



myGrid

Katy Wolstencroft
University of Manchester





Background

- myGrid middleware components to support *in silico* experiments in biology
- Originally designed to support bioinformatics
 - chemoinformatics
 - health informatics
 - medical imaging
 - integrative biology



History

EPSRC funded UK eScience Program Pilot Project





myGrid in OMII-UK

myGrid



OMII Stack



OGSA-DAI

March 2006





Virtual Grid of Resources

- Biology knowledge-rich
- Applying prior knowledge to new data
- myGrid middleware to enable interoperation between distributed data and resources – a grid of data – not a grid of resources



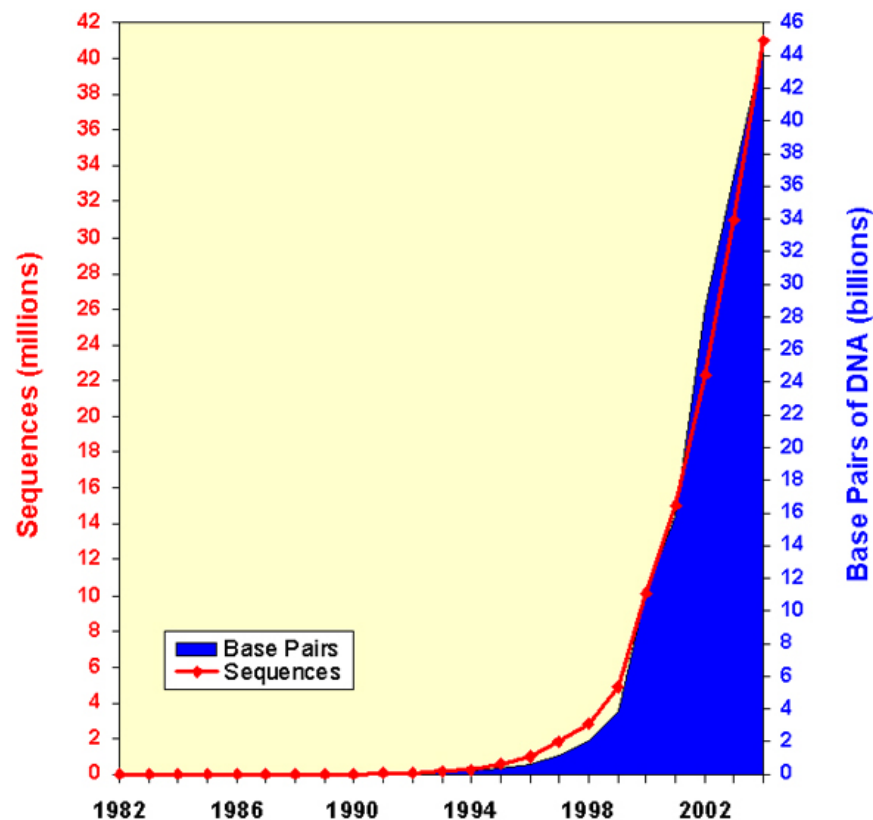


Lots of Resources

NAR 2005 – over 700 databases



Growth of GenBank
(1982 - 2004)





The User Community

Bioinformatics is an open Community

- Open access to data
- Open access to resources
- Open access to tools
- Open access to applications

Global *in silico* biological research





The User Community Problems

- Everything is Distributed
 - Data, Resources and Scientists
- Heterogeneous data
- Very few standards
 - I/O formats, data representation, annotation
 - Everything is a string!

Integration of data and interoperability of resources is difficult

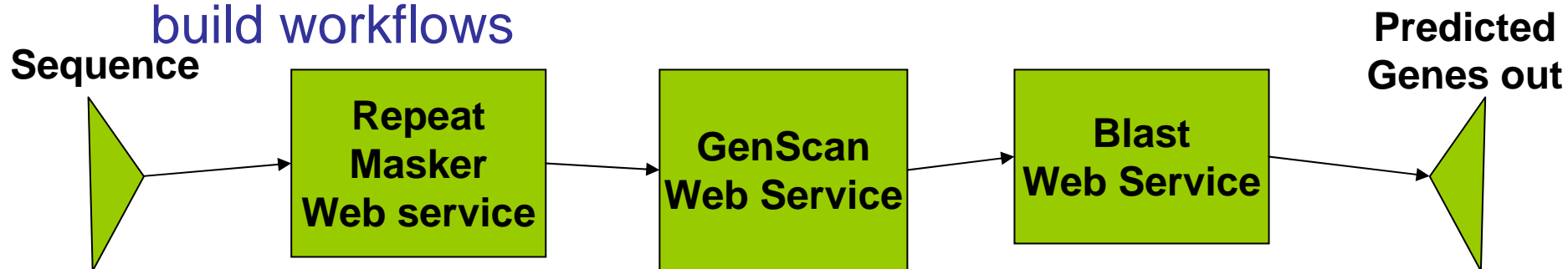


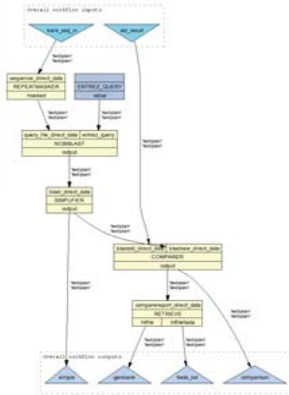


myGrid Approach - Workflows

General technique for describing and enacting a process describes *what* you want to do, not *how* you want to do it
Simple language specifies how bioinformatics processes fit together – *processes are web services*

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





SCUFL

Taverna Workbench



Scufi + Workflow Object Model

Application data flow layer
Scufi graph + service introspection



Workflow Execution

Execution flow layer
List management; implicit iteration mechanism; MIME & semantic type decoration; fault management; service alternates



Freefluo Workflow enactor



Processor



Processor



Processor



Processor



Processor



Processor



Processor

Processor invocation layer



Bio MART



Seq Hound



Plain Web Service



Soap lab



Bio MOBY



Local App



Enactor



Taverna Workflow Components



The screenshot shows the Taverna Workbench interface with several windows:

- Workflow diagram:** A flowchart starting with 'Workflow Inputs' leading to a 'seq' process, then 'seqfile_direct_data', followed by a 'tess' process (subdivided into 'tessstabular', 'tessbigjava', 'tesssmalljava', 'tessmodel'), and finally 'Workflow Outputs' (table, big, small, model).
- Enactor invocation:** A table showing process completion status and timestamps.
- Advanced model explorer:** A tree view of workflow objects and available services.
- Run Workflow:** A window for loading inputs and running the workflow.

Freefluo

Freefluo
Workflow
engine to run
workflows

Scufl Simple Conceptual Unified Flow Language
Taverna Writing, running workflows & examining results
SOAPLAB Makes applications available

Web Service	e.g. DDBJ BLAST
SOAPLAB Web Service	Any Application



So many services – semantic discovery

Over 3000 services

SeqHound –

Database of biological sequences and tools

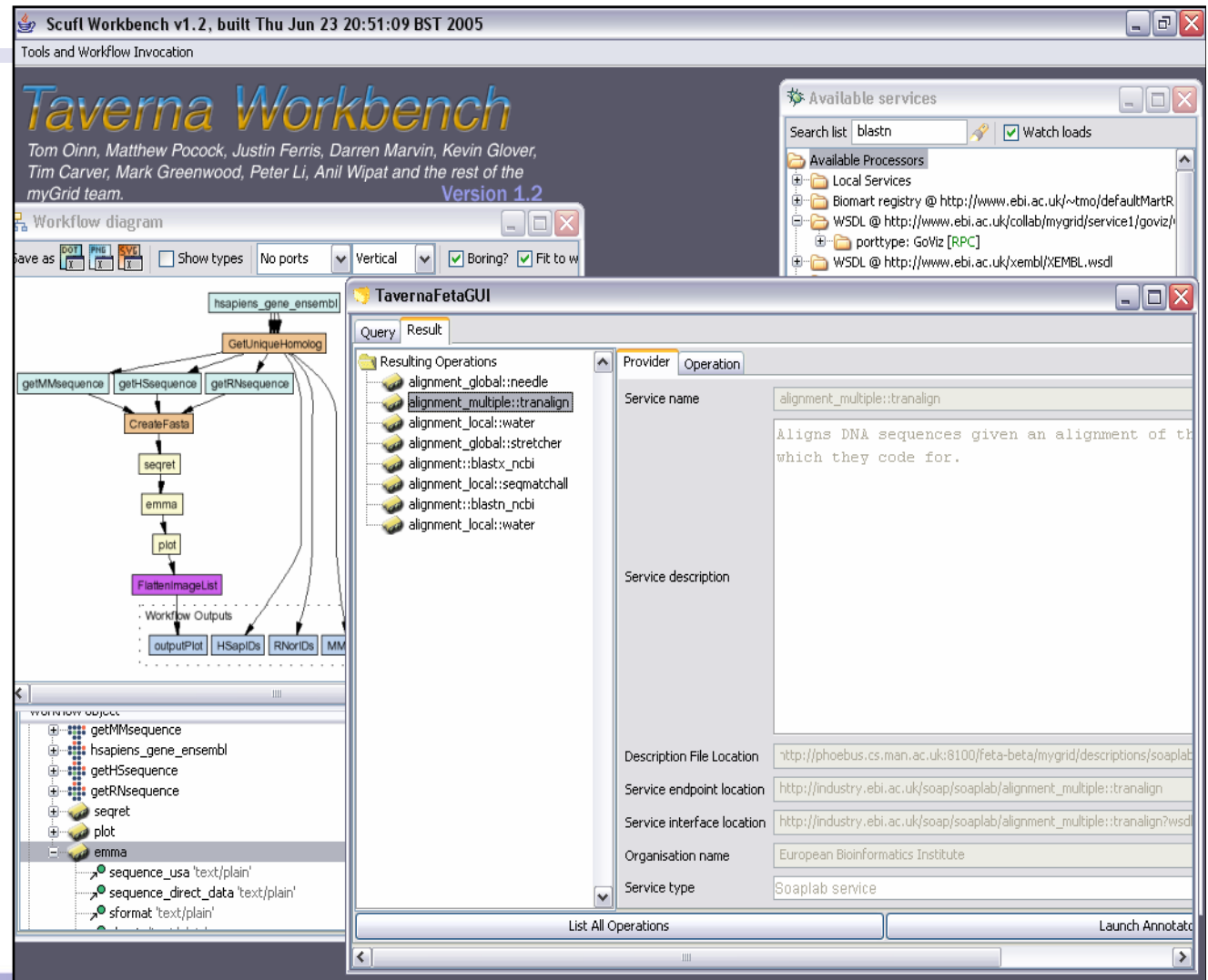
BioMart – Federated query system

EMBOSS – Sequence analysis tools

BioMoby – Collection of web services

EBI SOAPLAB –

Collection of supported services

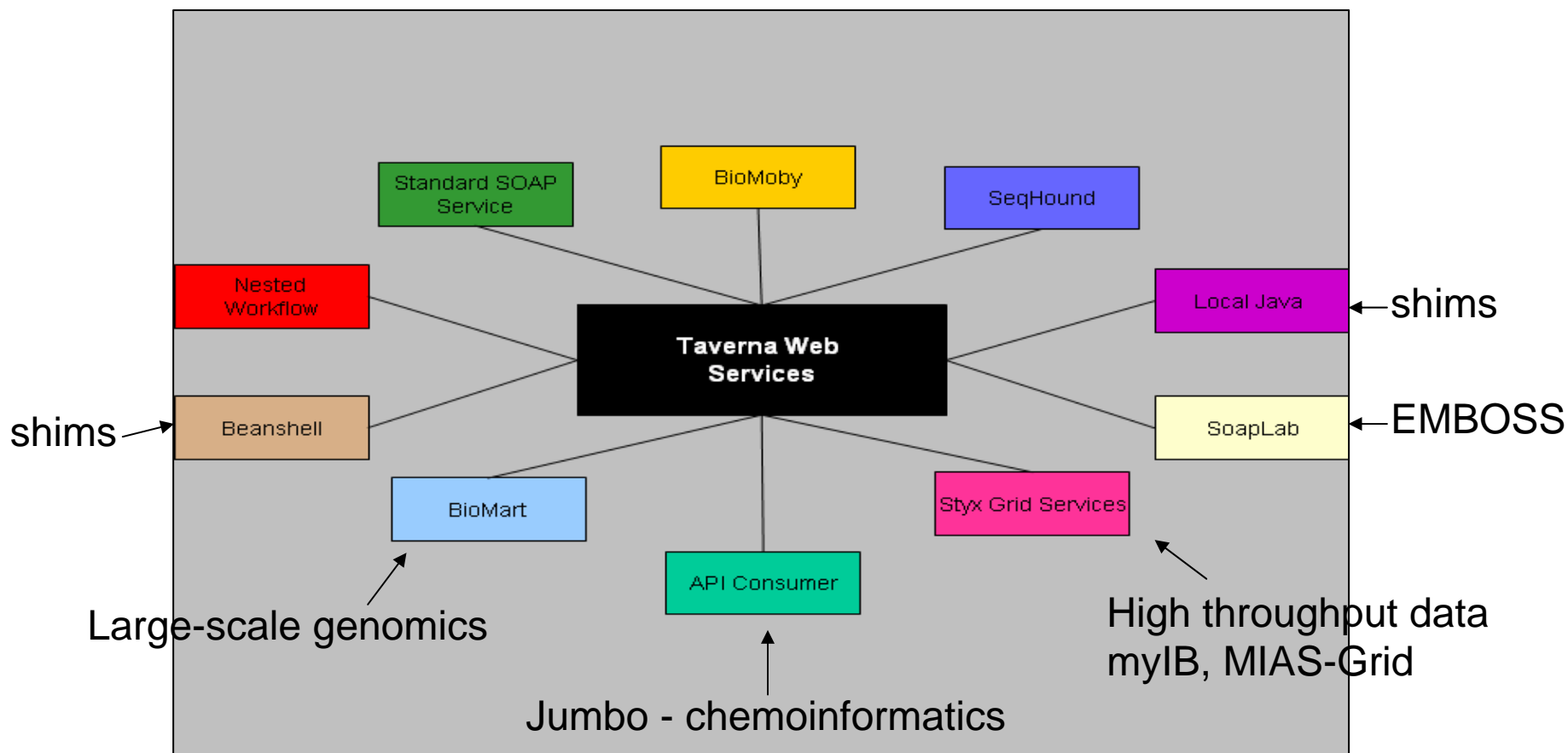


The screenshot displays the Taverna Workbench v1.2 interface. At the top, it says "Taverna Workbench" and lists the authors: Tom Oinn, Matthew Pocock, Justin Ferris, Darren Marvin, Kevin Glover, Tim Carver, Mark Greenwood, Peter Li, Anil Wipat and the rest of the myGrid team. The main window is divided into several panes:

- Workflow diagram:** Shows a flowchart starting with "hsapiens_gene_ensembl" leading to "GetUniqueHomolog", which then branches into "getMMsequence", "getH5sequence", and "getRNsequence". These lead to "CreateFasta", "seqret", "emma", and "plot". The final outputs are "outputPlot", "HSapIDs", "RNorIDs", and "MM".
- Available services:** A search window with "blastn" entered. It lists several services, including "Local Services", "Biomart registry", and "WSDL @ http://www.ebi.ac.uk/collab/mygrid/service1/goviz/".
- TavernaFetaGUI:** A window showing "Resulting Operations" for the query "alignment_multiple::tralign". It lists several providers like "alignment_global::needle", "alignment_multiple::tralign", "alignment_local::water", etc.
- Service details:** A detailed view of the "alignment_multiple::tralign" service, showing its description: "Aligns DNA sequences given an alignment of the sequences which they code for." It also lists the description file location, service endpoint location, service interface location, organization name (European Bioinformatics Institute), and service type (Soaplab service).



What Services we Support





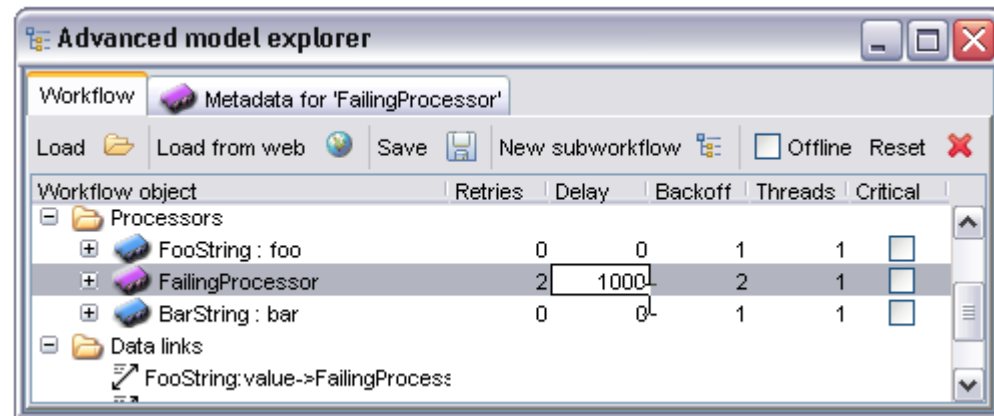
What shall I do when a service fails?

- Most services are owned by other people
- No control over service failure
- Some are research level

Workflows are only as good as the services they connect!

To help - Taverna can:

- Notify failures
- Instigate retries
- Set criticality
- Substitute services
- Instigate checkpoints for long



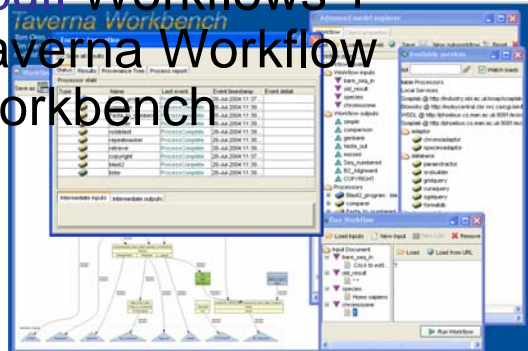


Data Management

- Workflows can generate vast amount of data - how can we manage and track it?
- Data **AND** metadata **AND** experiment provenance
- LSIDs - to identify objects
- Semantic Web technologies (RDF, Ontologies)
 - To store knowledge provenance
- Taverna workflow workbench & plugins
 - Ensure automated recording

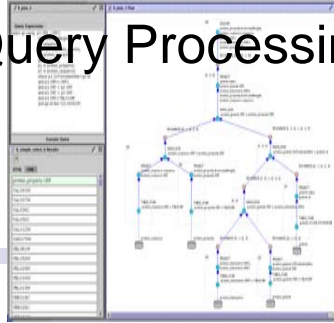


Scuff Workflows + Taverna Workflow Workbench

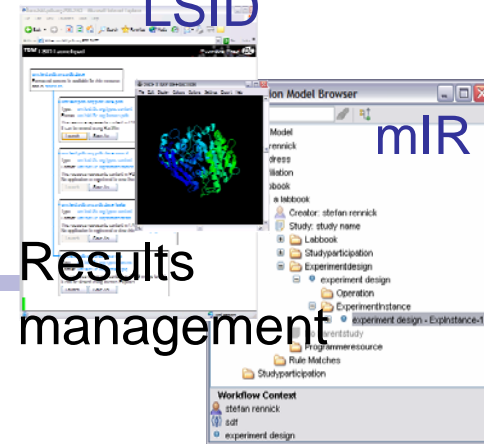


The U of Ma

OGSA-Distributed Query Processing



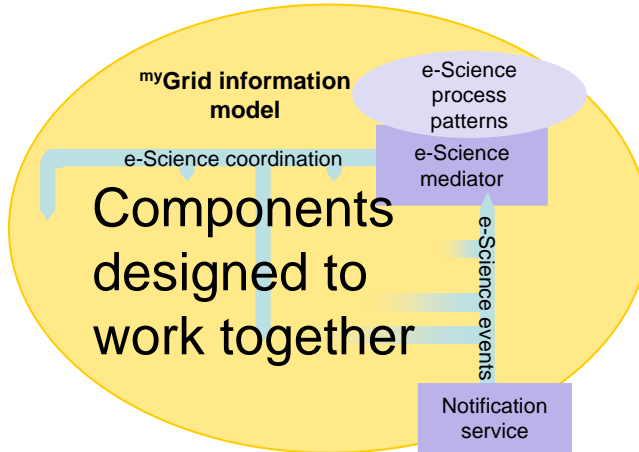
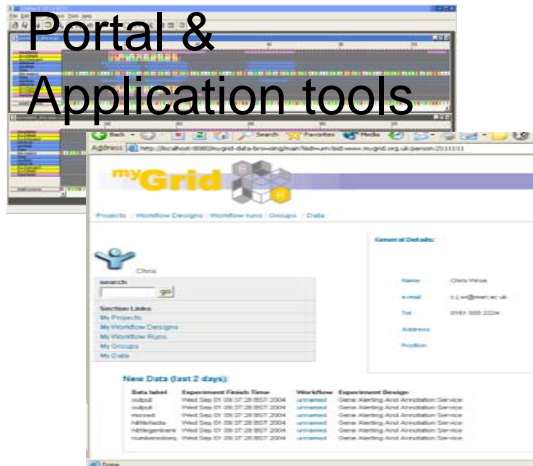
LSID



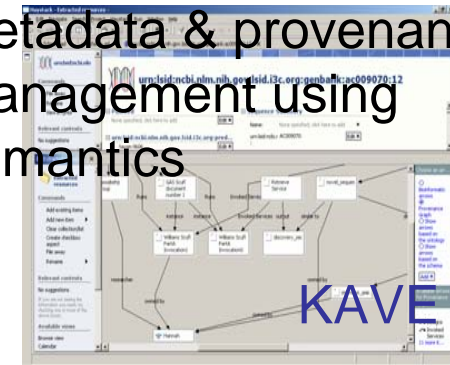
mIR

Results management

Portal & Application tools



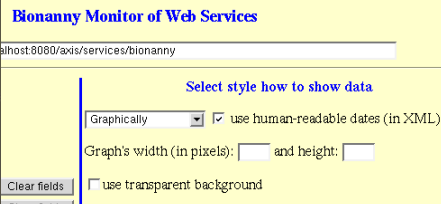
Metadata & provenance management using semantics



KAVE



Service management

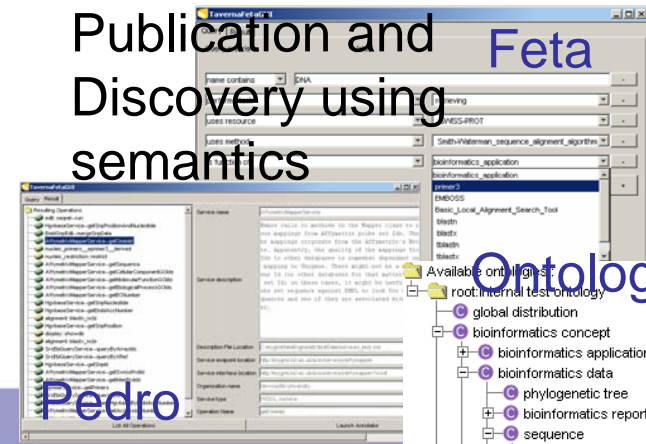


Text Mining Services



Termino

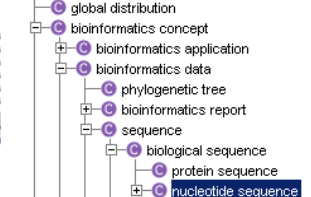
Publication and Discovery using semantics



Feta

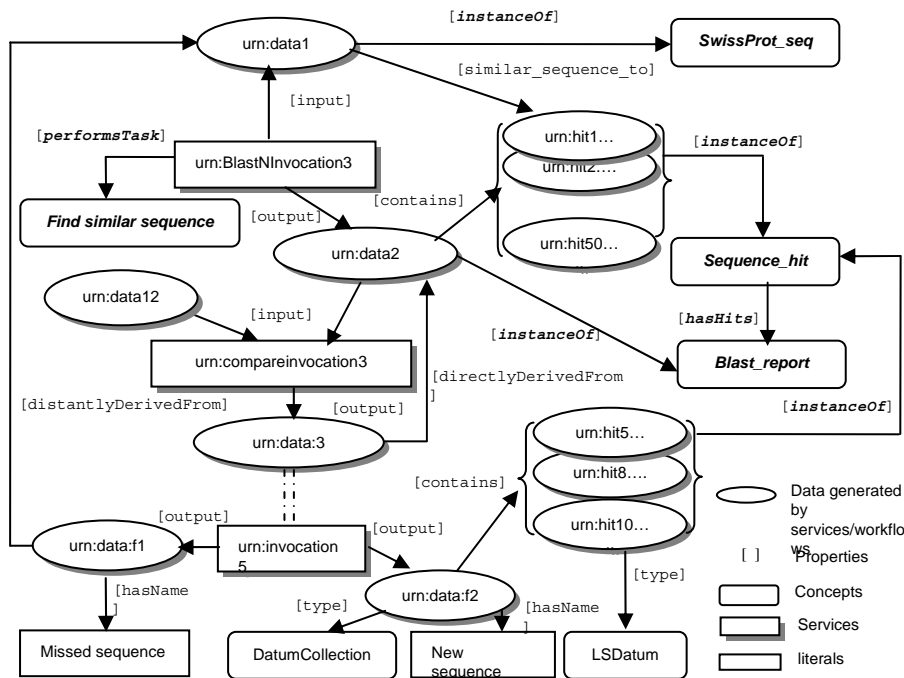
Pedro

Ontology





KAVE Data and metadata management



- Life Science Identifiers
- Information Model
- File management
- Support for custom database building
- Provenance metadata capture using RDF
- SRB integration
- OGSA-DAI integration



Provenance Browsing in Taverna

Provenance Browser (Left Screenshot)

Workflow Instances

Workflow ID	Date	Author
Fetch PDB flatfile from RCSB server	3/10 14:19:34	Tom Oinn
TEY24033SM10	4/10 11:16:22	
Workflow:	4/10 11:16:22	
Fetch PDB flatfile from RCSB server	2/10 17:47:47	
created by:	3/10 11:11:54	
Tom Oinn	2/10 17:38:59	
	3/10 14:15:15	
	4/10 10:56:46	

Description:
Given an identifier such as '1crn'

Status:
fetches the PDB format flatfile from the RCSB

Processor status

Type	Name	Event End Time	Event detail
	AddPrefixToID	4/10 11:16:22	ProcessCompleted
	AddSuffix	4/10 11:16:22	ProcessCompleted
	RCSBPrefix	4/10 11:16:22	ProcessCompleted
	FetchPage	4/10 11:16:22	ProcessCompleted
	RCSBSuffix	4/10 11:16:22	ProcessCompleted

Intermediate outputs

name
text/plain,chemical/x-pdb,text/html
Click to view...
[urn:lsid:www.mygrid.org.uk:8080/Isdocument:YE00LSOZMQ5](http://www.mygrid.org.uk:8080/Isdocument:YE00LSOZMQ5)

Provenance Browser (Right Screenshot)

Workflow Instances

Workflow ID	Date	Author
Compare functions of genes on human Y chromosome to those on X	28/9 14:52:49	Tom Oinn
SXUNSS64083	3/10 12:20:51	
BAE16NL8NE0	3/10 12:30:49	
CXWPBUYMRJ0		
Cluster		Tom Oinn
Example of a conditional execution workflow		Tom Oinn
Williams PartA - version to match ISMB paper		Hannah J. Tipney adapted by ...
CXWPBUYMRJ2	3/10 12:50:41	

Status

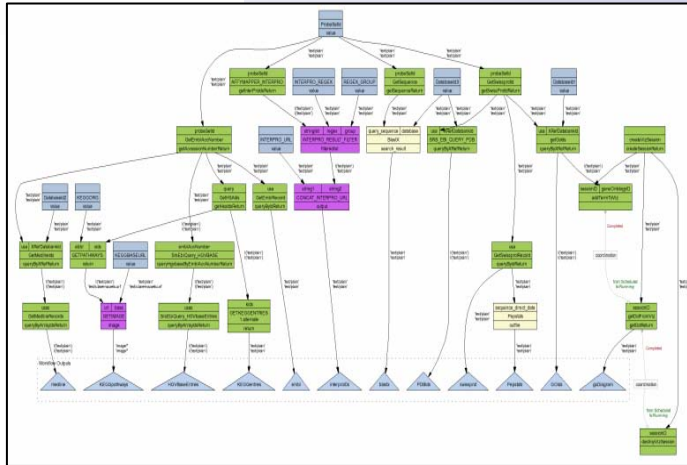
Processor status

Type	Name	Event End Time	Event detail
	showCommon	3/10 12:20:35	ProcessCompleted
	showAnnotatedNonUnique	3/10 12:20:47	ProcessCompleted
	showCommon	3/10 12:20:21	ProcessCompleted
	getDot	3/10 12:20:51	ProcessCompleted
	getCommonAncestors	3/10 12:18:34	ProcessCompleted
	addTerm_collection	3/10 12:18:2	ProcessCompleted
	getCommonAncestors	3/10 12:19:5	ProcessCompleted
	showCommon	3/10 12:20:31	ProcessCompleted

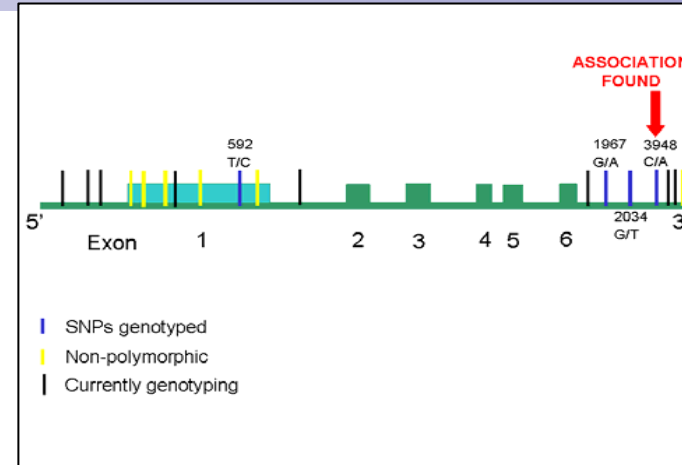
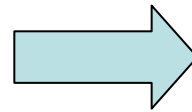
Intermediate outputs

name
text/plain,text/x-graphviz
Click to view...
[urn:lsid:www.mygrid.org.uk:8080/Isdocument:A428DHQKPZ5615](http://www.mygrid.org.uk:8080/Isdocument:A428DHQKPZ5615)

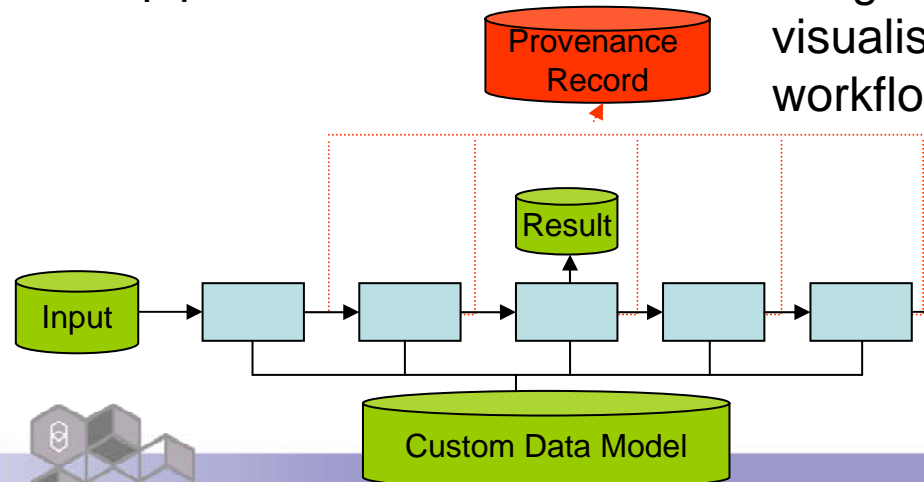
Results Integration



Gene annotation pipeline workflow



Integration and visualisation of GD annotation workflow results



Smarter workflow design
incorporating visualisation
VBI collaboration



Visualisation

The screenshot shows the SeqVista interface for the protein 1FAS (TOXIN). At the top, the amino acid sequence is displayed: T M C Y S H T T T S R A I L T N C G E N S C Y R K S R R H P P K M V L G R G C G C P P G D D Y L. Below the sequence, there are tracks for Secondary Structure, Topohydrophobic Residues, Most Interacting Residues, and Tight End Fragments. The main window shows a 3D ball-and-stick model of the protein structure. On the right, there is a workspace for annotations and options.

SeqVista

Utopia

The screenshot shows the Utopia software interface. The top menu includes File, Visualization, Analysis, Comparison, Advanced, and Help. The location is set to 's:\Hu\Java\Gene\SeqVISTA\SampleInput\Alad_RM.gv'. The main window displays a gene model for the H.sapiens ALAD gene for porphobilinogen synthase. A feature labeled 'Protein_Bind' is highlighted with a blue box and an orange oval. Below the gene model, a sequence alignment is shown with the following text:

H.sapiens ALAD gene for porphobilinogen synthase.

2051	agcctcaaga	tcctcttgc	cctgcacatc	tccaatctcc	ataaagacct
2101	ttgatcggat	ctatcattgt	acctatcata	ggtctgatgc	ccctatcaag
2151	acttggagtt	ttcctaaacg	cccattgtct	tcaatcaca	tctctcaact
2201	catagcattg	cgtaccctg	agaaataaat	gaagctgga	caaattttct

The 'acctatcata' sequence in the second row is circled in orange, with arrows pointing to the 'Protein_Bind' feature in the gene model above.

Applications

Resistance to trypanosomiasis in cattle in Kenya

Andy Brass, Paul Fisher – University of Manchester

Microarray

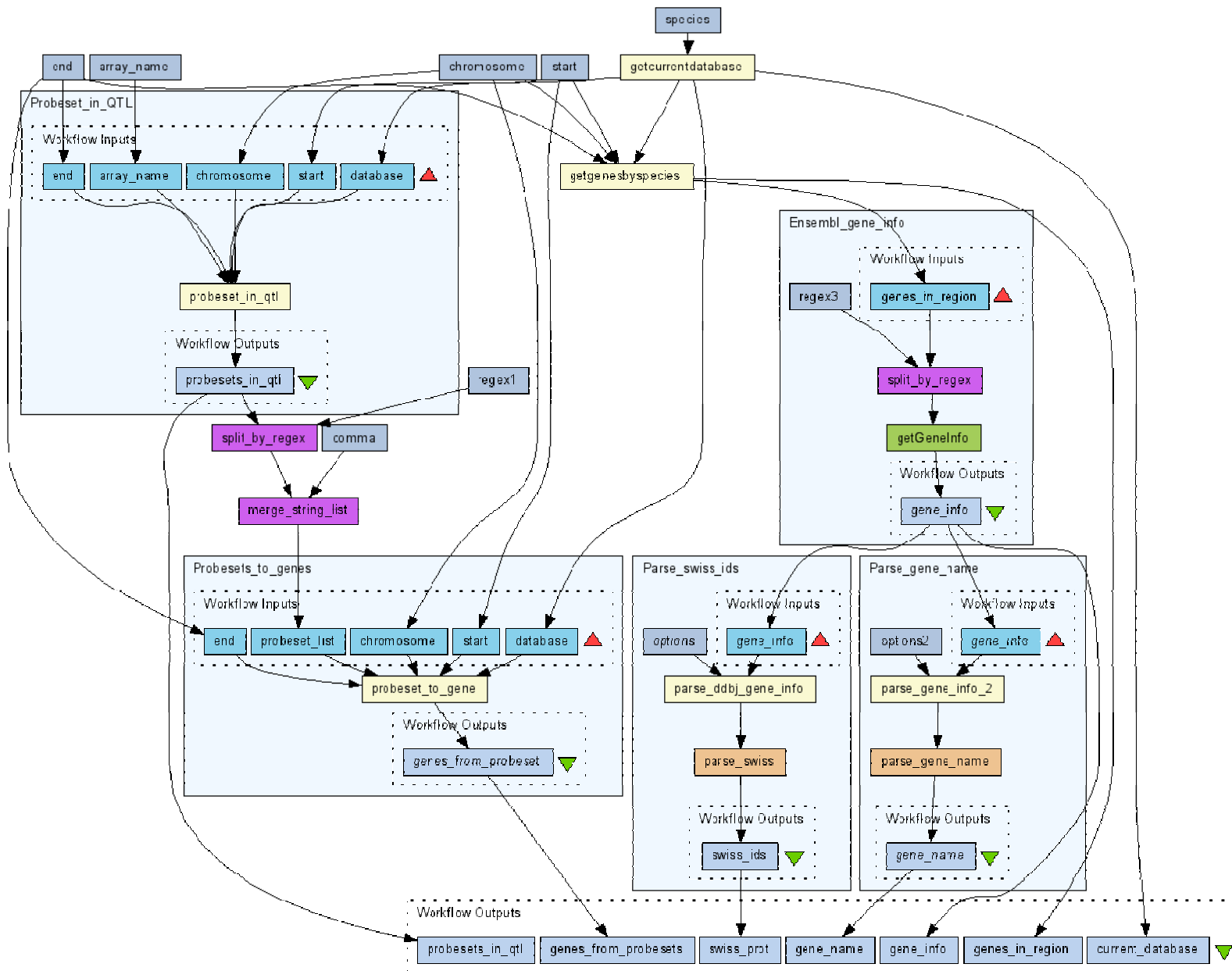
QTL

SNPs

Metabolic pathway analysis



Need to access microarray data, genomic sequence information, pathway databases AND integrate the results





myGrid Alliance: Applications

Large user community – over 15600 downloads

PsyGrid



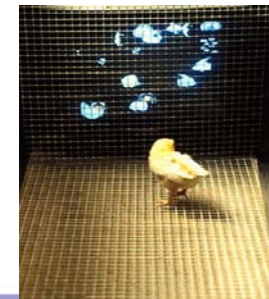
Small molecules,
Murray-Rust, Cambridge



UK - NZ Partnership

Mias-Grid

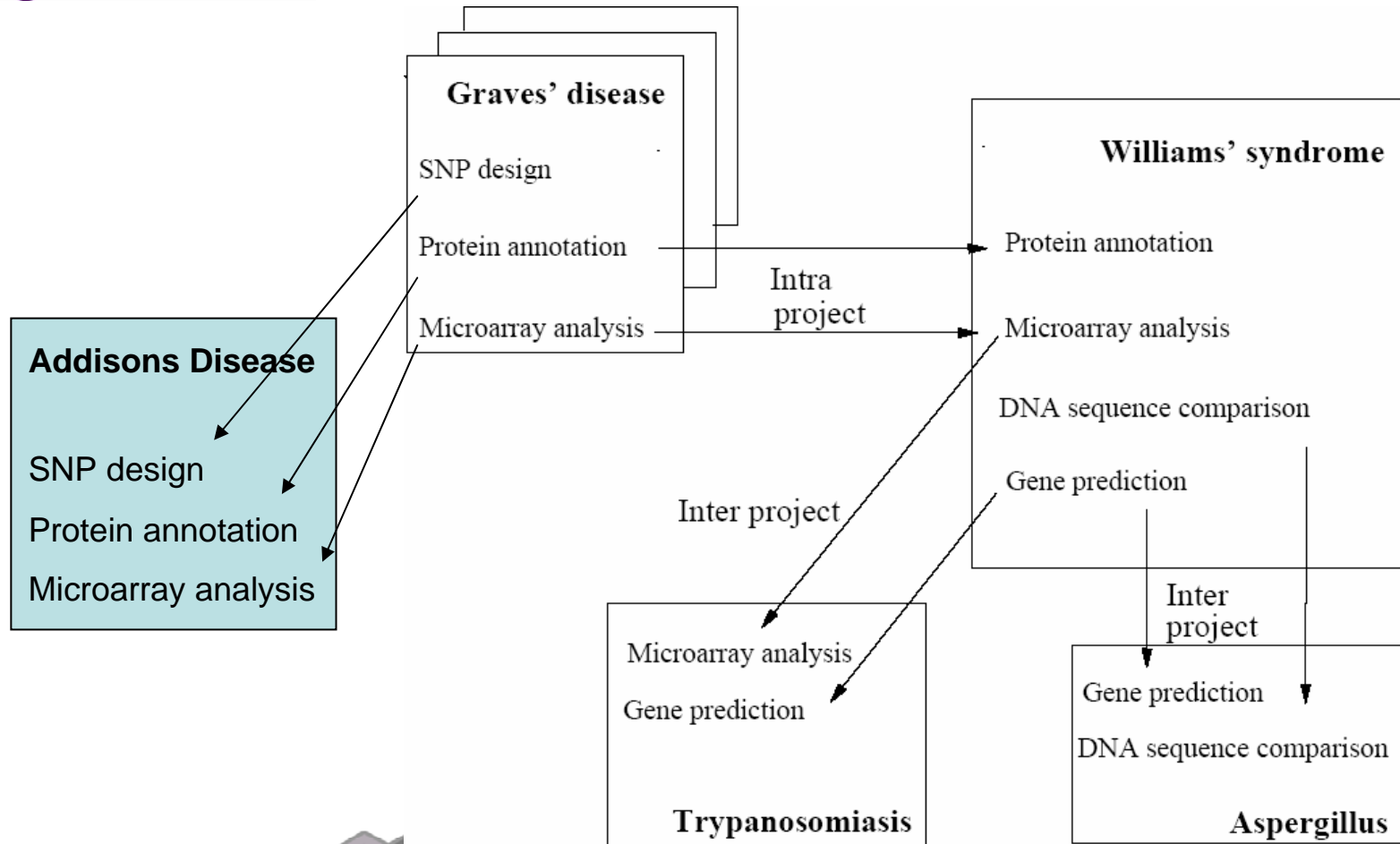
Chicken genome
Roslin Institute



myGrid



Workflow Reuse





Taverna is now OMII-UK

- Taverna 1.3.1 production Sept 2006
 - Packaging, Installation, Deployment, Maintenance, Testing
 - GridSAM, GRIMOIRES, BioMOBY integration
 - Semantic content for registry
 - Smoothed integration of discovery and metadata management
 - Security AA for KAVE data and metadata management
- Taverna 2.0 Spring 2007
 - Redevelopment of the plug in and enactor framework, improved iteration events, data management
- Close collaboration with pioneers
- Incremental rollouts to early adopters





Taverna in OMII-UK

- Development of Taverna 2.0
 - reworking of the processor model to include dual execution semantics incorporating data and control flow
 - enhanced support for long-running workflows
 - large scale data transfer
 - improved provenance collection with nested workflows and complex iterations
 - fully distributed workflow enactment and authoring





Acknowledgements

- Carole Goble and the myGrid team
- OMII-UK
- All of our users

