# $^{my}$Grid : Workflow based *in silico* experiments in biology

**Katy Wolstencroft, Tom Oinn, Phillip Lord, Jun Zhao Carole Goble and the $^{my}$Grid team & May Tassabehji / Hannah Tipney Medical Genetics, St Mary's Hospital Manchester**

# What is ᵐʸGrid?

- e-Science pilot research project funded by EPSRC http://www.mygrid.org.uk

- Manchester, Newcastle, Sheffield, Southampton, Nottingham, EBI and RFCGR, also industrial partners.

- 'targeted to develop open source software to support personalised *in silico* experiments in biology on a grid.'

- Now - platform grant for ᵐʸGrid 2 and ᵐʸGrid is an OMII-UK (Open Middleware Infrastructure Institute) node

# What is <sup>my</sup>Grid?

"A comprehensive loosely-coupled suite of middleware components specifically to support data intensive *in silico* experiments in biology"

- Distributed computing
- Workflow Design and Enactment
- Provenance and Data management
- Semantic Discovery

# Collaborators



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

# Motivation

Bioinformatics is an open Community

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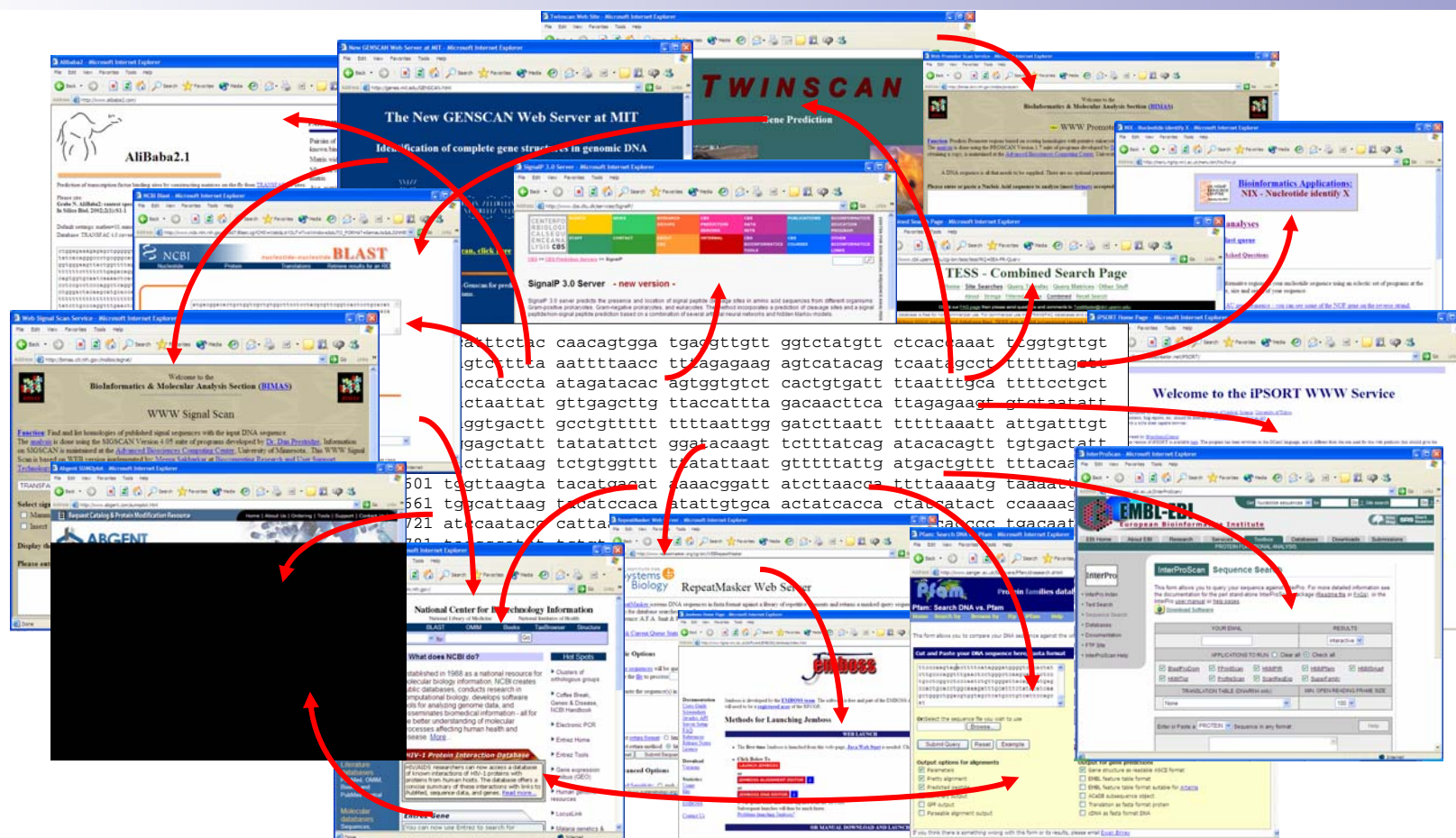- Open access to tools
- Open access to applications

Global Bioinformatics

# Problems

- Heterogeneous data
- Distributed resources
- Potentially requires supercomputing power
- Very few standards – I/O formats, data representation, annotation

Integration and interoperability between resources is difficult

# Traditional Approach

# Cutting and Pasting

- Advantages:
  - Low technology on both server and client side
  - Very robust: Hard to break
  - Data integration happens along the way
- Disadvantages:
  - Time consuming (and painful!)
    - Can be repeated rarely
    - Limited to small data sets
  - Error prone:
    - Poor repeatability

# Pipeline Programming

- Advantages
  - Repeatable
  - Allows automation
  - Quick, reliable, efficient

- Disadvantages
  - Requires programming skills
  - Difficult to modify
  - Requires local tool and database installation
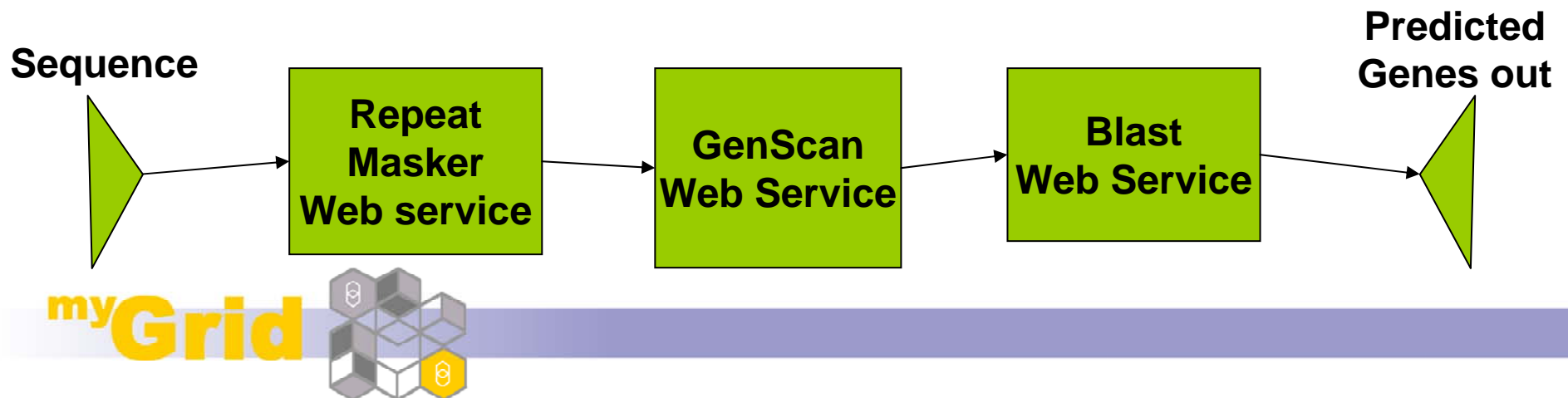  - Requires tool and database maintenance!!!

# <sup>my</sup>Grid Requirements

- Automation
- Reliability
- Repeatability
- Distributed resources
- Few programming skill required

# <sup>my</sup>Grid 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
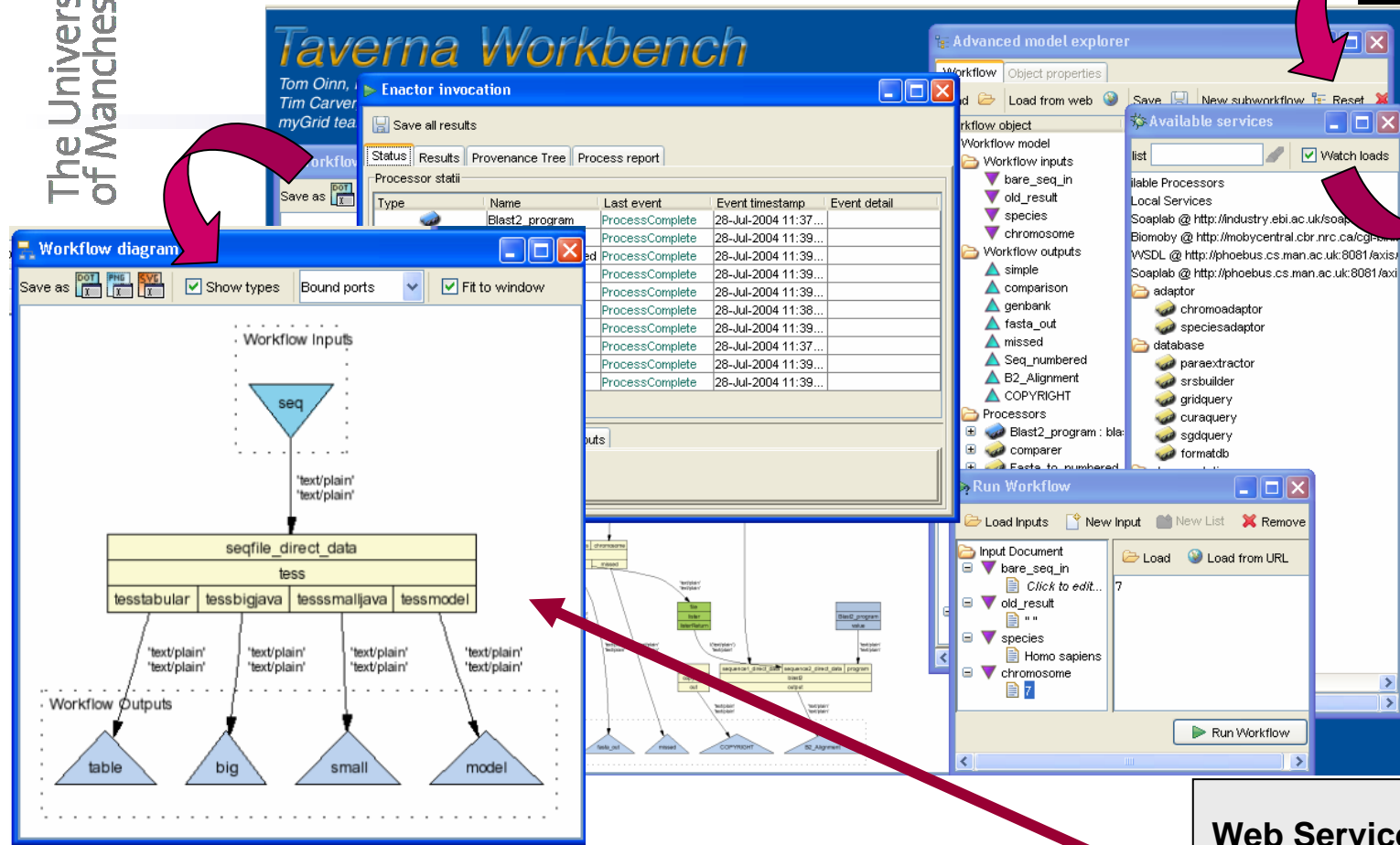   – therefore, you don't have to be a coder to build workflows

# Workflow Diagram

# Workflow Advantages

- Automation
  - Capturing processes in an explicit manner
  - Tedium! Computers don't get bored/distracted/hungry/impatient!
  - Saves repeated time and effort

- Modification, maintenance, substitution and personalisation

- Easy to share, explain, relocate, reuse and build

- Releases scientists/bioinformaticians to do other work

- Record
  - Provenance: what the data is like, where it came from, its quality
  - Management of data (LSID - Life Science Identifiers)

# Workflow Components



**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 |

# Workflow Services – Web Services

- automated programmatic internet access to applications
- Many bioinformatics resources provide web service versions of popular tools - e.g. NCBI BLAST

^myGrid > 1000 web services for bioinformatics

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

# Web Service Creation

Many service providers provide web service implementations of applications

## SoapLab

- For wrapping 'legacy' command-line applications for use in Taverna

## GowLab

- For wrapping html resources

potentially ANY bioinformatics tool / resource can be wrapped and used in myGrid workflows

# Use Cases:
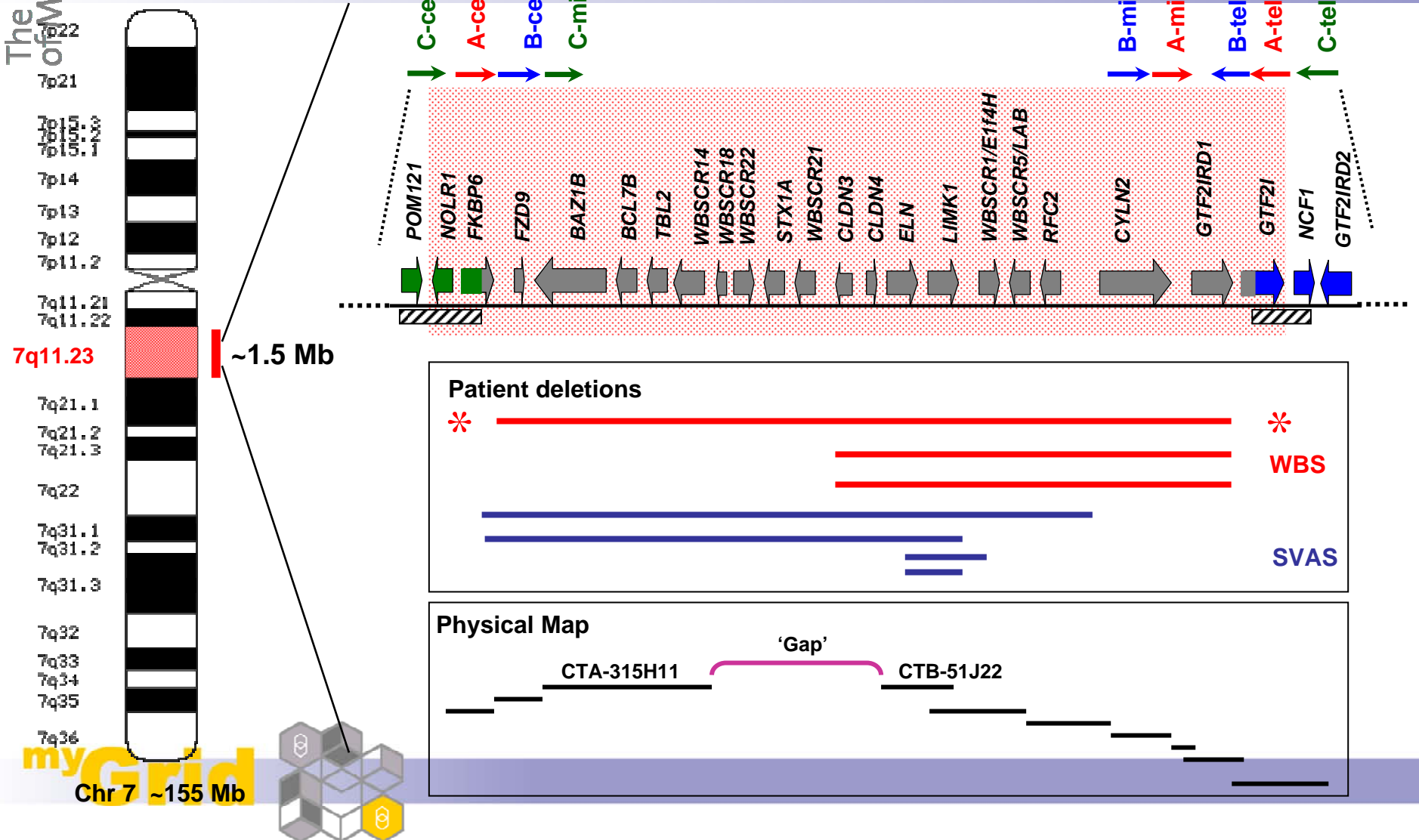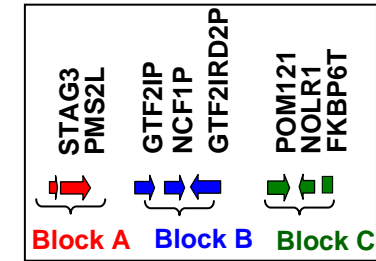# Williams-Beuren Syndrome (WBS)



- Contiguous sporadic gene deletion disorder
- 1/20,000 live births, caused by unequal crossover (homologous recombination) during meiosis
- Haploinsufficiency of the region results in the phenotype
- Multisystem phenotype – muscular, nervous, circulatory systems
- Characteristic facial features
- Unique cognitive profile
- Mental retardation (IQ 40-100, mean~60, 'normal' mean ~ 100 )
- Outgoing personality, friendly nature, 'charming'

# Williams-Beuren Syndrome Microdeletion

*Eicher E, Clark R & She, X  An Assessment of the Sequence Gaps: Unfinished Business in a Finished Human Genome.  Nature Genetics Reviews (2004) 5:345-354*
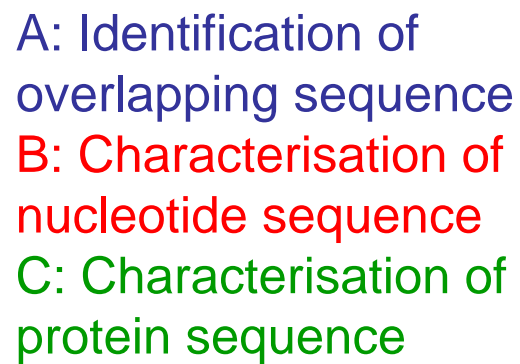*Hillier L et al. The DNA Sequence of Human Chromosome 7.  Nature (2003) 424:157-164*

MANCHESTER 1824

The University of Manchester

myGrid



STAG3 PMS2L | GTF2IP NCF1P GTF2IRD2P | POM121 NOLR1 FKBP6T

Block A    Block B    Block C

7q11.23    ~1.5 Mb

Chr 7  ~155 Mb

C-cen  A-cen  B-cen  C-mid

B-mid  A-mid  B-tel  A-tel  C-tel

POM121  NOLR1  FKBP6  FZD9  BAZ1B  BCL7B  TBL2  WBSCR14  WBSCR18  WBSCR22  STX1A  WBSCR21  CLDN3  CLDN4  ELN  LIMK1  WBSCR1/E1f4H  WBSCR5/LAB  RFC2  CYLN2  GTF2IRD1  GTF2I  NCF1  GTF2IRD2

## Patient deletions

\*    \*

WBS

SVAS

## Physical Map

'Gap'

CTA-315H11    CTB-51J22

# Filling a genomic gap *in silico*

- Identify new, overlapping sequence of interest

- Characterise the new sequence at nucleotide and amino acid level

- Frequently repeated – info rapidly added to public databases
- Time consuming and mundane
- Don't always get results
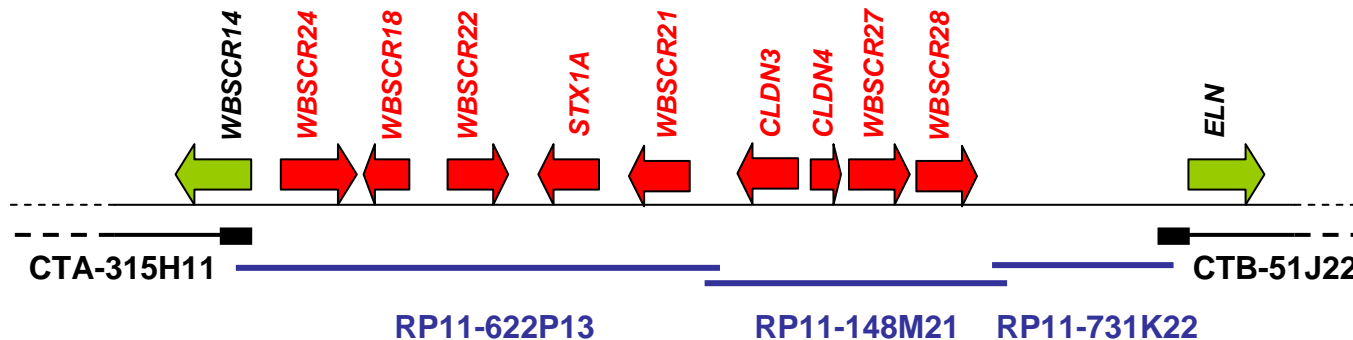- Huge amount of interrelated data is produced

**The Williams Workflows**

A

B

C

A: Identification of overlapping sequence
B: Characterisation of nucleotide sequence
C: Characterisation of protein sequence

# The Biological Results

**Four workflow cycles totalling ~ 10 hours**
**The gap was correctly closed and all known features identified**



314,004bp extension

All nine known genes identified
(40/45 exons identified)

# The Workflow Experience

Have workflows delivered on their promise?  **YES!**

- Correct and biologically meaningful results

- Automation
  - Saved time, increased productivity
  - Process split into three, you still require humans!

- Sharing
  - Other people have used and want to develop the workflows

- Change of work practises
  - *Post hoc* analysis. Don't analyse data piece by piece receive all data all at once
  - Data stored and collected in a more standardised manner
  - Results amplification

# Workflow Reuse



**Mouse genome**



## Trypanosomiasis in cattle



**Chicken genome**

# Web Service Issues

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

## However - Taverna can

- Notify users of service failures
- Instigate retries
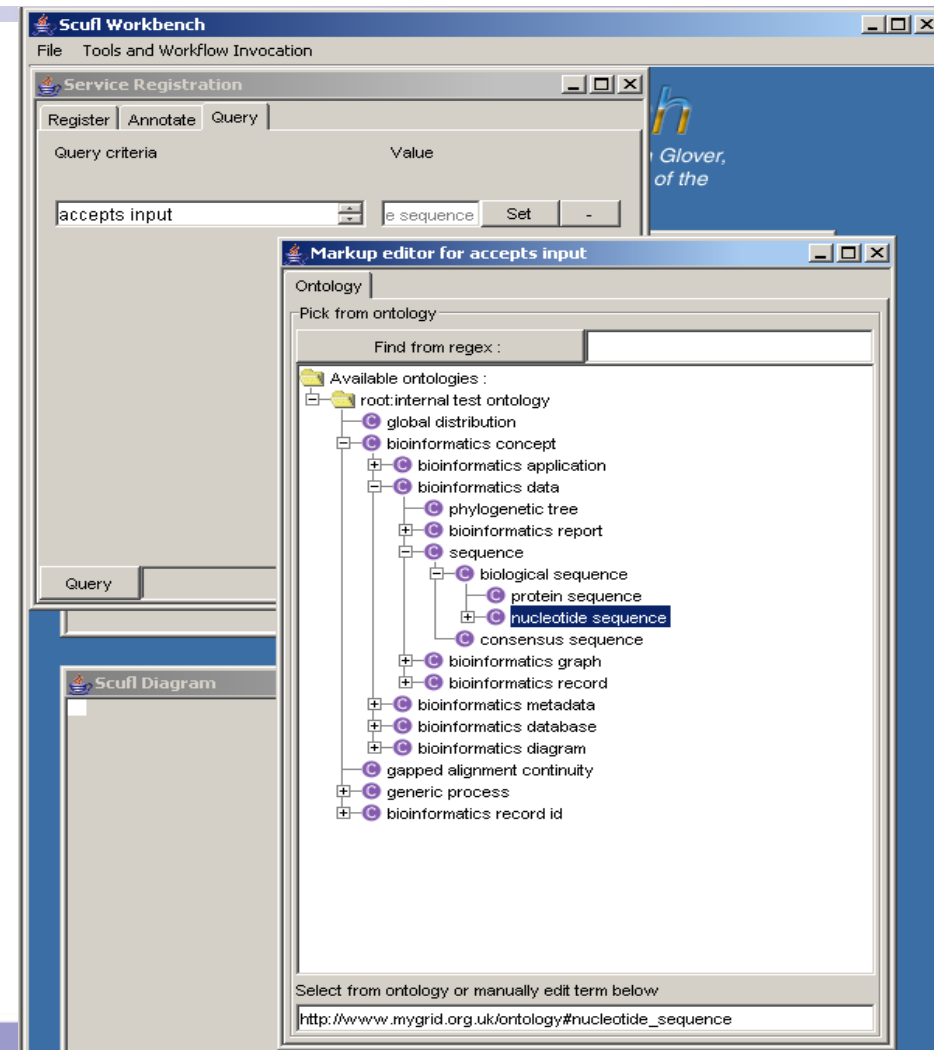- Substitute services from other sites

# Fault Tolerance

# Service Discovery – Feta

Services – only useful if the user can find them!

Feta semantic discovery – can find services by:

- **Bioinformatics task**
- **Resource used**
- **Bioinformatics method**
- **Input format**
- **Output format**

# 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

# Provenance

- Stored in RDF–Resource Description Framework
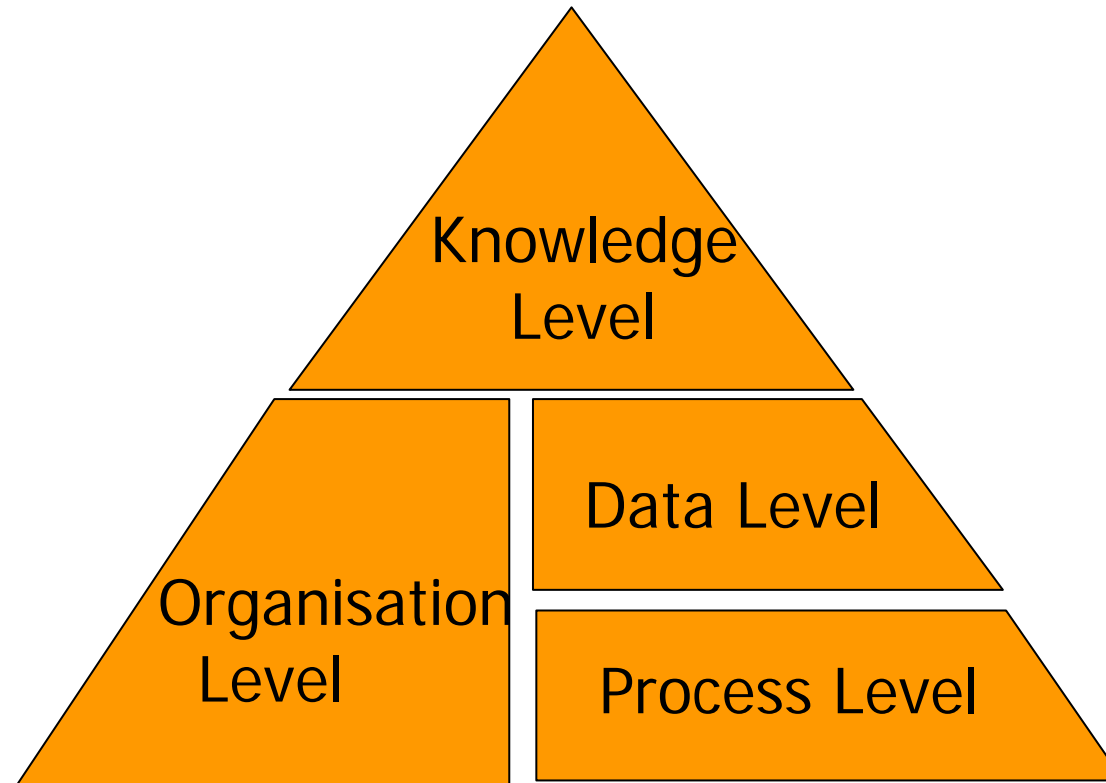
- < Subject, Predicate, Object>

e.g. *<urn:actualInputParameter0, isAbout: my:provenance>*

Identified by a URI and typed by an ontology

- Inputs and outputs and ᵐʸGrid services can be defined using ᵐʸGrid Service Ontology – describes bioinformatics processes

# Life Science Identifiers

- Life Science Identifiers (LSIDs) are the standard adopted by the Object Management Group (OMG) for the identification of life science data objects

  urn:lsid:ncbi.nlm.nlh.gov.lsid.biopathways.org:genbank_gi:7717376

- LSIDs used throughout $^{my}$Grid to ID data objects from external sources as well as internally created data.

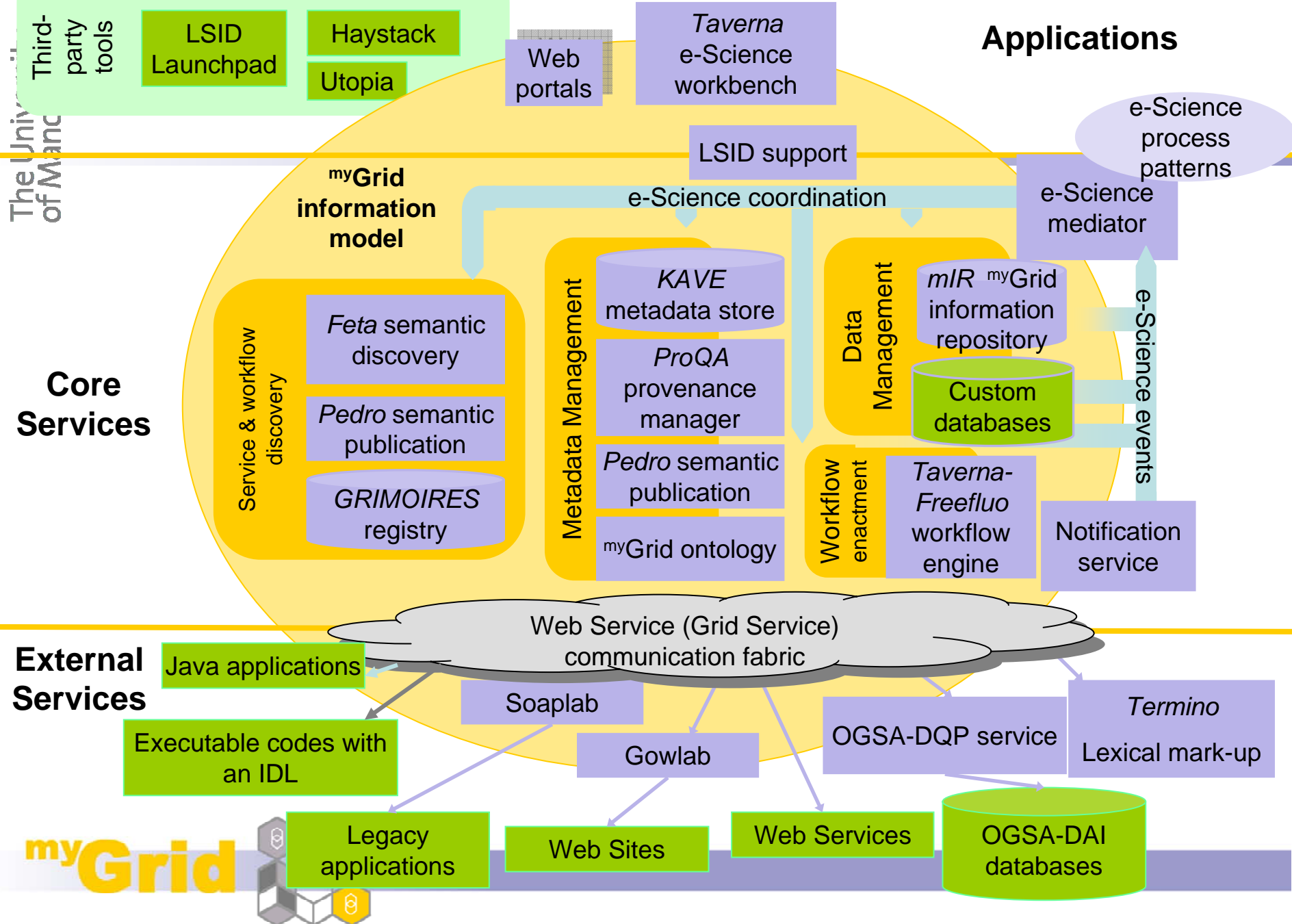- Using a standard mechanism for identification allows for more efficient and cohesive exchanges between $^{my}$Grid components

# <sup>my</sup>Grid Architecture

In keeping with the bioinformatics community

- Open architecture
  - Service Oriented Architecture
  - Loosely coupled
  - Web services based
  - Assemble your own components
  - Designed to work together

# <sup>my</sup>Grid Users

## Widespread uptake

Bioinformatics

Systems Biology

Chemistry

Medical Physics

Many new e-Science projects using <sup>my</sup>Grid platform
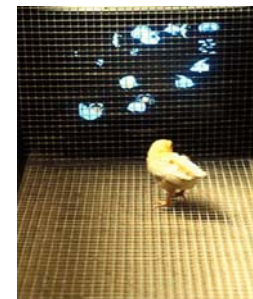
# <sup>my</sup>Grid Alliance: Application

PsyGrid

Small molecules,
Murray-Rust, Cambridge

Chicken genome
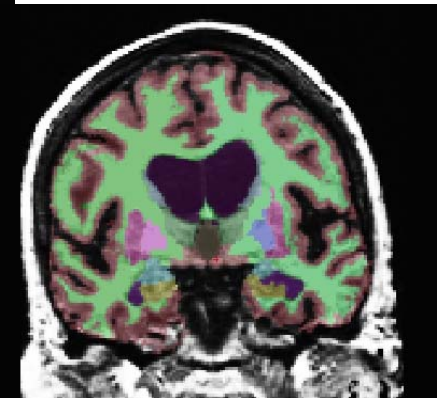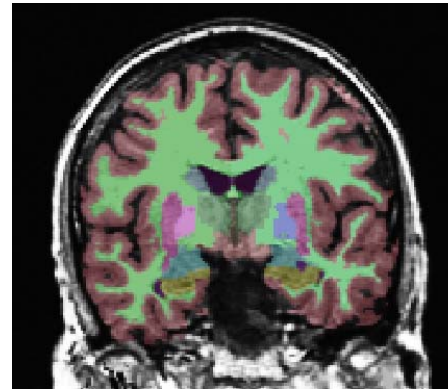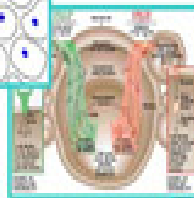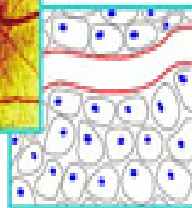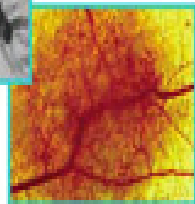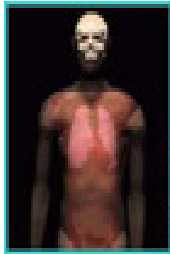Roslin Institute

# New projects - New directions for <sup>my</sup>Grid

Bioinformatics data – mostly sequence based, workflow processes take several hours

Other Fields
- More data / larger data
- Long running workflows – weeks/months

# New Directions



Long Running Workflows

Large Data Sets

MIAS-Grid

PsyGrid

Computational steerage of
heart simulation codes

# Conclusions

- $^{my}$Grid enables interoperability and integration of bioinformatics resources
- Legacy applications and new tools can be made available in the workbench
- User driven approach to development
- $^{my}$Grid offers a whole suite of components to design and enact workflows and trace experiments

# myGrid and WBS People

**Core**

Matthew Addis, Nedim Alpdemir, Pinar Alper, Tim Carver, Rich Cawley, Neil Davis, Alvaro Fernandes, Justin Ferris, Robert Gaizaukaus, Kevin Glover, Carole Goble, Chris Greenhalgh, Mark Greenwood, Yikun Guo, Ananth Krishna, Peter Li, Phillip Lord, Darren Marvin, Simon Miles, Luc Moreau, Arijit Mukherjee, Tom Oinn, Juri Papay, Savas Parastatidis, Norman Paton, Terry Payne, Matthew Pockock Milena Radenkovic, Stefan Rennick-Egglestone, Peter Rice, Martin Senger, Nick Sharman, Robert Stevens, Victor Tan, Daniele Turi, Anil Wipat, Paul Watson, Katy Wolstencroft and Chris Wroe.

**Users**

Simon Pearce and Claire Jennings, Institute of Human Genetics School of Clinical Medical Sciences, University of Newcastle, UK

Hannah Tipney, May Tassabehji, Andy Brass, St Mary's Hospital, Manchester, UK

**Postgraduates**

Martin Szomszor, Duncan Hull, Jun Zhao, John Dickman, Keith Flanagan, Antoon Goderis, Tracy Craddock, Alastair Hampshire

**Industrial**

Dennis Quan, Sean Martin, Michael Niemi, Syd Chapman (IBM)

Robin McEntire (GSK)