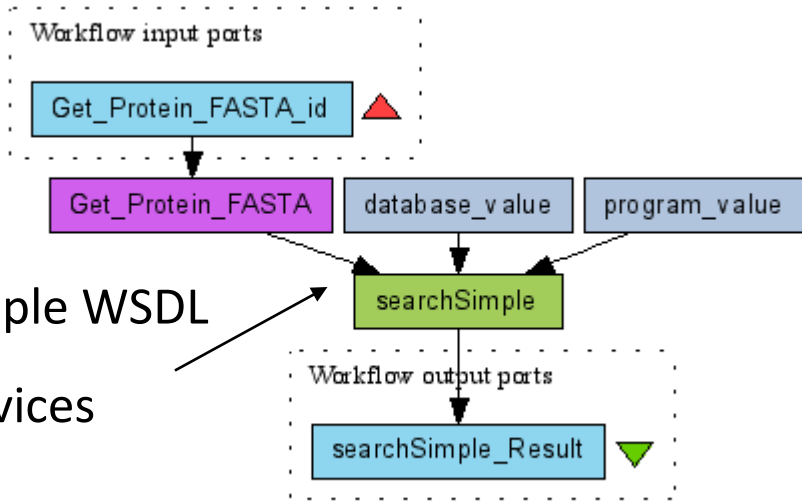


SeqHound

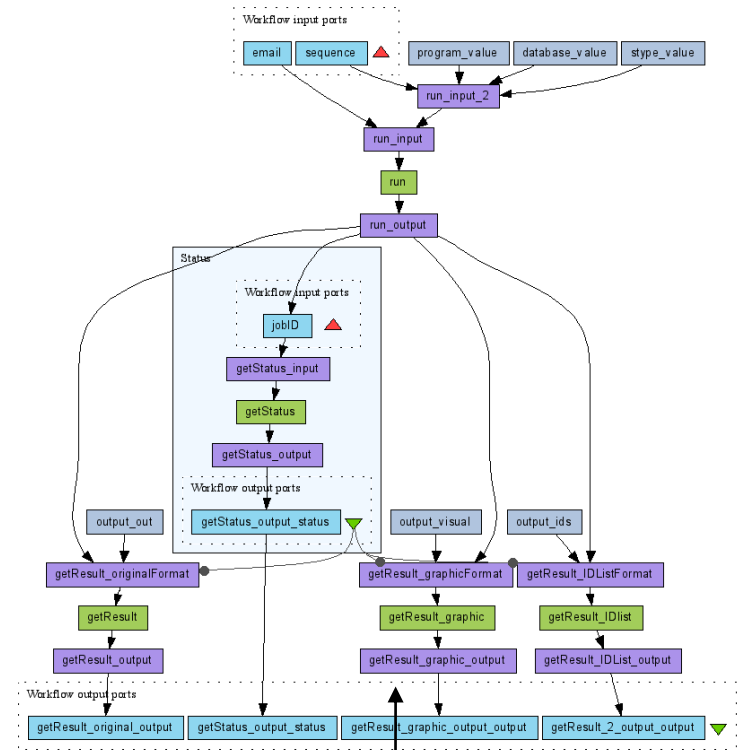
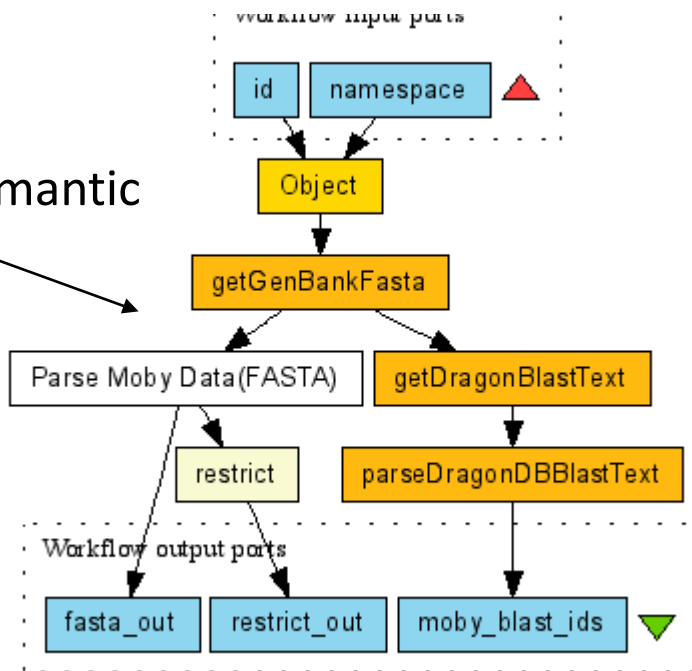


How do you use the services?

Simple WSDL
services



BioMoby Semantic
Services



Asynchronous services
(Submit, Wait, Fetch)



Search:



[Home](#)

[Services](#)

[Register a Service](#)

[Providers](#)



[Home](#) »










SHARE

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

-  [yasunun](#) **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 \(REST\)](#)
- [F&STM \(REST\)](#)



Home » Services » InterProScan (REST)

SHARE [social icons]

InterProScan (REST) REST

99 0

aka JDispatcher aka InterProScan

Monitoring

Categories: Function Prediction

Annotations: 30 0 30 0

Overview REST Endpoints (6) Monitoring News

Provider: [European Bioinformatics Institute \(EBI\)](#)

Provider

Location: UNITED KINGDOM

Submitter

Submitter / Source: [Hamish McWilliam](#) (about 1 month ago)

Service Description

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

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)

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**

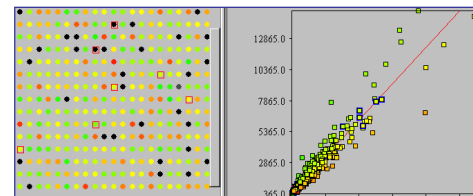
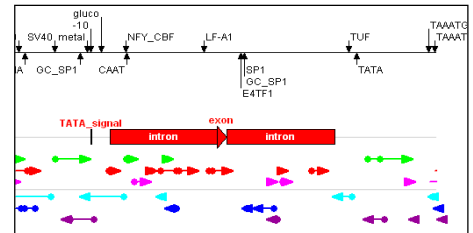
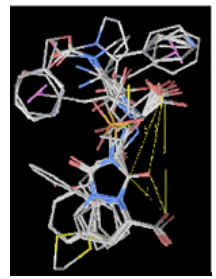
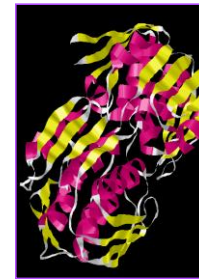
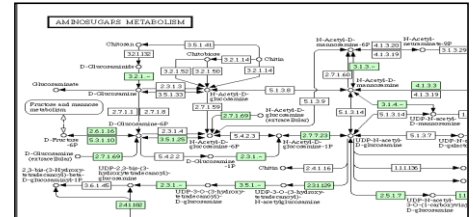
High throughput screening Proteomics **Music**

Phenotypical studies Text mining **Meteorology**

Public Health care epidemiology **Social Science**

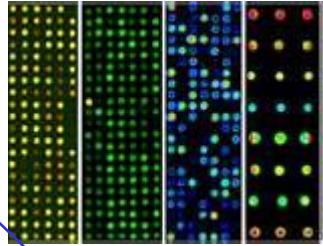
Medical image analysis QTL studies **Cheminformatics**

QSAR studies Genome Wide Association Studies



Lymphoma Prediction Workflow

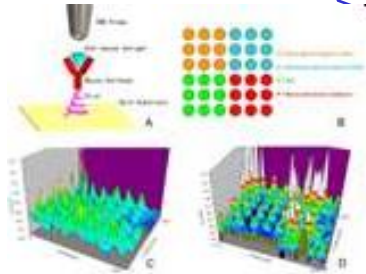
MicroArray from
tumor tissue



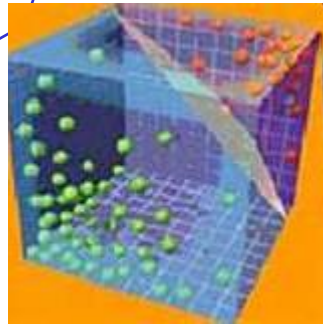
caArray

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

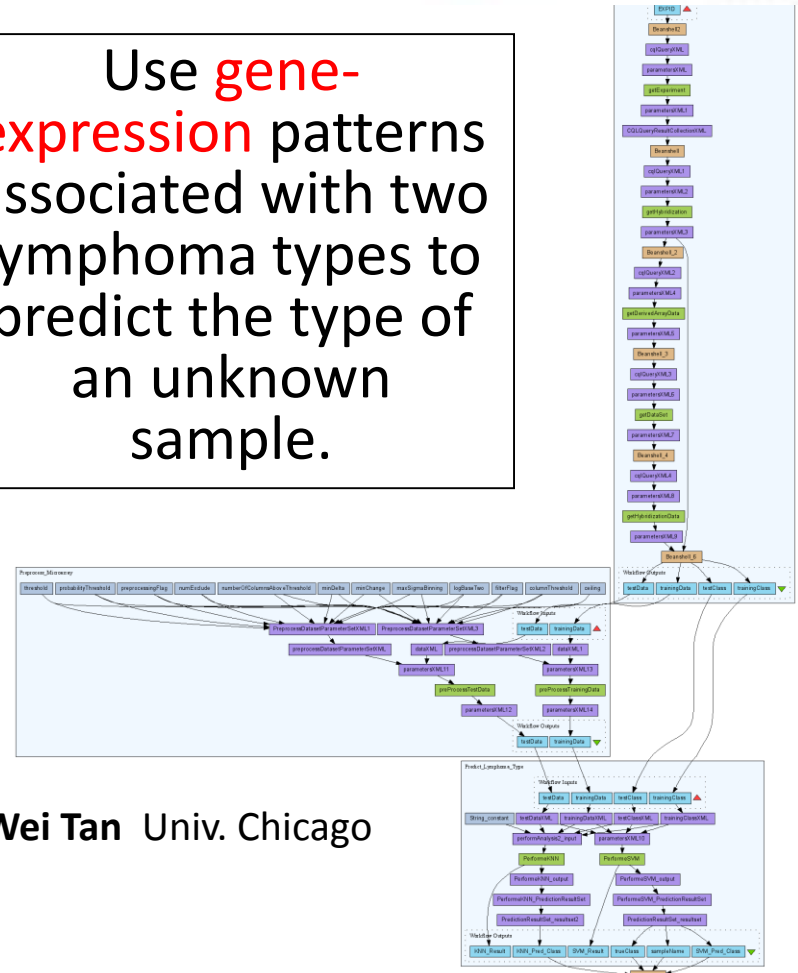
Microarray
preprocessing



Lymphoma
prediction



GenePattern

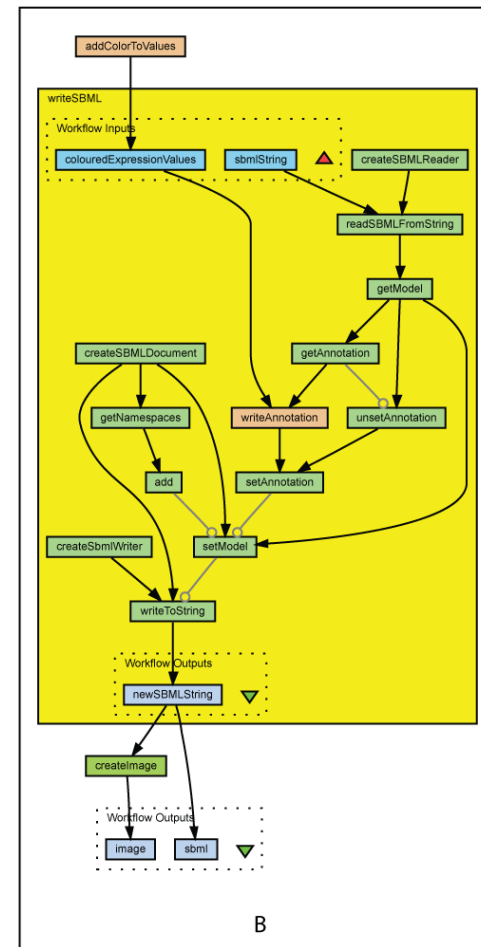
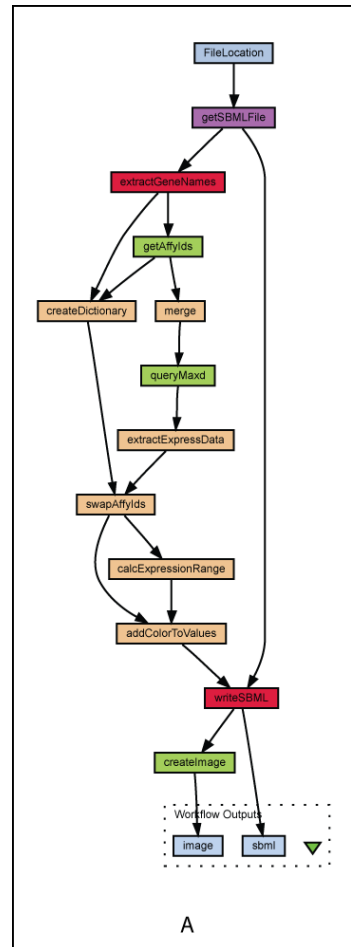
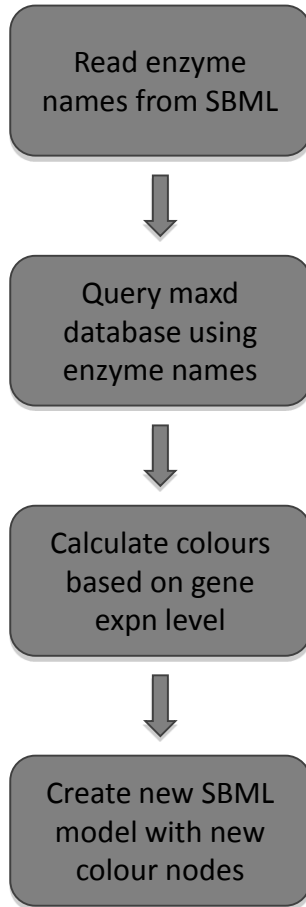


Wei Tan Univ. Chicago



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

Systems Biology Data Integration



Mapping transcriptomics data onto SBML models

Peter Li, Doug Kell, U Manchester



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



Workflows

New/Upload

Workflow



[Katy Wolstencroft](#)

[My Profile](#) [edit]

[My Messages](#)

[My Memberships \(5\)](#)

[My History](#)

[My News](#)

4 new friendship requests

[Onlyhakanboz](#)

[mihaionita_me](#)

[Pankaj chauhan](#)

[Hanny](#)

Search filter terms

« previous **1** 2 3 ... 191 next »

Sort by: Rank

Showing 1905 results. Use the filters on the left and the search box below to refine the results.

Filter by type

- Taverna 2 863
- Taverna 1 645
- RapidMiner 171
- Kepler 43
- Bioclipse Scri... 34
- GWorkflowDL 24
- LONI Pipeline 22
- BioExtract Server 16
- Trident (Packa... 10
- LabTrove Tem... 9

Filter by tag

Taverna 2 **Pathways and Gene annotations for QTL region** View Download (v7)

(v7)

Created: 19/11/09 @ 18:18:52 | Last updated: 02/09/11 @ 11:44:57

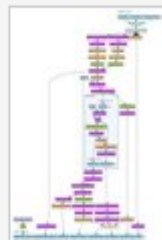
Credits: Paul Fisher

License: Creative Commons Attribution-Share Alike 3.0 Unported License

Original Uploader



[Paul Fisher](#)



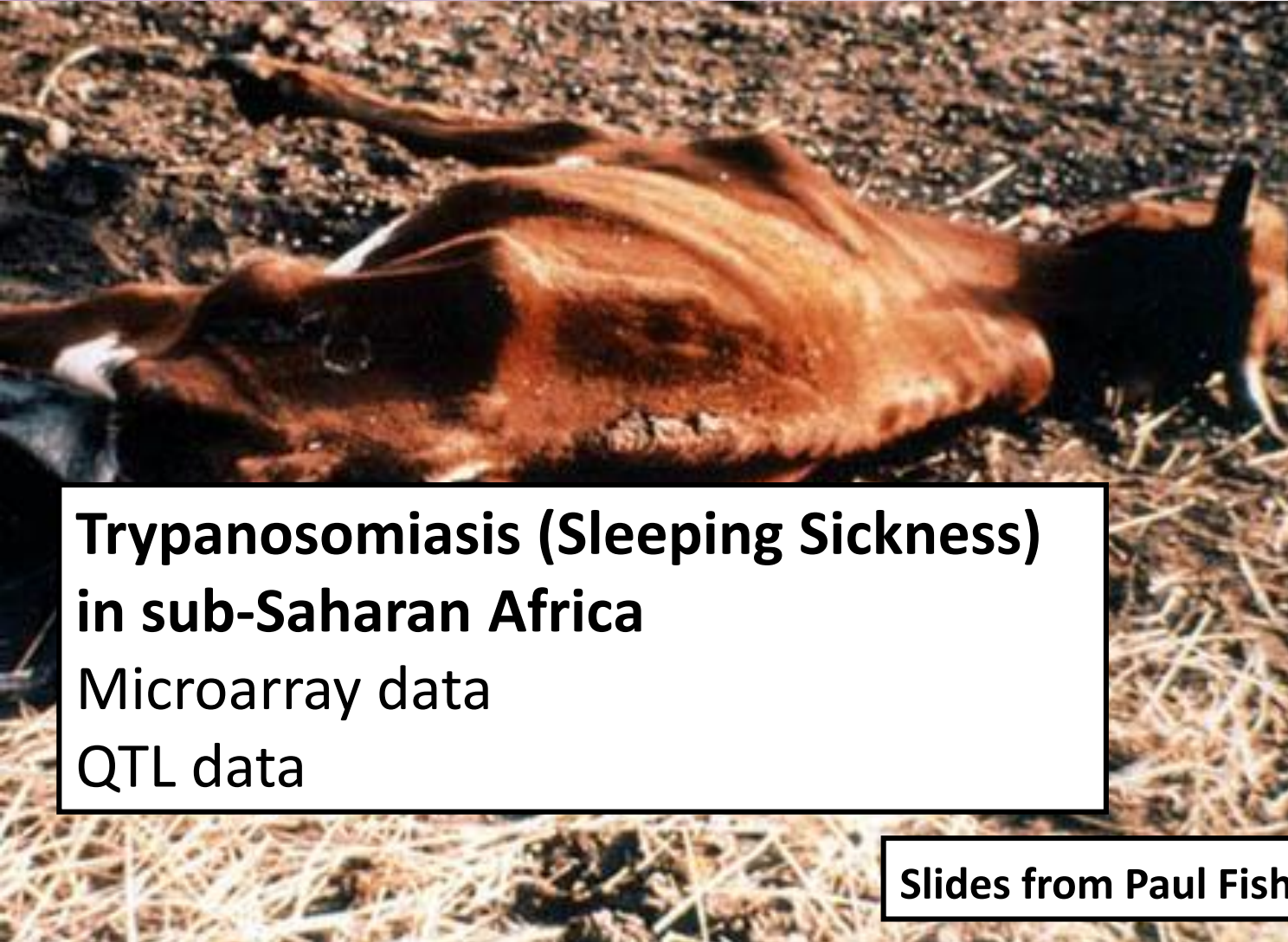
This workflow searches for genes which reside in a QTL (Quantitative Trait Loci) region in the mouse, *Mus musculus*. The workflow requires an input of: a chromosome name or number; a QTL start base pair position; QTL end base pair position. Data is then extracted from BioMart to annotate each of the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

Just Enough Sharing....

- myExperiment can provide a central location for workflows from one community/group
- myExperiment allows you to say
 - Who can look at your workflow
 - Who can download your workflow
 - Who can modify your workflow
 - Who can run your workflow
- Ownership and attribution



The Wellcome Trust Funded Host-Pathogen Project



Trypanosomiasis (Sleeping Sickness) in sub-Saharan Africa

Microarray data
QTL data



Steve Kemp



Andy Brass



Paul Fisher

Slides from Paul Fisher



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



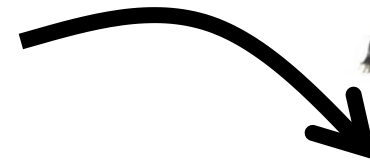
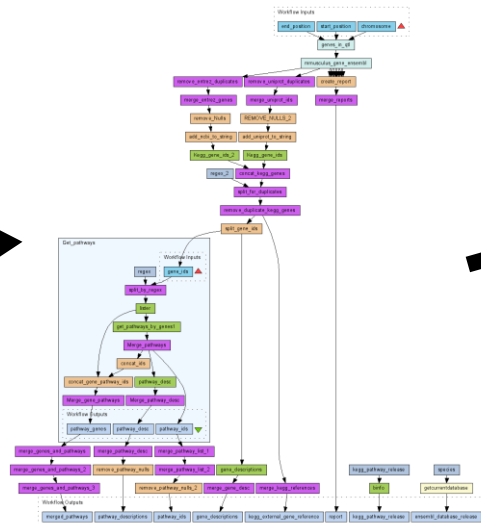
Reuse, Recycle, Repurpose Workflows



Dr Paul Fisher



Identify *QTg* and pathways implicated in resistance to Trypanosomiasis in cattle



Dr Jo Pennock

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



PubMed ID: 20687192

[doi:10.1186/1471-2164-14-127](https://doi.org/10.1186/1471-2164-14-127)

my



Another 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 Grecis
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



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?”



Advanced users design and build workflows (informaticians)

The screenshot shows the 'myexperiment' web interface. At the top, there are navigation tabs for 'Users', 'Groups', 'Workflows', 'Files', 'Blogs', and 'Forums'. Below this, a search bar and a 'New/Upload' button are visible. The main content area displays a 'Workflow Entry: Microarray CEL file'. It includes a 'Version: 2 (latest)' indicator, a 'Change to:' dropdown, and a 'Title: Microarray CEL file to candidate pathways'. The 'Version created on:' is 'Wednesday 03 October 2007 @ 18:35:55 (GMT)'. A 'Diagram: (click to expand)' button is present. Below the title, there is a small thumbnail of the workflow diagram. A 'Description:' section follows, explaining that the workflow takes a URL and returns a list of top differentially expressed genes. A 'Download' button is located at the bottom left. On the right side, a 'Workflow diagram' is displayed, showing a complex flow of services and data. A 'Service panel' is also visible, listing various services like 'Beanshell', 'Nested workflow', 'Rshell', 'Spreadsheetimport', and 'String constant'. A 'Workflow explorer' and 'Validation report' are also shown.

The screenshot shows the Taverna Workbench 2.2.0 software interface. It features a 'Service panel' on the left with a search filter and a list of 'Available services' including 'Beanshell', 'Nested workflow', 'Rshell', 'Spreadsheetimport', and 'String constant'. Below this is a 'Workflow explorer' showing a tree view of the workflow components. The main area displays a 'Workflow diagram' with a complex flow of services and data. The diagram includes components like 'compound', 'getMolarVolumeAndRefractivity', 'getAbrahamDescriptors', 'molarVolume', 'molarRefractivity', 'refractivity', 'volume', 'chemspiderid', and 'abrahamDescriptors'. The workflow starts with 'Workflow input ports' leading to 'compound', which then flows through several other services and components, eventually leading to 'Workflow output ports'.

Intermediate users reuse and modify existing workflows or components

Others “replay” workflows through web page



A Collection of Tools

Workflow Repository



Service Catalogue



Activity and Service Plug-in Manager

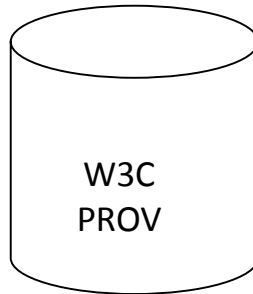


Workflow GUI Workbench and 3rd party plug-ins



Taverna

Provenance Store



Workflow Server



Secure Service Access, and Programming APIs

Client User Interfaces



Web Portals

E-Laboratories

Programming and APIs



Summary – Workflow Advantages

- Informatics often relies on **data integration** and large-scale **data analysis**
- Workflows are a mechanism for **linking** together resources and analyses
- Promote **reproducible** research
- Find and use successful analysis methods *developed by others* with myExperiment



More Information

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



- Using Taverna to design and build workflows
- Reusing workflows from myExperiment
- Finding and using different services:
REST, Xpath, Beanshell, R, ...
- Exploring the workflow engine: iteration, looping, retries, parallel invocation
- Web: Taverna Online, Taverna Player
- Interactions
- Components

