# Going with the Flow Distributed Computing for Systems Biology Using Taverna

Prof Carole Goble

The University of Manchester, UK
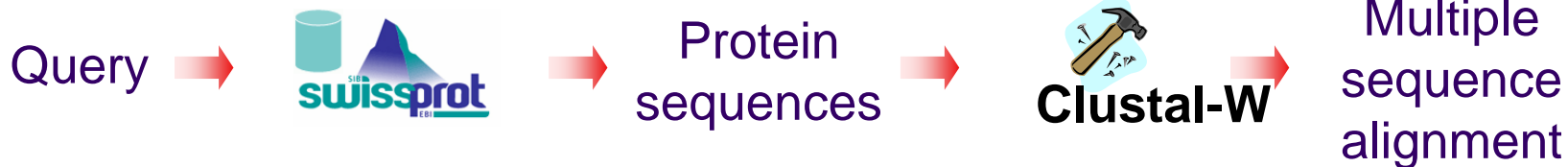
http://www.mygrid.org.uk
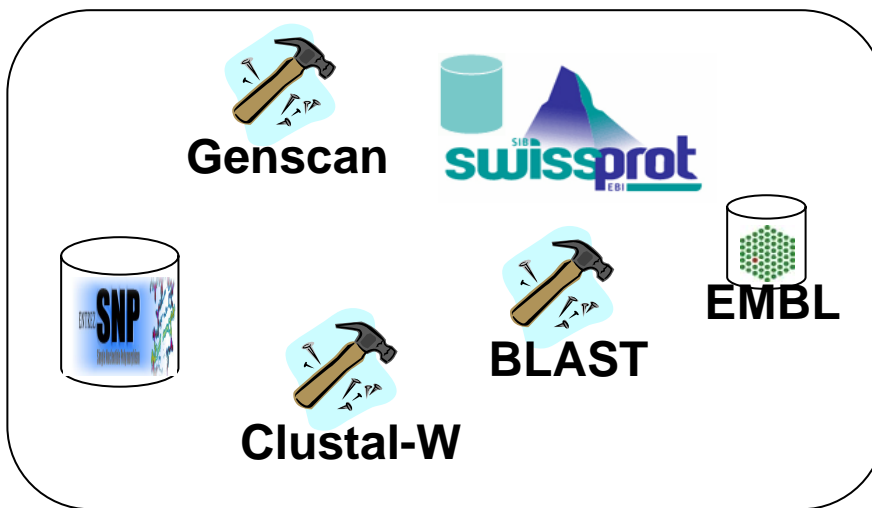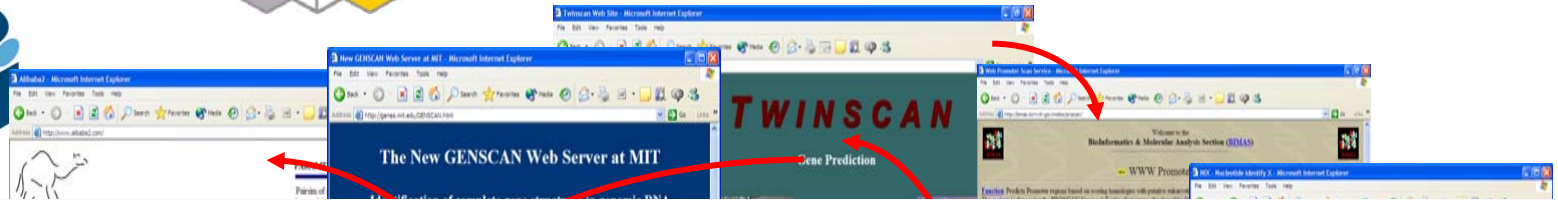http://www.omii.ac.uk

# Data pipelines in bioinformatics

Resources /Services

**Genscan**

**swissprot**

**SNP**

**Clustal-W**

**BLAST**

**EMBL**

Query → **swissprot** → Protein sequences → **Clustal-W** → Multiple sequence alignment
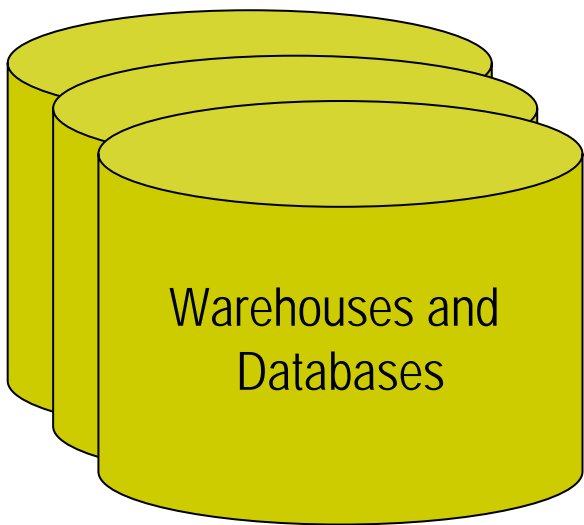
Example *in silico* experiment: Investigate the evolutionary relationships between proteins

[Peter Li]

- Manual creation
- Semi-automation using bespoke software
- Issues:
  - Volatility of data in life sciences
  - Data and metadata storage
  - Integration of heterogeneous biological data
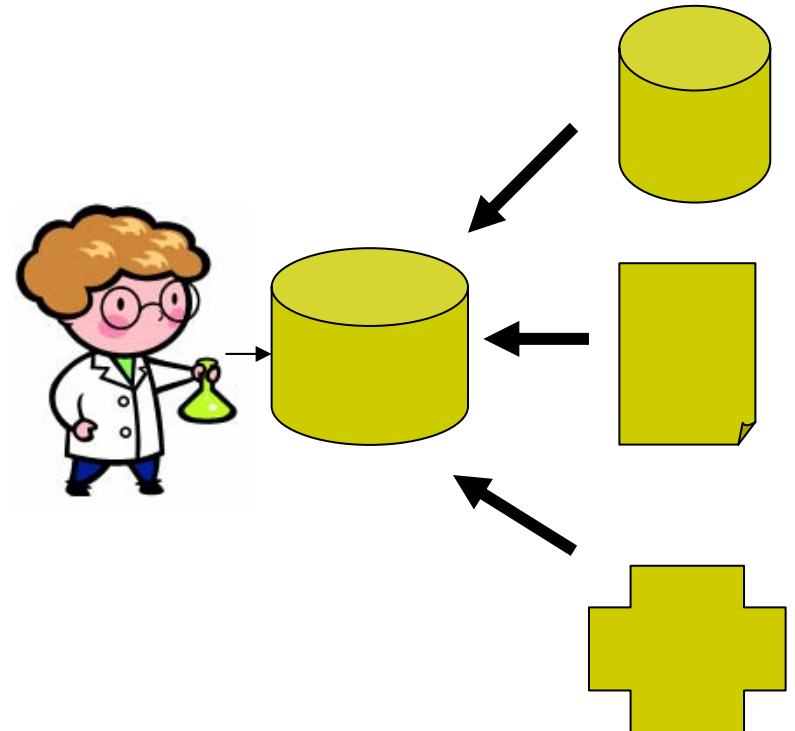  - Visualisation of models
  - Brittleness
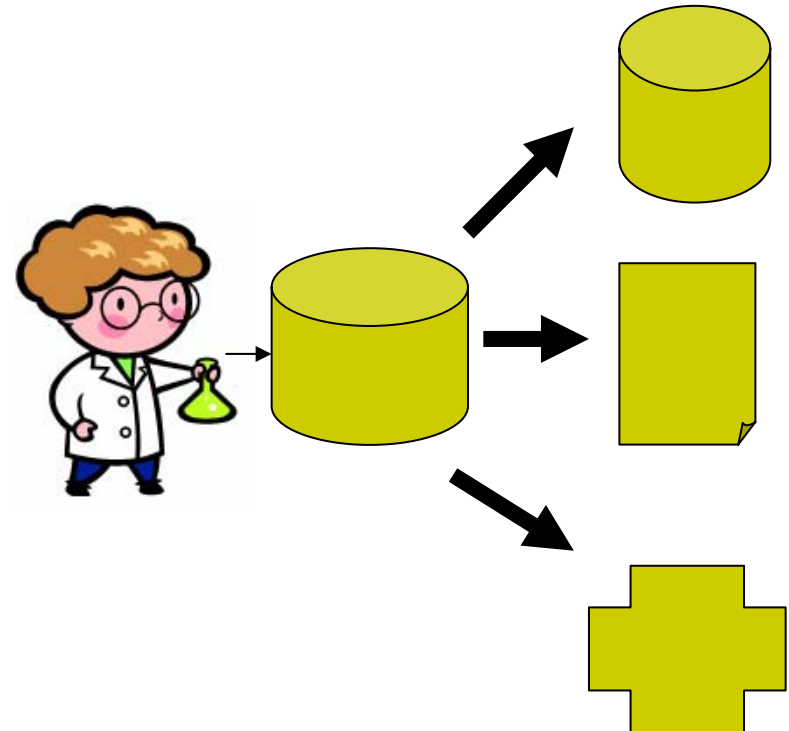
Warehouses and Databases

# Data Warehouse

- Copy the data sets
- Combine them into a pre-determined model before query
- Query that model
- Clean data
- Refresh, Fixed
- High cost, Front loaded
- You can only use what has been set up for you.

# Distributed Database Integration

- Marshal the data sets
- Combine it into a pre-determined model when you query
- Always fresh
- Map from model to databases dynamically
- More flexible but still depends on model
- High cost
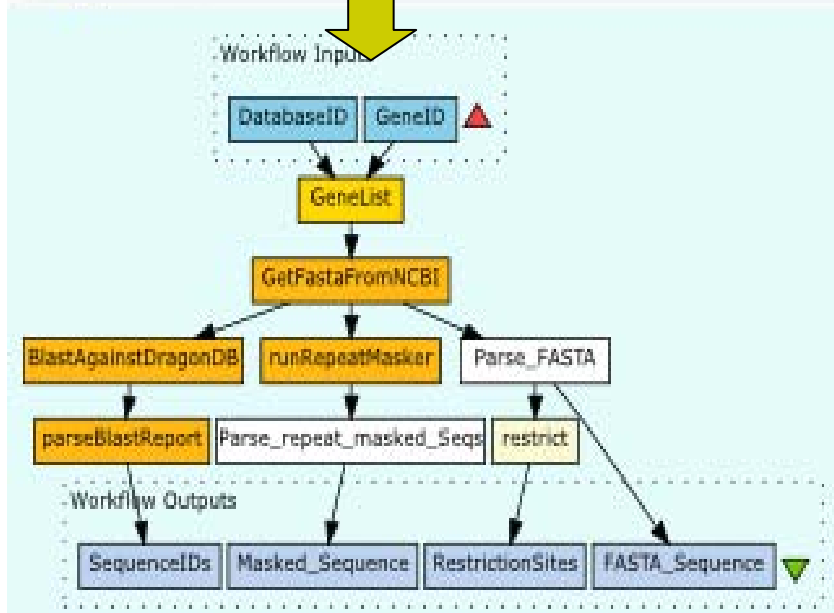- You can only use what has been set up

Protocol

Create a gene list in Excel
Go to NCBI
Retrieve FASTA for each gene
DragonDB Blast each sequence
Copy/paste IDs into a spreadsheet
Run Repeat Masker on each sequence
copy/paste masked sequence into Excel
Run MacVector cut each seq with EcoRI

Warehouses and Databases

[Mark Wilkinson, 2006 BioMOBY]

© omi.

# It would be good if you could systematically automate…

Make data sets / resources / tools / codes / models accessible to a computer.

And cope when they change

And run them where they are hosted ….

# It would be good if you could systematically automate…

….Link together resources

Automate the protocol so I don't have to do it every time I need to repeat the search or re-run the analysis.

And do it accurately and systematically every time without mistakes. And not get bored and sloppy. And be comprehensive too….

# It would be good if you could systematically automate…

…Rerun it over and over and over and over again. Automatically. And keep the log of what actually happened. Automatically.

Manage the results of the protocol. Not just the data results, but the evidence for the results, the source of the data, the log of what you did and why. Helpful when you publish!....

# It would be good if you could systematically automate…

…Record this protocol, share it with colleagues.

Fiddle with it.

Remember what it was 2 weeks later.

Adapt a colleague's or expert's to suit you

Give your protocol to a colleague…

# It would be good if you could systematically automate…

And do it in my lab without having to have a lot of systems administrators and developers building databases for me.

Or writing Perl.

And it runs on my crappy laptop.

# And be …

- Un-biased and Unambiguous in my science

- # Systematic

- Efficient

- Scalable

- Flexible

- Customisable

- Transparent in my scientific method

# The Two W's

- Web Services
    - Technology and standard for exposing code / database with an API that can be consumed by a third party remotely.
    - Describes how to interact with it.
- Workflows
    - General technique for describing and enacting a process
    - Describes *what* you want to do, not *how* you want to do it

# Workflows

Workflow language specifies how bioinformatics processes fit together.

High level workflow diagram separated from any lower level coding – you don't have to be a coder to build workflows.

Workflow is a kind of script or protocol that you configure when you run it.

Easier to explain, share, relocate, reuse and repurpose.

© omii

- myGrid  http://www.mygrid.org.uk

- UK e-Science pilot project since 2001

- Build middleware for Life Scientists that enables them to undertake *in silico* experiments and share those experiments and their results.

- Individual scientists, in under-resourced labs, who use other people's applications.

- Open source.

- Workflows.

- Data flows. Ad hoc & exploratory

**Taverna Workflow Workbench**
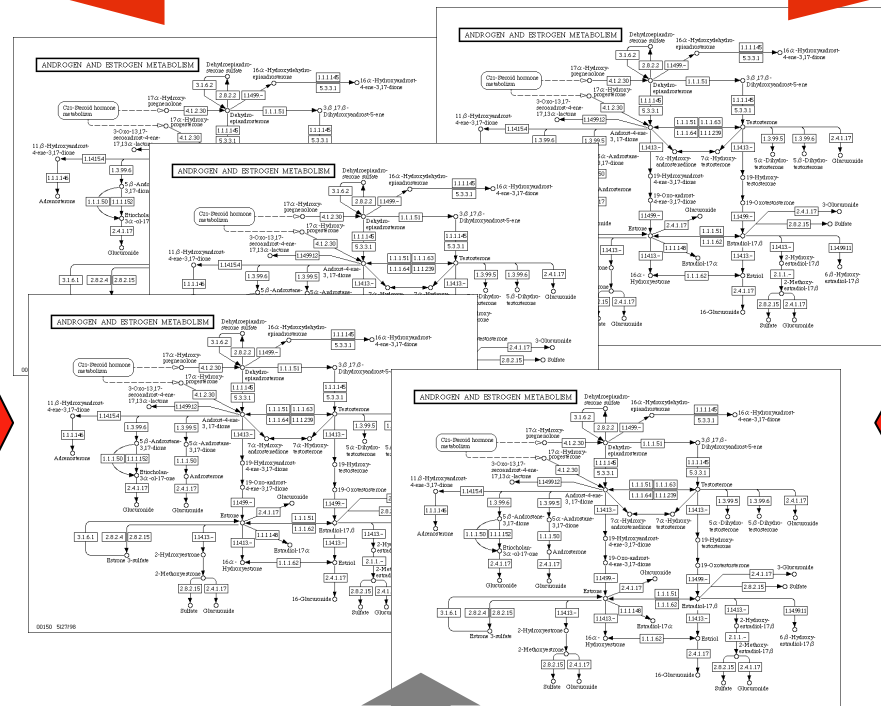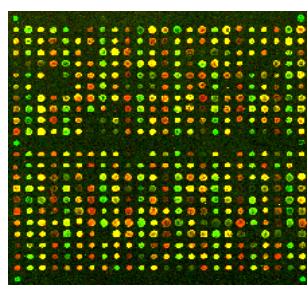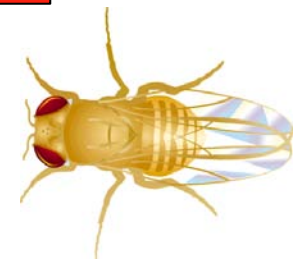**http://taverna.sourceforge.net**

**Genotype** ⟷ **Phenotype**

**200**

**?**

Genes captured in microarray experiment and present in QTL region

Microarray + QTL

Phenotypic response investigated using microarray in form of expressed genes or evidence provided through QTL mapping

19

Key:

**A** – Retrieve genes in QTL region

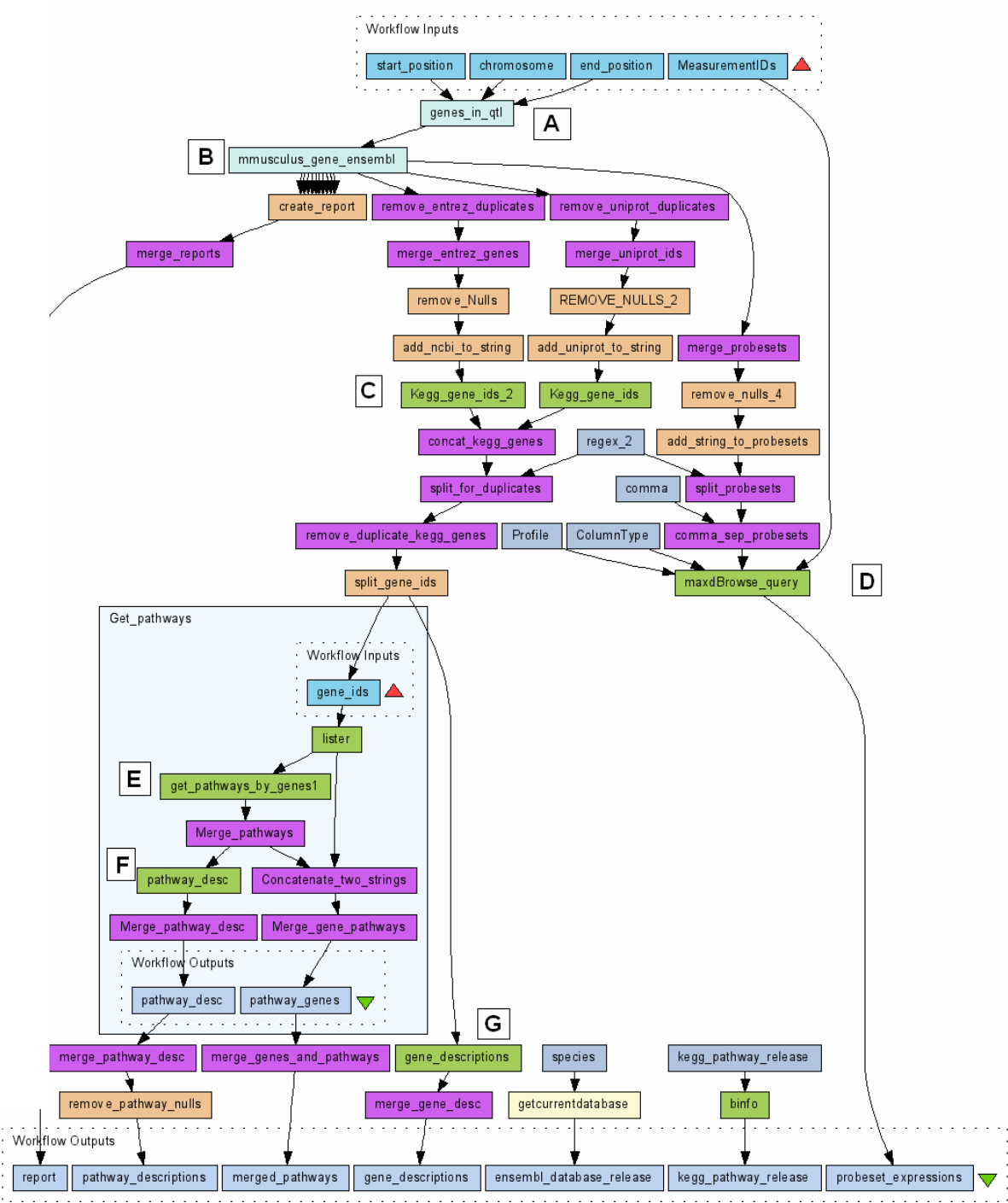**B** – Annotate genes with external database Ids

**C** – Cross-reference Ids with KEGG gene ids

**D** – Retrieve microarray data from MaxD database

**E** – For each KEGG gene get the pathways it's involved in

**F** – For each pathway get a description of what it does

**G** – For each KEGG gene get a description of what it does

[Andy Brass, Steve Kemp, Paul Fisher, 2006]

# Result

- Captured the pathways returned by QTL and Microarray workflows over the MaxD microarray database

- Identified a pathway for which its correlating gene (Daxx) is believed to play a role in trypanosomiasis resistance.

- Manually analysis on the microarray and QTL data had failed to identify this gene as a candidate.
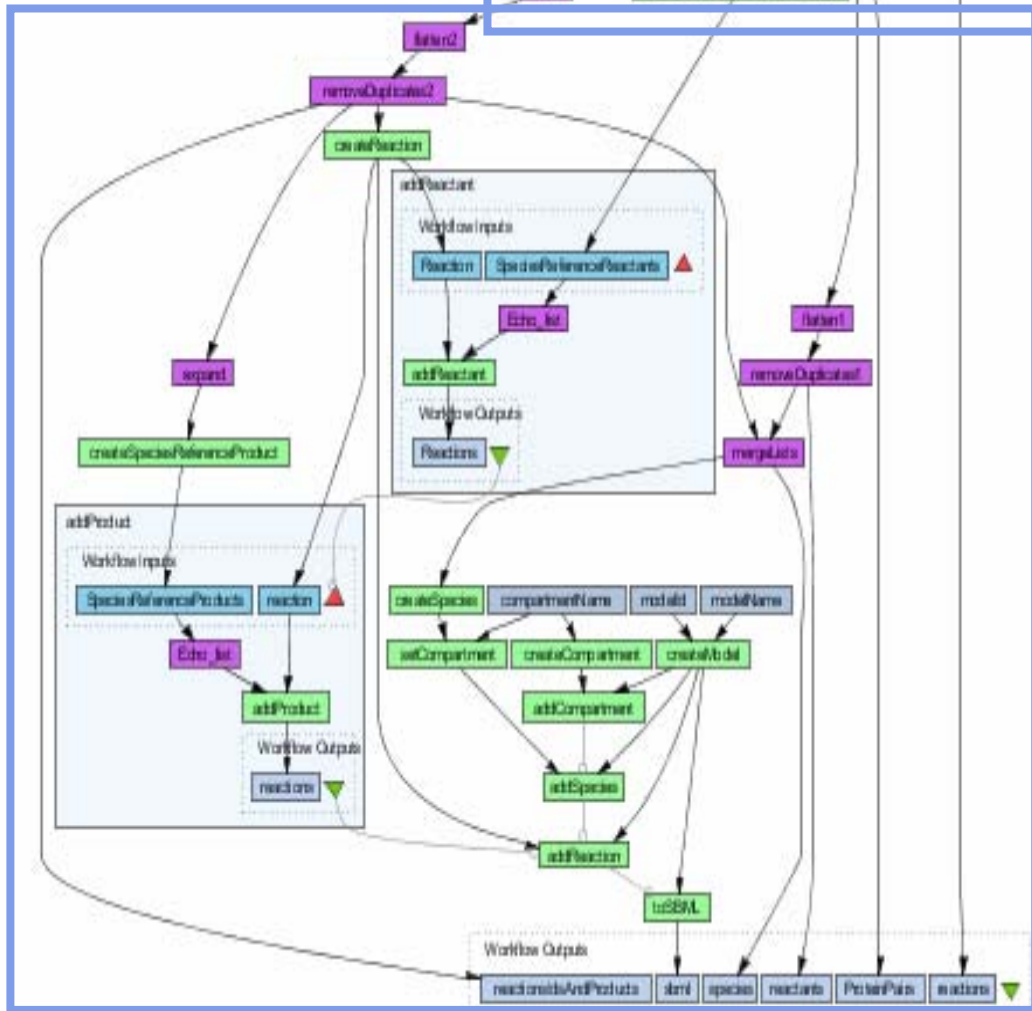
© omii

[Andy Brass, Steve Kemp, Paul Fisher, 2006]

# **Trichuris muris (mouse whipworm) infection**

- 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

- Workflows: trypanosomiasis cattle experiment, was reused without change.

- Analysis of the data by a biologist found the processes in a couple of days.

© omii

[Joanne Pennock, Paul Fisher, 2006]

# Pull Public Databases
# +
# inHouse Data
# =
# Model

Core SBML model construction workflow

# Visualise results using routine SBML tools

SharkView – interactive SBML viewer

[Peter Li, Doug Kell, 2006]
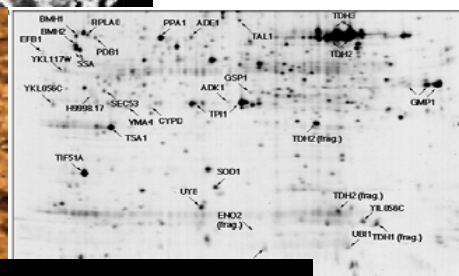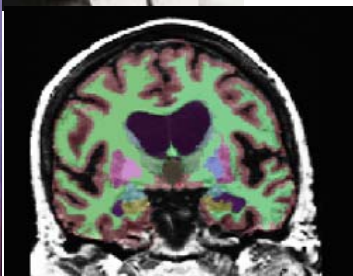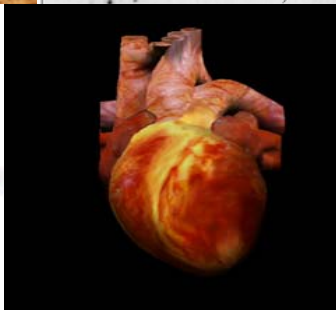
# Model construction: Post-Taverna

- Captures the scientific process of model construction as workflows

- Workflows enacted 'on demand' to construct most up-to-date models using the latest data

- Models are pushed into a data model of choice

- Provide various ways of visualising models

# Multi-disp.

- ~20000 downloads
- Users in US, Singapore, UK, Europe, Australia,
- Systems biology
- Proteomics
- Gene/protein annotation
- Microarray data analysis
- Medical image analysis
- Heart simulations
- High throughput screening
- Phenotypical studies
- Plants, Mouse, Human
- Astronomy
- Dilbert Cartoons

26

**A workflow marketplace**

Finding and Sharing Tools

Taverna Workbench

3rd Party Applications and Portals

myExperiment

DAS

Feta

Utopia

Workflow enactor

Clients

Service Management

LSIDs

Log Meta data

Default Data Store

Custom Store

© omii

KAVE

BAKLAVA

Results Management

# Transparency

## Illuminating the black box

Note to biologists: submissions to *Nature* should contain complete descriptions of materials and reagents used.

This journal aims to publish papers that are not only interesting and thought-provoking, but reproducible and useful. In order to do this, novel materials and reagents need to be carefully described and readily available to interested scientists.

That might seem obvious. But despite the efforts of our editors and referees, papers in the biological sciences are still being submitted — and occasionally published — that do not adequately describe the reagents used. Unless efforts are redoubled to eliminate this practice, we could see an era of 'black box' biology, in which outside researchers cannot work out what was done in an experiment.

established didn't want the author to reveal the sequences, as this would jeopardize its *raison d'être*. This kind of stalemate matters, because it prevents the replication of experiments and inhibits the

29

# Provenance

- ## Who, What, Where, When, Why?, How?
- Context
- Interpretation
- Logging & Debugging
- Reproducibility and repeatability
- Evidence & Audit
- Non-repudiation
- Credit
- Credibility
- Accurate reuse and interpretation
- Smart re-running
- Cross experiment mining
- Just good scientific practice

© omii



BioMOBY

Protocol

Create a gene list in Excel
Go to NCBI
Retrieve FASTA for each gene
DragonDB Blast each sequence
Copy/paste IDs into a spreadsheet
Run Repeat Masker on each sequence
copy/paste masked sequence into Excel
Run MacVector cut each seq with EcoRI

# Tracking

From which Ensembl gene does pathway come from mmu004620 ?

# Workflows over Results



**Automatically** backtrack through the data provenance graph

# An Open World

- **Open** domain services and resources.
- Taverna accesses 3000+ operation.
- Third party.
- All the major providers
  - NCBI, DDBJ, EBI …
- Enforce NO common data model.

- Quality Web Services considered desirable

.

# If you don't provide a Web Service Interface…

- SoapLab



http://www.ebi.ac.uk/soaplab/

- Java API Consumer



import Java API of libSBML as workflow components

© 

# Shield the Scientist

Bury the complexity



**Workflow enactor**

| Processor | Processor | Processor | Processor | Processor | Processor | Processor | Styx | Processor | **...** |

| Bio MART | WSRF | Plain Web Service | Soap lab | Bio MOBY | Local Java App | Enactor | Styx client | R package | **...** |

# User Interaction

- Allows a workflow to call out to an expert human user

- E.g. Used to embed the Artemis annotation editor within an otherwise automated genome annotation pipeline

[University of Bergen]

# No miracles here.

- Building good workflows
  - Pattern books
  - Best practice
  - Workflow packs
- Data integration
  - Still have to think about building models of results
- Services
  - Properly computer-accessible (Web) services
  - Maintenance

# Changes to Scientific Practice

- Systematic and comprehensive automation.
  - Eliminated user bias and premature filtering of datasets and results leading to single sided, expert-driven hypotheses
- Dry people hypothesise, wet people validate.
  - "make sense of this data" -> "does this make sense?"
- Workflow factories.
  - Different dataset, different result
- Workflow market.
- Accurate provenance.

# Conclusions

Distributed computing

1. Web Services

 Make your data or your code accessible to be a component in a …

2. Workflow

 For flexible, transparent and systematic encoding of protocols for linking services/processes up

Taverna http://taverna.sourceforge.net

ᵐʸGrid http://www.mygrid.org.uk

OMII-UK http://www.omii.ac.uk

© omii

# **Acknowledgements**

- Phase1 ᵐʸGrid researchers, Phase2 OMII-UK, ᵐʸGrid Research Team
- Tom Oinn (EBI), Martin Senger, Katy Wolstencroft
- Peter Li, Paul Fisher, Andy Brass, Robert Stevens, Mark Wilkinson
- EPSRC, Wellcome Foundation

Katy Wolsencroft

Tom Oinn



open middleware
infrastructure institute uk
www.omii.ac.uk