# Introduction

Dr Katy Wolstencroft

# Introduction to the Tutorial

- Using scientific workflows for bioinformatics data analysis, integration and scaling
    - Galaxy – local tools workflows
    - Taverna – distributed computing workflows
    - WS-PGRADE - grid workflows
- RNA-Seq analysis pipeline
- 3 hands-on session
- Discussion
- Workflows surgery session

# Introduction to Scientific Workflows

# Typical Bioinformatics Experiments

1. Collect data / find data in repositories
2. Run analysis algorithm on data
3. Run different analysis algorithm on results of (2) or get more/different data and run (2) again
4. Repeat for any number of sequential algorithms
5. Analyse results
6. Search for supporting information to help explain results and generate new hypotheses
7. Start cycle again with new hypotheses

# The Practical Reality

- Sequential use of distributed tools
- Incompatible input and output formats
- Analysis of large data sets by multiple researchers
- Difficult to record parameter selections
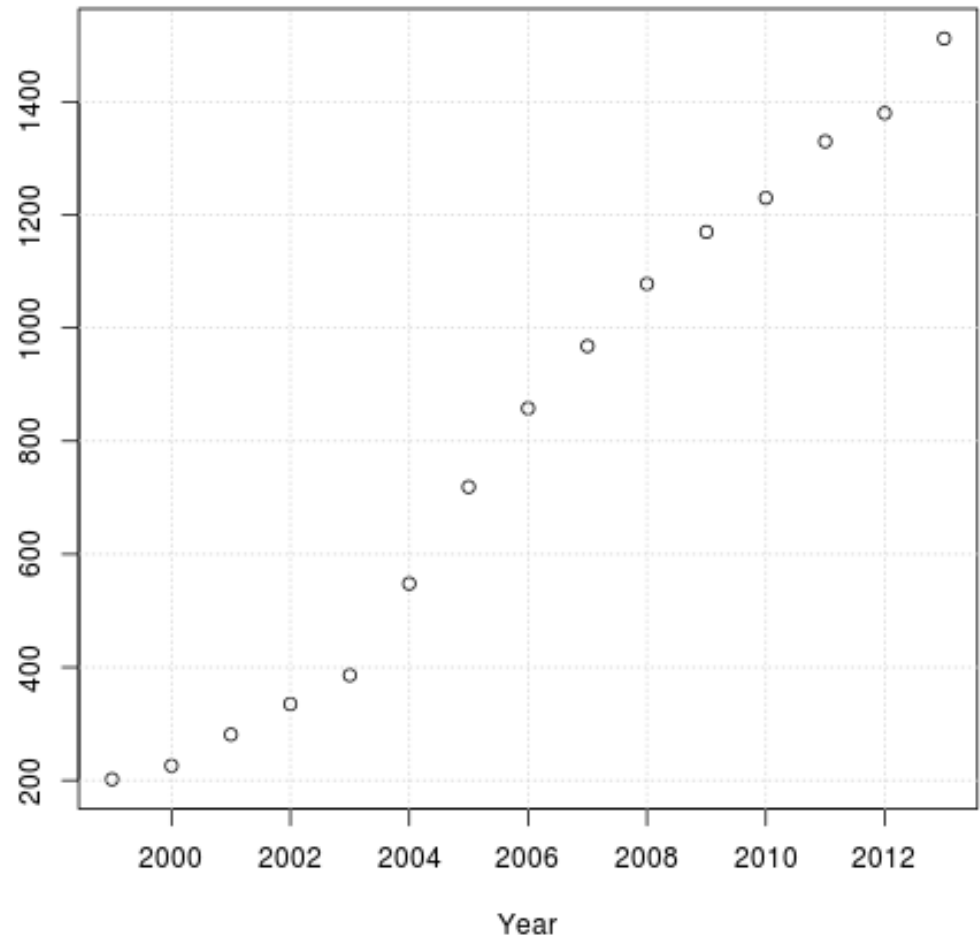- Difficult to reproduce analyses

# Lots of Resources

**NAR 2014 – 1552 databases**


Growth of Biological Databases

# Sources of Data

- In repositories run by major service providers (e.g. NCBI, EBI, DDBJ)

- Group/Institute web sites

- On ftp servers

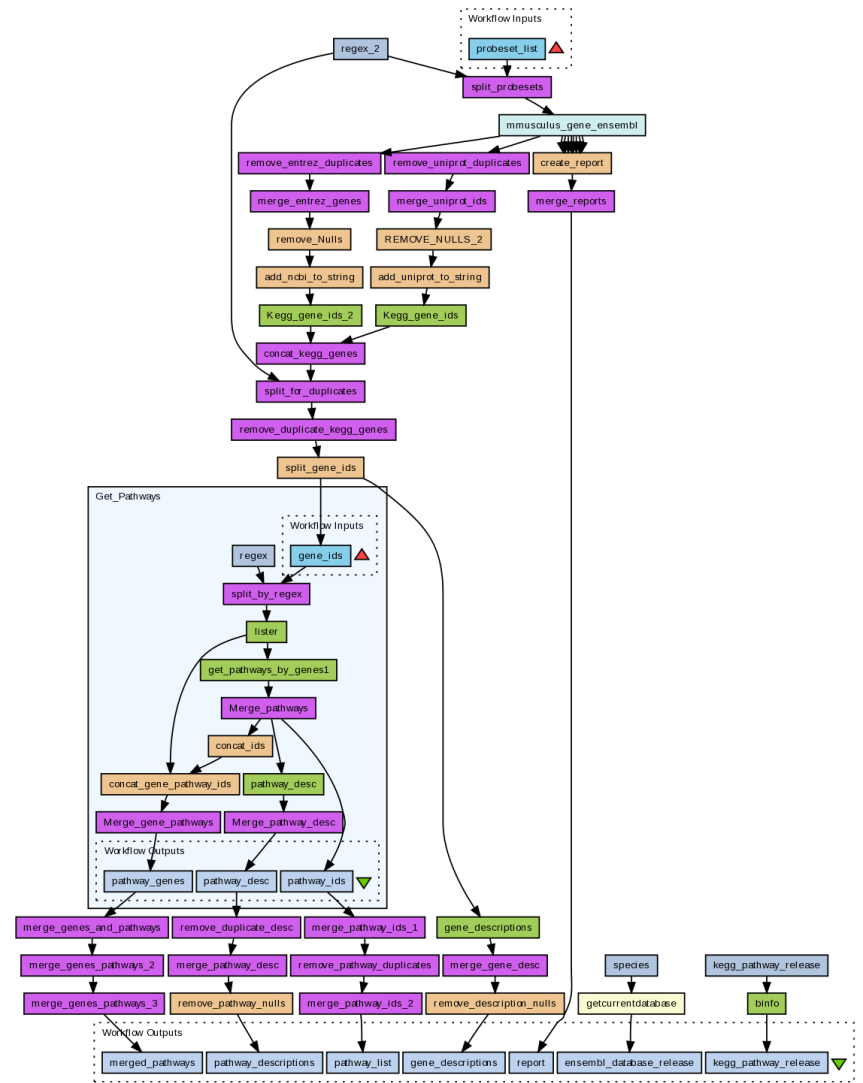- In local project stores

- Few defined formats

- Inconsistent metadata

# Workflow as a Solution

- Automating the process
- Sophisticated analysis pipelines
- A set of **services** to analyse or manage data (either local or remote)
- Data flow through services
- Control of service invocation
- Iteration
- Encapsulates experimental method

# Promoting Reproducible Research

Informatics involves

- Complex, multi-step analyses

- Lots of data as inputs

- Lots of data generated

- Workflows encapsulate the methods used and their parameters and the input data

  - Easier to repeat

- Workflows allow you to visualise the methods

  - Easier to assess

**Misconduct in science**

# An array of errors

Investigations into a case of alleged scientific misconduct have revealed numerous holes in the oversight of science and scientific publishing

Sep 10th 2011 | From the print edition

Timekeeper   Like 977   Tweet 219

http://www.economist.com/node/21528593



Nature Oncology Reviews

# Preventing Irreproducible Research

- Duke University, 2006 -Prediction of the course of a patient's lung cancer using expression arrays

- *Nature Medicine* - **12**, 1294 - 1300 (2006)

- Recommendations on different chemotherapies from cell cultures

- Clinical trials of personalised treatments

- Major advances in personalised medicine?

# Preventing Irreproducible Research

- 3 different groups **could not reproduce** the results and **uncovered mistakes** in the original work

- Requested raw data and more information about all the methods

- Re-created whole analysis methods

- *Nature Medicine* - **12**, 1294 - 1300 (2006) Published online: 22 October 2006; Corrected online: 27 October 2006; Corrected online: 21 July 2008; Retracted: 07 January 2011 | doi: 10.1038/nm1491
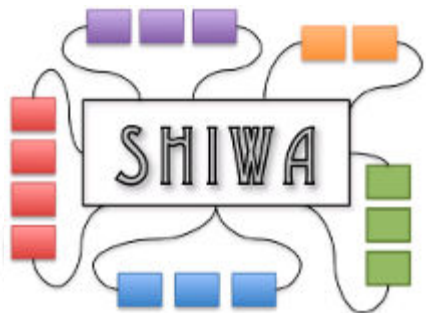
# If the Analyses were done using Workflows.....

- Reviewers could re-run experiments and see results for themselves

- Methods could be properly examined and criticised

- Mistakes could be pinpointed
  - At what point in the analysis pipeline did it start to go wrong
  - What were the inputs and outputs for each step

# 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**

Published workflows

# Workflows are ideal for...

- High throughput analysis
  - Transcriptomics, proteomics, next gen sequencing
- Data integration, data interoperation
- Data management
  - Model construction
  - Data format manipulation
  - Database population
  - Semantic integration
  - Visualisation

# SCIENTIFIC WORKFLOW MANAGEMENT SYSTEMS

# Different Workflow Systems

**VisTrails**

**Kepler**

**Triana**

**Ptolemy II**

**Taverna**

**BPEL**

**Pipeline Pilot**

**Galaxy**

myGrid

# Example Workflow Systems



**https://usegalaxy.org/**

Penn State and Emory University James Taylor and Anton Nekrutenko

UK Consortium, led by Carole Goble,
University of Manchester



**http://www.taverna.org.uk/**



**http://guse.hu/about/architecture/ws-pgrade**

Supported by the ER-Flow project

# Different Types of Workflows

- Two main types of workflows:

  - ## Data flow workflows
    A task is invoked once its **expected data** has been received. When complete, it passes any resulting data downstream

  - ## Control flow workflows
    A task is invoked once its **dependant tasks** have been completed

  - ## Combination
    A task is invoked once its **expected data** has been received AND **dependant tasks** have been completed

# Possible Data Flow Structures

**Sequence**
*Store intermediate results*

**Parallel**
*Apply multiple components to a set of data*

**Choice**
*Decisions at runtime*

**Iteration**
*Loop through datasets*

myGrid

# Workflow engine features

- Implicit iterations
- Parallelisation
  - Run as soon as data is available
- Streaming
  - Process partial iteration results early
- Retries, looping
  - For stability and conditional testing

# Workflows Supporting *in silico* Science

# SERVICES IN WORKFLOWS

# Services in workflows

- Web Services
  - WSDL
  - REST
- Local services
- Grid Services
- Cloud Services
- Workflows
- Secure Services

# Different Approaches to Service Connections

- Closed – connect to services installed locally to the server,designed specifically to work together
  - Less heterogeneity, but fewer services
  - Harder to add new services
  - Galaxy server, Knime
- Open – connect to ANY service regardless of type and structure, hosted anywhere in the world
  - More services, but more heterogeneity
  - Easy to add new services
  - Taverna, Kepler

myGrid

# Different Approaches to Service Connections

- On-the-fly service execution
  - Data and legacy code (i.e. Command-line application) are submitted to a workflow engine for execution
  - A language or interactive mechanism to describe the code that needs to be executed, its dependencies, and the arguments. The system in some cases also takes
  - Transports the data to the computing resource where the job is executed

# Using and Making Galaxy Services

- Tools installed on the Galaxy server (admin)
- Create Galaxy tool definition file
- Create an entry in the Galaxy tool registry tool_conf.xml
- Finite, but compatible collection of resources
- Focus on genomics and NGS, but expanding
- Includes built-in genome browsers

# Using Services with Taverna

- Open domain
- Web Services
  - WSDL or REST - 8000 +
- Local services
- Grid Services
- Cloud Services
- Secure Services
- Workflows
- Specialised services
  - BioMart , R

No installation required!

- Third party – we don't own them – we didn't build them
- All the major providers
  - NCBI, DDBJ, EBI …
- Enforce NO common data model

**my Grid**

# Understanding how services work

# Using and Making WS-Pgrade Services

- On-the-fly service composition

- Execution of web services

- Execution of workflows

# Data and Provenance

# Data and Provenance

- Workflows can generate vast amount of data - how can we manage and track it?

- We need to manage data AND metadata AND experimental provenance

  - Input data

  - Parameter selection

  - Versions of the analysis tools used

  - Intermediate results

  - Intermediate parameter selection

  - Final results

  - What happened, what went wrong, how long it took

# Data and Provenance

Scientists need to:

- Check back over past results, compare workflow runs and share workflow runs with colleagues
- Look at intermediate results when designing and debugging

# Histories



- What you did
- With what data
- In what order

- Automatically converts histories to workflows

# Where Workflows are Executed

# Where workflows are executed

- Local Execution
  - Client supports complex workflow design
  - Download and install client
  - Easy to access and use local data and tools
  - Easy to store results

- Server execution
  - Workflows executed through web interface or other client
  - Workflows run on server, grid or cloud
  - Better for larger, long running workflows
  - Better for scaling-up
  - Serves finished workflows to users

myGrid

# Spectrum of Users

Intermediate users reuse and modify existing workflows



Advanced users design and build workflows (informaticians)



http://www.myexperiment.org

Load Data:

Run Workflow

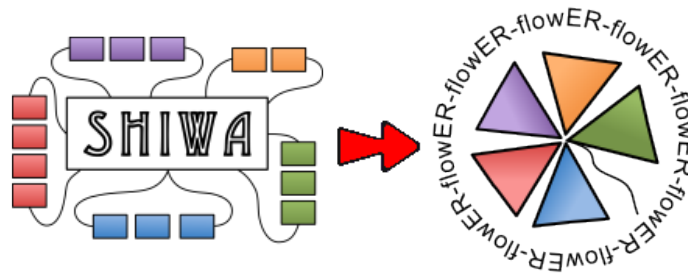Others "replay" workflows through a web interface or Taverna Lite
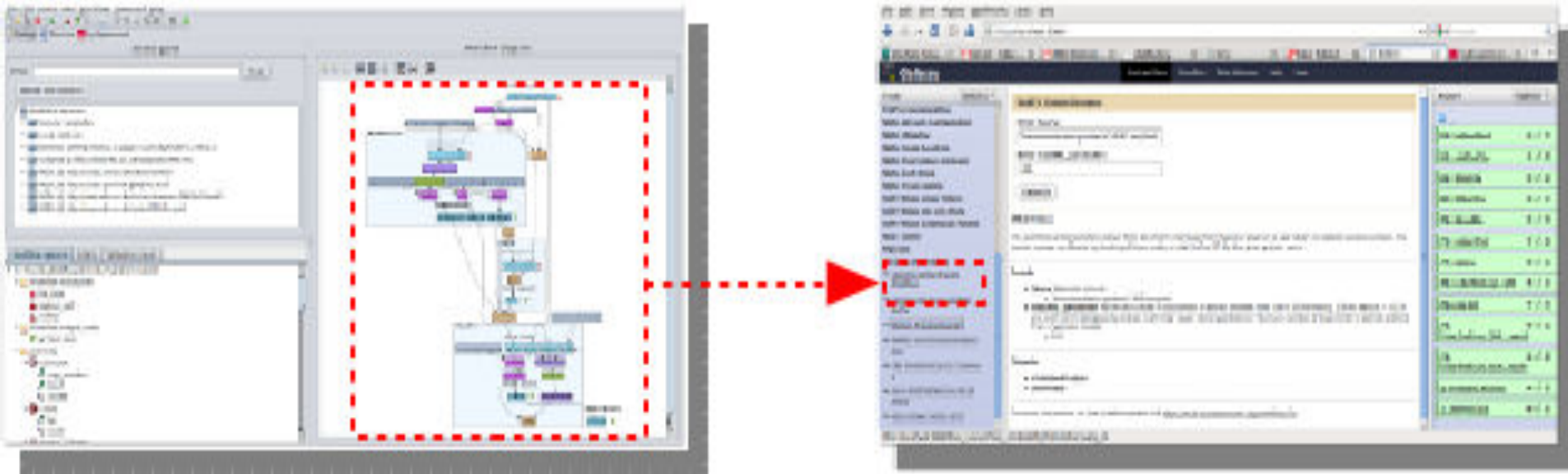
# Workflow Interoperability

# SHIWA: Sharing Interoperable Workflows

- Development of workflow interoperability technologies

  - Workflow development, testing and validation is a time consuming process and it requires specific expertise

  - Workflows developed for one workflow system  not compatible with workflows of other systems

# Converting Taverna Workflows into Galaxy Tools



## Why?

- Connect the Taverna and Galaxy communities
- Combine the power of Taverna with the simplicity of Galaxy
- Any Taverna service can be made available to Galaxy users

# **Summary**

- Informatics often relies on data integration and large-scale data analysis
- Workflows are a mechanism for linking together resources and analyses
- Automation
- Large data manipulation
- Promote reproducible research
- Easy to find and use successful analysis methods
- Allows scaling-up of computational resources on HPC, cloud and Grid

# More Information

Galaxy
- http://usegalaxy.org

- Taverna
  - http://www.taverna.org.uk

- WS-PGrade
  - http://www.myexperiment.org