# Taverna and ᵐʸGrid

## Open Workflow for Life Sciences

Tom Oinn

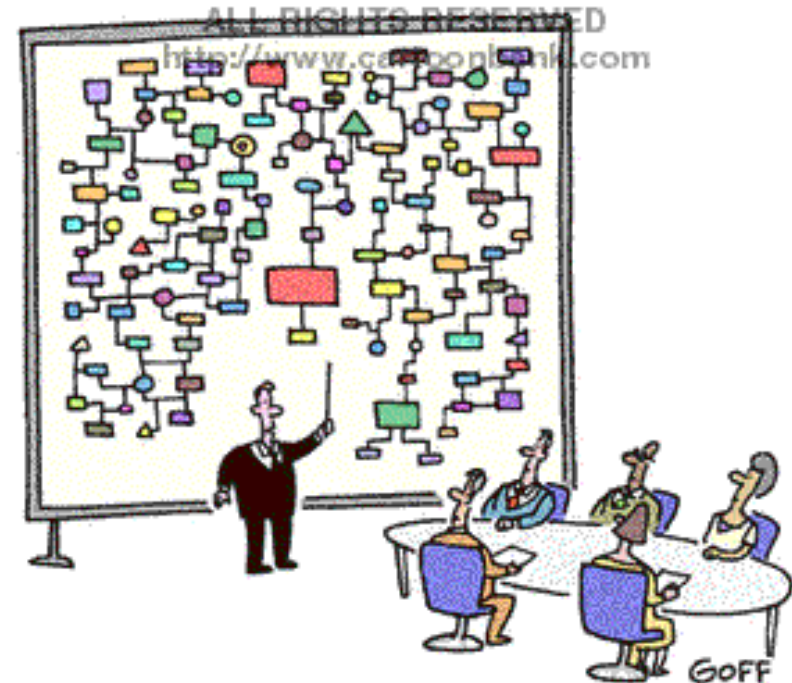tmo@ebi.ac.uk

# What, who, why?

- Taverna – a workflow development and enactment environment
- Who – part of myGrid, an EPSRC funded UK eScience Pilot project coordinated by Carole Goble at Manchester University
- Why – because bioinformatics is hard enough without turning users into web spiders ☺

GOFF
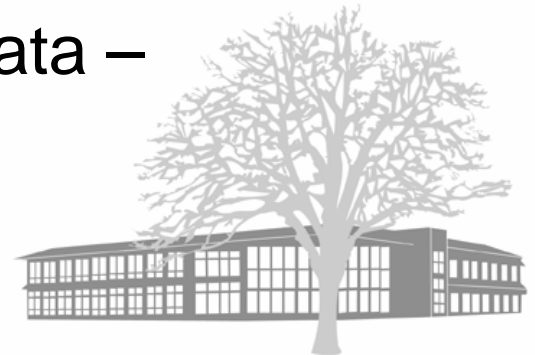
"And that's why we need a computer."

# Old approach

- Cut and paste, cgi, shell scripting, ftp, excel
- Time intensive
- Manual process, fails to scale sensibly
- Hard to document and reproduce
    - Good scientific discipline hard to maintain
- Boring, waste of highly trained scientists

# Our approach

- Capture the scientific method as a formal process model

- Allow users to construct such models from libraries of available components in a graphical editing environment with semantic support

- Publish process definitions as scientific methods, enact and automatically scale to large data sets, multiple runs

- Automatically collect enactment metadata – workflow provenance.

Tools and Workflow Invocation

## Advanced model explorer

Workflow | Remote resource usage

Load | Load from web | Save | New subworkflow | ☐ Offline | Reset ✖

| Workflow object | Retries | Delay | Back... | Thre... | Critical |
|---|---|---|---|---|---|
| Workflow model | | | | | |
|   Workflow inputs | | | | | |
|   Workflow outputs | | | | | |
|     Graph | | | | | |
|   Processors | | | | | |
|     Green : chartreuse3 | 0 | 0 | | | |
|     PassAllTerms | 0 | 0 | | | |
|     Red : crimson | 0 | 0 | | | |
|     PassUniqueTerms | 0 | 0 | | | |
|     GetUniqueIDs | 0 | 0 | | | |
|       in I('text/plain') | | | | | |
|       out I('text/plain') | | | | | |
|     GenericSetOperations | 0 | 0 | | | |
|     ShowOnlyUniqueTerms : true | 0 | 0 | | | |
|     GetUniqueIDs1 | 0 | 0 | | | |
|     GenericSetOperations1 | 0 | 0 | | | |
|     Purple : purple | 0 | 0 | | | |
|     Yellow : gold | 0 | 0 | | | |
|     GetUniqueIDs2 | 0 | 0 | | | |
|     FlattenList | 0 | 0 | | | |

## Advanced model explorer

Workflow | Remote resource usage

Save HTML description

2. *Terms in red are those implied by those in purple and not explicitly or implicitly mapped to any genes in X*
3. *Yellow terms are those explicitly mapped in X*
4. *Green terms are those implied by terms mapped in both X and Y*
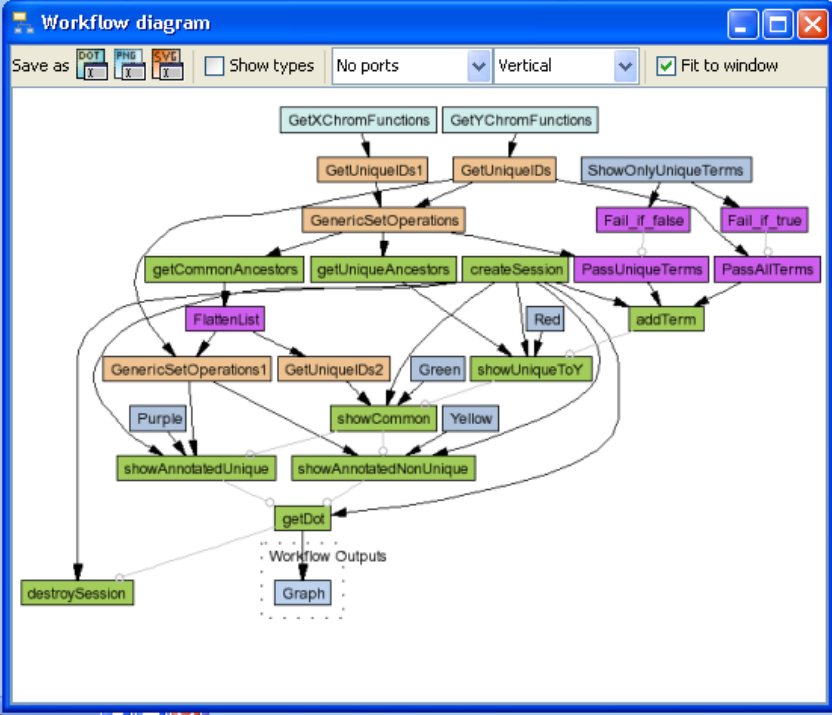
### Resource usage report

This display shows the various external resources used by the current workflow. It does not show resources such as local operations or string constants which are run within the enactment engine. Services are categorized by resource host and type, and the name of the instance of each service shown to the right.

Resources on **martdb.ebi.ac.uk**, 2 instances.

| | Processor |
|---|---|
| _ensembl | GetYChromFunctions |

| | Processor |
|---|---|
| _ensembl | GetXChromFunctions |

k, 10 instances.

at */collab/mygrid/service1/goviz/GoViz.jws?wsdl*

## Available services

Search list | db | 🔍 | ☑ Watch loads

- 📁 Available Processors
  - 📁 Local Services
  - 📁 WSDL @ http://www.ebi.ac.uk/collab/mygrid/service1/goviz/G
  - 📁 WSDL @ http://www.ebi.ac.uk/xembl/XEMBL.wsdl
  - WSDL @ http://soap.genome.jp/KEGG.wsdl
    - porttype: KEGGPortType [RPC]
      - list_databases
      - list_organisms
      - list_pathways
      - list_ko_classes
      - binfo
      - bget
      - bfind
      - btit
      - get_linkdb_by_entry
      - get_best_neighbors_by_gene
      - get_best_best_neighbors_by_gene
      - get_reverse_best_neighbors_by_gene
      - get_paralogs_by_gene
      - get_motifs_by_gene
      - get_genes_by_motifs
      - get_ko_by_gene
      - get_ko_by_ko_class
      - get_genes_by_ko
      - get_genes_by_ko_class

⏸ Pause | ⏹ Stop

## Workflow diagram

Save as | DOT | PNG | SVG | ☐ Show types | No ports | Vertical | ☑ Fit to window



| ... | Name | Last event | Event timestamp | Event detail | ... |
|---|---|---|---|---|---|
| | GetUniqueIDs1 | ProcessComplete | 03-Jun-2005 09:... | | |
| | addTerm | InvokingWithIteratio | 03-Jun-2005 09:... | IterationNumber='1' IterationTot... | |
| | createSession | ProcessComplete | 03-Jun-2005 09:... | | |
| | getCommonAnce... | InvokingWithIteratio | 03-Jun-2005 09:... | IterationNumber='1' IterationTot... | |
| | showAnnotatedU... | ProcessScheduled | 03-Jun-2005 09:... | | |
| | getDot | ProcessScheduled | 03-Jun-2005 09:... | | |

Graph | Intermediate inputs | Intermediate outputs
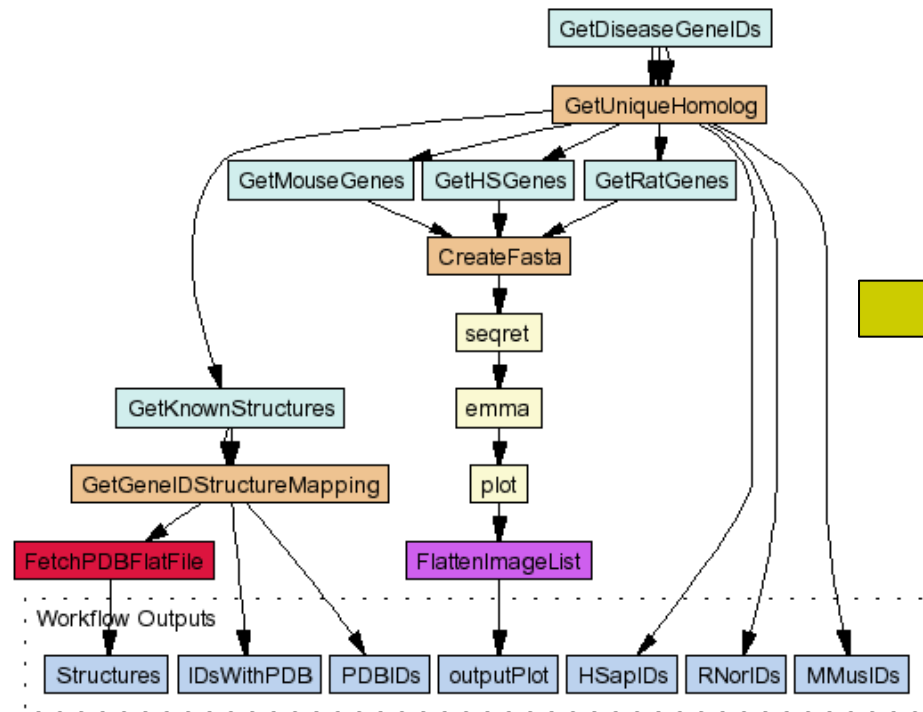
set2 | set1

| List | Self |
|---|---|
| urn:lsid:www.mygrid.org.uk:documentco | GO:0004674 |
| text/plain | GO:0005554 |
| GO:0004221 | GO:0000785 |
| urn:lsid:www.mygrid.org.uk:lsdoc | GO:0003682 |

🦗 Confi...

# What can we integrate?

- Web services defined by WSDL
  - *Pathport, BIND, Gene Ontology, DBFetch, FASTA, InterproScan, NCBI eUtils…*
- Complex analysis services conforming to Life Science Analysis Engine (LSAE) specification
  - *EMBOSS, Jess, any arbitrary legacy C, PERL or Shell script*
- BioMoby services (www.biomoby.org)
  - *PlaNeT, IRI, Spanish Bioinformatics Network, Genome Prairie…*
- Biomart Database Queries
  - *Ensembl, DbSNP, VEGA…*
- Local embedded scripts via Java, Perl, Python, Ruby etc.
- Seqhound Genomic data warehouse
  - *Genbank, LocusLink, GO*
- Styx Grid Service
  - *Environmental eScience, ocean temperature analysis etc*
- Arbitrary 3rd Party APIs i.e. BioJava, JUMBO, caBIG

# Comparative Genomics



*BiomartAndEMBOSSAnalysis.xml*

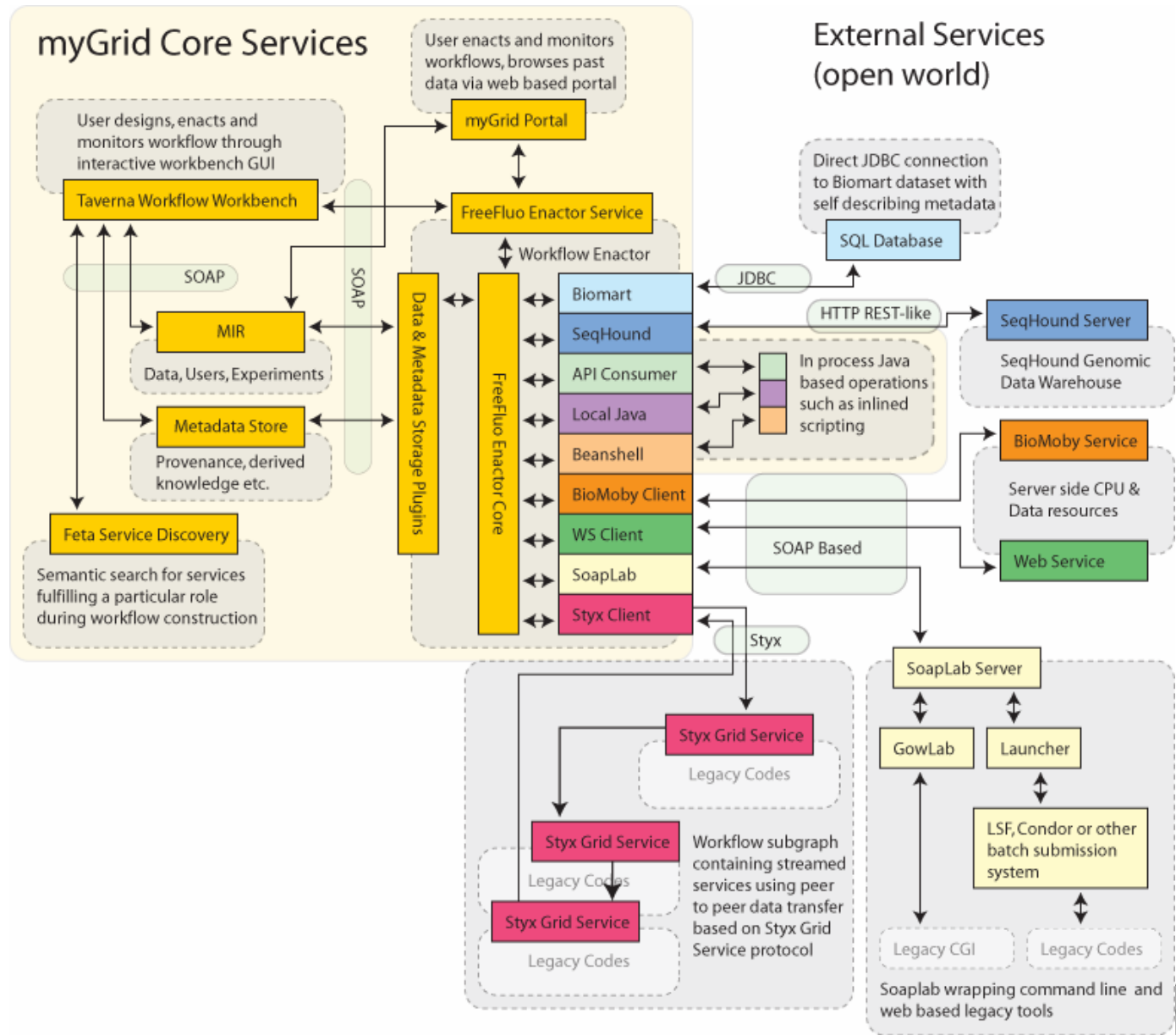# Functional Analysis Workflow…



*CompareXandYFunctions.xml*

# …and result

# Philosophy

- Open world approach for services
  - Do not require service providers to change
  - Maximize interoperability
  - Extend on demand
    - Minimalist functional core, declarative language, many plugin extension points
- Open development approach as well
  - LGPL License
  - Transparent, public development process
  - CVS, Mailing lists, website are all public at all times
  - Avoid institutional 'ownership' of code to safeguard long term future development

Taverna network architecture diagram

# Implicit Iteration

- Allows services to consume collections of items without service modification
- Equivalent to higher order map functions
- Graphical configuration
- Intuitively understood by our user community
- Scares computer scientists ☺

# Workflow summary views

- Diagram and HTML report of the structure of and resources used by the workflow

- Intended to be added to papers, websites etc.

- Can be used by portals, workflow repositories

- Supports reuse – very important!

# Semantic and Naïve Search

- Find services by name or…

- …by function, input types, resources

# Successful?

- Over 1200 downloads of the workbench software for release 1.0
  - Averaging 10-15 downloads / day for release 1.1
  - Slightly scary 220 downloads in three days for 1.2 ☺
- Over 100 active mailing list participants
- Over 1300 available services
- Used across the world in widely differing projects, mostly but not all in bioinformatics (some cheminformatics)
- Active external developer community!

# Taverna User Support

- Taverna has a self supporting user community

- Access help from other users and from the project developers via our mailing lists

- All accessible from http://taverna.sf.net

- We have a user manual! Please use it ☺

# Where next?

- Funding
  - Core myGrid project has completed (3 years)
  - Follow-on platform grant for core team until 2008
  - Associated consumer / helper projects
    - Comparagrid, EMBRACE, iSpider…
- Will be used to…
  - Enhance the scalability of the workflow core
  - Investigate new interfaces (Dalec, Data driven workbench…)

# Schedule

- 1.3 Release in September
  - Final version 1 release
- Moving to 2.0 with new workflow core by end 2005

# Acknowledgements

<sup>my</sup>Grid is an EPSRC funded UK eScience Program Pilot Project

Particular thanks to the other members of the
Taverna project, http://taverna.sf.net