



myExperiment 2.0 – Preserving digital Research Objects using the Wf4Ever architecture

Stian Soiland-Reyes
myGrid, University of Manchester

*EGI/SHIWA Workshops on e-Science Workflows
Budapest, 2012-02-10*





Taverna

<http://www.taverna.org.uk/>

my experiment

<http://www.myexperiment.org/>



UNIVERSITY OF
Southampton
School of Electronics
and Computer Science

Taverna - Scientific Workflow Management System

~85000 downloads

~EU projects: SCAPE, BioVeL, HELIO, e-Lico, VPH-SHARE, EGI-INSPIRE....

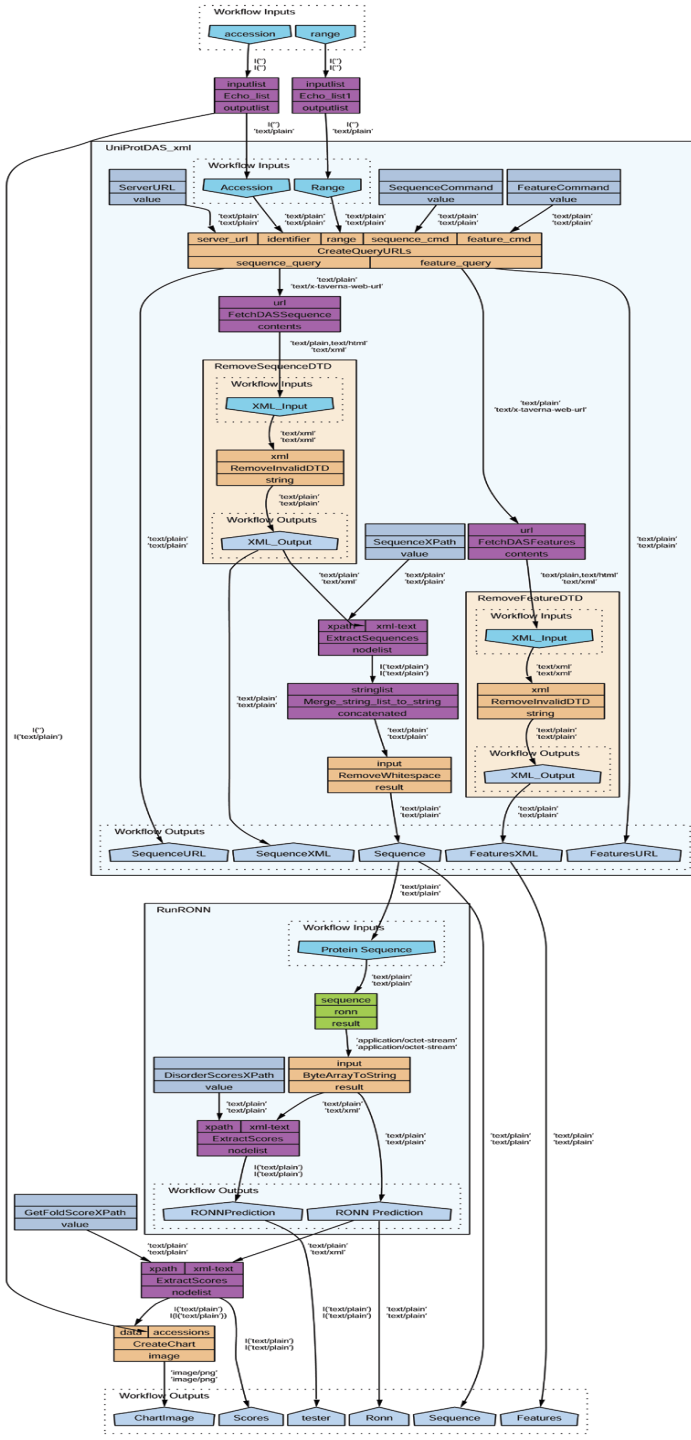
myExperiment - Web 3.0 virtual environment, library and social network for workflows

~5000 registered users

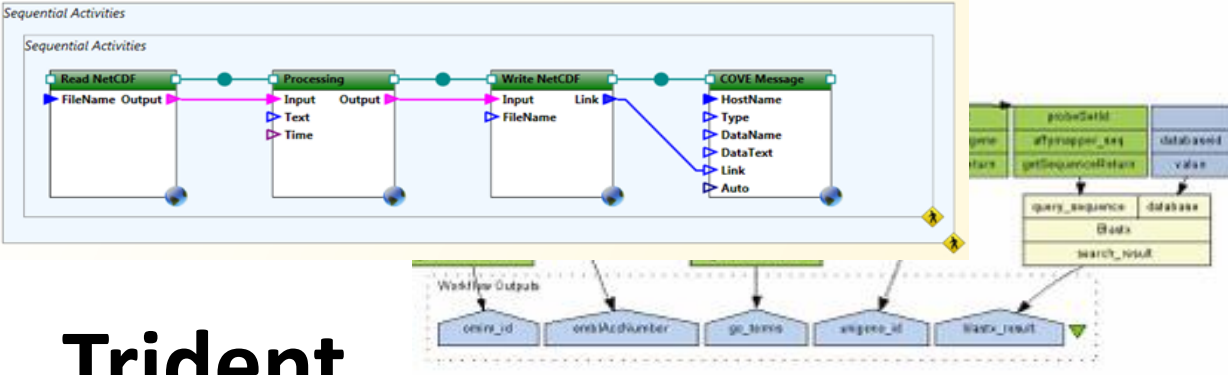
~2200 workflows

~21 different systems

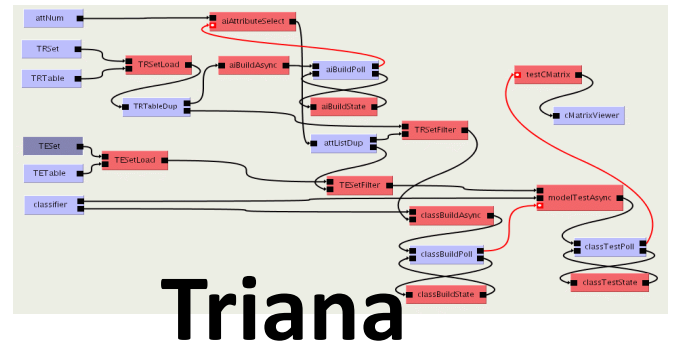
What is a Scientific Workflow?



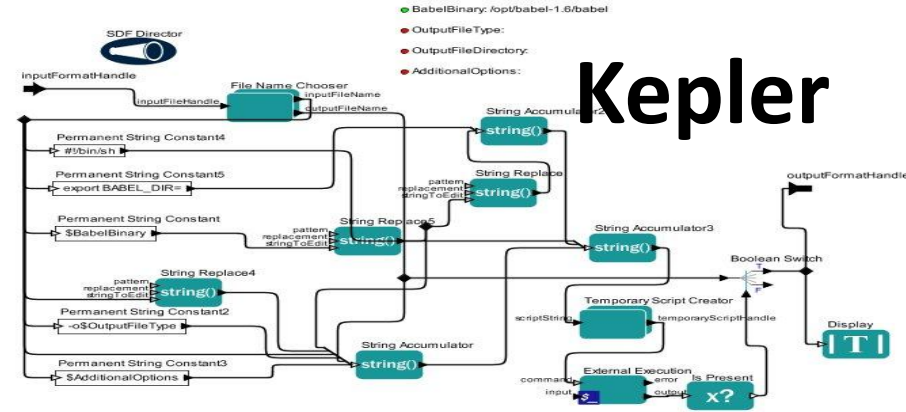
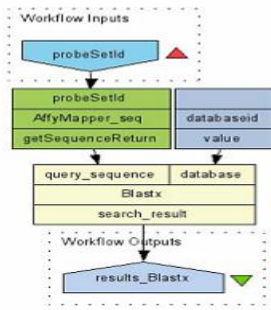
- » Workflows coordinate the execution of **services** and **link** together resources.
- » **Data-driven** rather than process-driven: «*Send output from A to B and C*»
- » Semi-automated computational execution in scientific problem-solving
→ *repeatable, reproducible, reusable*
- » The implementation of a **scientific method**



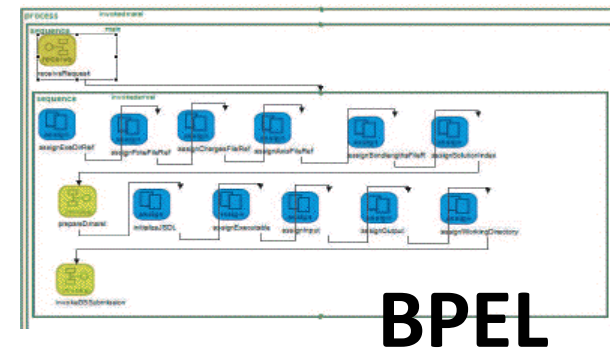
Trident



Triana

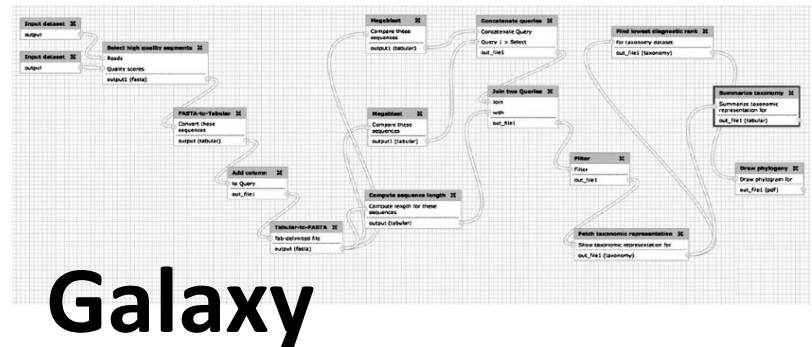
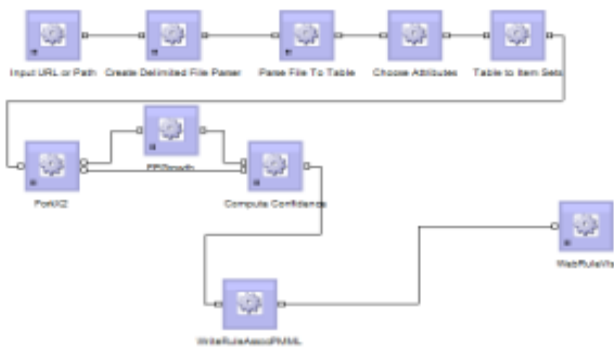


Kepler



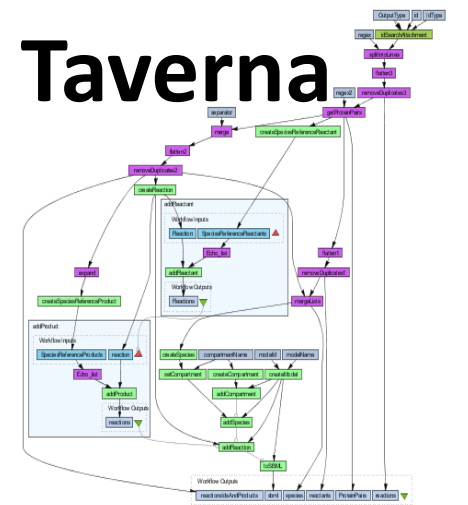
BPEL

Meandre

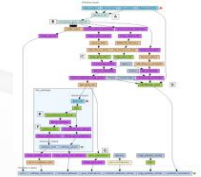


Galaxy

Taverna



- Paul writes **workflows** for identifying *biological pathways* implicated in *resistance to Trypanosomiasis* in **cattle**
- Paul meets Jo. Jo is investigating *whipworm* in **mouse**.
- Jo reuses one of Paul's workflow **without change**.
- Jo identifies the *biological pathways* involved in sex dependence in the mouse model, believed to be involved in the ability of mice to *expel the parasite*.
- Previously a manual **two year study** by Jo had failed to do this.



“A biologist would rather share their toothbrush than their gene name”



*Mike Ashburner and others
Professor in Dept of Genetics,
University of Cambridge, UK*

my experiment

<http://www.myexperiment.org/>

- “Facebook for Scientists”
...but different to Facebook!
- A repository of research methods
- A social network of people and things
- A Social Virtual Research Environment
- A probe into researcher behaviour
- Open source (BSD) Ruby on Rails app
- REST and SPARQL, Linked Data
- Influenced BioCatalogue, MethodBox and SysMO-SEEK

myExperiment currently has 5378 members, 292 groups, 2273 workflows, 534 files and 217 packs



[Home](#)
[Users](#)
[Groups](#)
[Workflows](#)
[Files](#)
[Packs](#)
[Services](#)
[Topics](#)
 Workflows ▼ Search

[Home](#) » [Workflows](#)
 BOOKMARK    ...

Workflows

Search filter terms

 « previous **1** [2](#) [3](#) ... [174](#) next »

 Sort by: Rank ▼

Showing 1739 results. Use the filters on the left and the search box below to refine the results.


Filter by type

- Taverna 2 752
- Taverna 1 562
- RapidMiner 192
- Kepler 43
- Bioclipse Script... 34
- GWorkflowDL 24
- LONI Pipeline 23
- BioExtract Server 16
- Trident (Packa... 10
- LabTrove Tem... 9

Filter by tag

- example 215
- mygrid 103

Taverna 2


Pathways and Gene annotations for QTL region
 [View](#)
 [Download \(v7\)](#)

(v7)

Original Uploader
Created: 19/11/09 @ 18:18:52 | **Last updated:** 02/09/11 @ 11:44:57

Credits:  Paul Fisher

License: Creative Commons Attribution-Share Alike 3.0 Unported License

 Paul Fisher


This workflow searches for genes which reside in a QTL (Quantitative Trait Loci) region in the mouse, *Mus musculus*. The workflow requires an input of: a chromosome name or number; a QTL start base pair position; QTL end base pair position. Data is then extracted from BioMart to

Workflow Entry: Success-Abandonment-Classification

Created at: 06/02/08 @ 14:35:41 Last updated: 02/07/08 @ 17:15:25

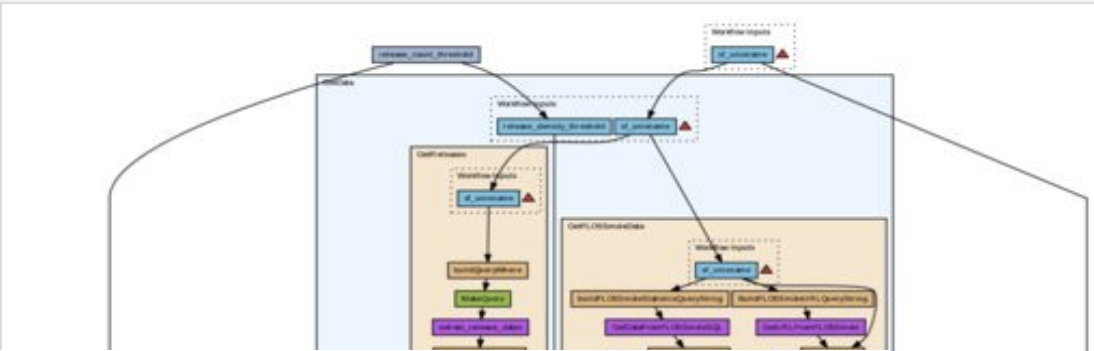
License | Credits (2) | Attributions (0) | Tags (6) | Featured in Packs (1) | Ratings (1) | Attributed By (0) | Favourited By (2) | Citations (0) | Version History | Reviews (0) | Comments (0)

Version 3 (latest) (of 3) View version: 3 (latest)
Version created on: 06/02/08 @ 14:35:41 by: Andrea Wiggins | Revision comments
Last edited on: 02/07/08 @ 17:15:25 by: Andrea Wiggins

Title: Success-Abandonment-Classification
Type: Taverna 1

Preview

(Click on the image to get the full size)



Workflow Type: Taverna 1
Original Uploader: Andrea Wiggins
License: All versions of this Workflow are licensed under: CC BY SA
Credits (2): Andrea Wiggins, James Howison
Attributions (0)

New/Upload
Workflow GO

David De Roure

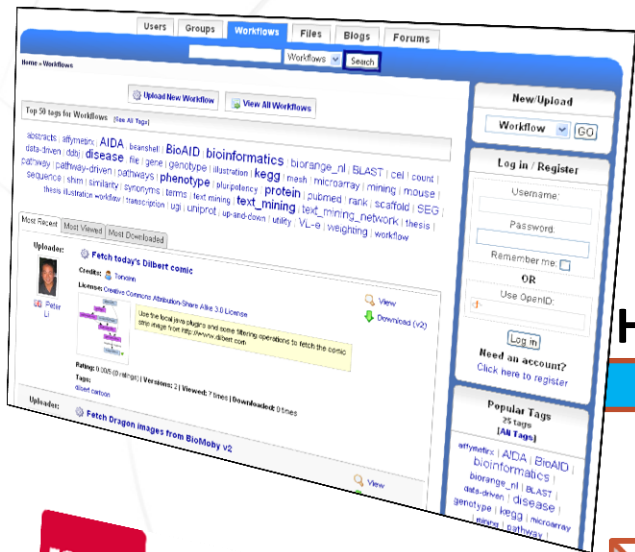
- My Profile [edit]
My Messages (3)
My Memberships (1)
My History
My News

Manage Announcements

- 3 new messages
Hi David
Sean Bechhofer is n...
Invitation to 'Wf4E...

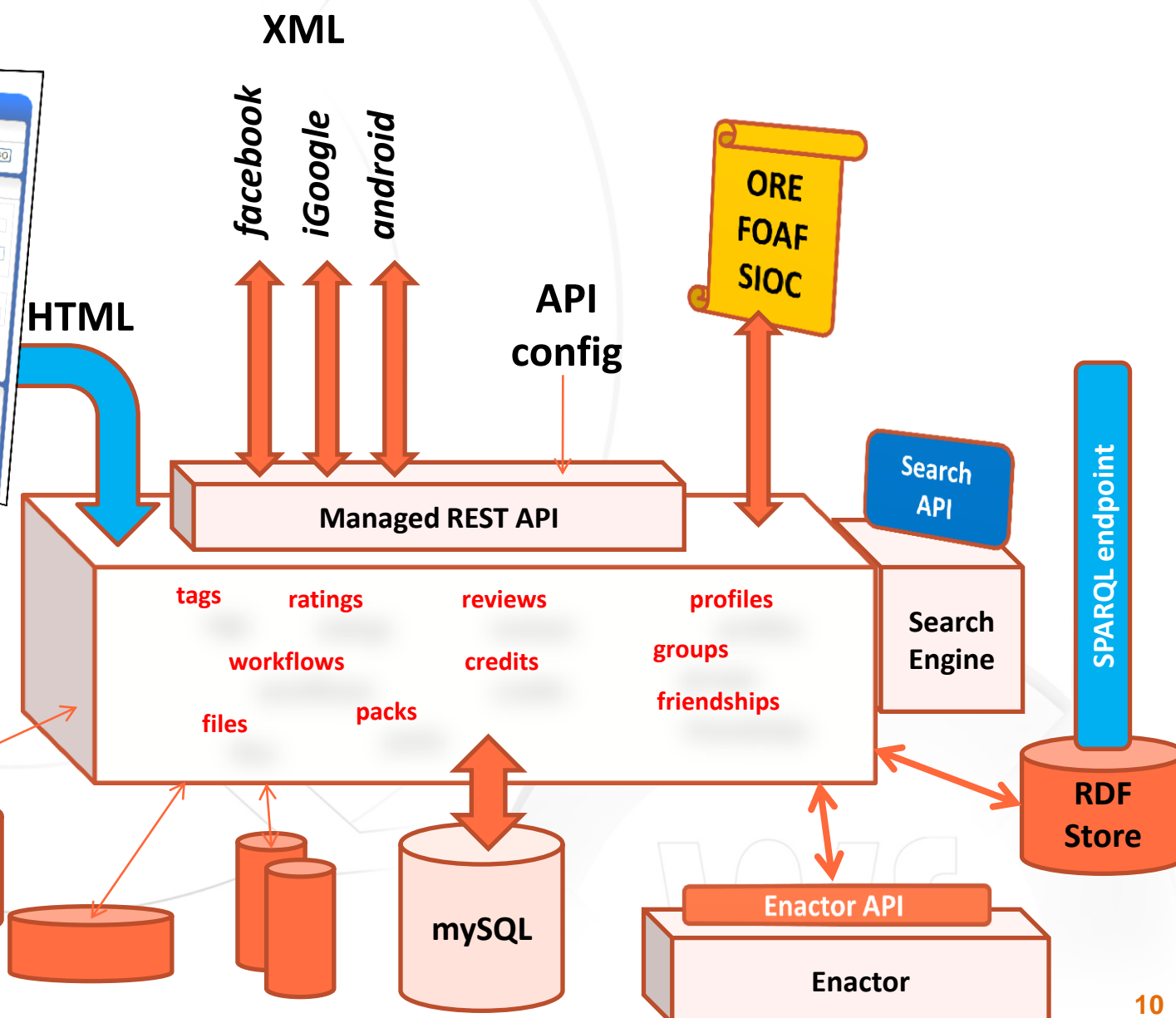
- 2 new friendship requests
Yehia El-khatib
mihailionita_me

- 1 new group request
From Pique
(for Group: Wf4Ever)



my experiment

Eprints
Dropbox
...



- » myExperiment plugin for Taverna
 - › Browse myExperiment workflows
 - My workflows
 - Tags
 - Search
 - › Open workflow
 - + Embed in existing workflow
 - › Upload workflows
 - Provide metadata

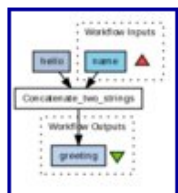


Taverna

Hello Anyone (version 1)

Uploader: Stian Soiland-Reyes

Type: Taverna 2



An extension to [helloworld.t2flow](#) - this workflow takes a workflow input "name" which is used to produce a string using the local produced string

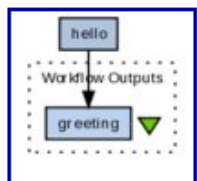
Open in myExperiment

Preview Download

Hello World (version 1)

Uploader: Stian Soiland-Reyes

Type: Taverna 2



One of the simplest single workflow produced by Taverna

Open in myExperiment

Preview Download

Example of explicit looping

Uploader: Alan Williams

Workflow Entry: Pathways and Gene annotations for QTL Phenotype annotated

Type: Taverna 1

Uploader: Stian Soiland-Reyes
Created at: Tue Oct 30 13:04:02 GMT
License: Creative Commons Attribution-NoDerivs

Import workflow

Import method

- Add a **nested workflow** into the destination workflow as a single service. The workflow is imported *separately*, but is shown expanded in the diagram of the parent workflow. In this method, you can connect directly to the input and output ports of the nested workflow.
- Merge** a workflow by copying all services, ports and links directly into the destination workflow. This method is more beneficial for merging smaller workflow fragments. For inclusion of larger workflows you may prefer the nested workflow method.

Workflow destination

- Already opened workflow: C:\Users\stain\Data\Desktop\helloanyone.t2flow

Import

Prefix:

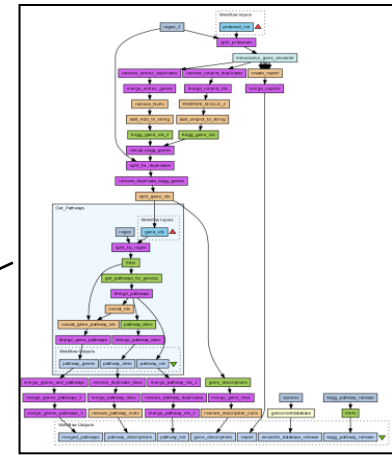
Optional prefix to be prepended to the name of the inserted services and workflow. Duplicate names will be resolved by adding numbers, for instance my_service_2

Paul's Research Object

Workflow 16

QTL

Results



```

path:mmu04060 Cytokine-cytokine receptor interaction - Mus musculus (mouse)
path:mmu00970 Aminoacyl-tRNA biosynthesis - Mus musculus (mouse)
path:mmu00240 Pyrimidine metabolism - Mus musculus (mouse)
path:mmu03010 Ribosome - Mus musculus (mouse)
path:mmu04080 Neuroactive ligand-receptor interaction - Mus musculus (mouse)
path:mmu05220 Apoptosis - Mus musculus (mouse)
path:mmu05220 Chronic myeloid leukemia - Mus musculus (mouse)
path:mmu04612 Antigen processing and presentation - Mus musculus (mouse)
path:mmu00271 Methionine metabolism - Mus musculus (mouse)
path:mmu04912 GnRH signaling pathway - Mus musculus (mouse)
path:mmu04330 Notch signaling pathway - Mus musculus (mouse)
path:mmu04640 Hematopoietic cell lineage - Mus musculus (mouse)
path:mmu00561 Glycerolipid metabolism - Mus musculus (mouse)
path:mmu04110 Cell cycle - Mus musculus (mouse)
path:mmu04530 Tight junction - Mus musculus (mouse)
path:mmu02010 ABC transporters - General - Mus musculus (mouse)
    
```

produces

Included in

Analysis Protocol for Candidate Genes and Pathways

This protocol is aimed at providing a guide to the interpretation of the results obtained from both the QTL and microarray workflows. Each workflow provides a series of text files, which are to be used as a means of obtaining the pathways which relate to differentially expressed genes in the microarray study and genes located within the chosen QTL region.

The output from each workflow consists of the following files:

- ensembl_database_release.text
- pathway_descriptions.text
- gene_descriptions.text
- pathway_descriptions.text
- merged_pathways.text
- kegg_external_gene_reference.text
- report.text
- pathway_list.text

ensembl_database_release.text

The current release of the Ensembl dataset for the chosen species, e.g. *Mus musculus*. Although this uses the programmatic interface of Ensembl, it can be used to identify which release was used to generate the list of genes in the QTL region or mapping of Affymetrix probe IDs.

Logs

produces

Metadata

Tags (19)

Creator tags

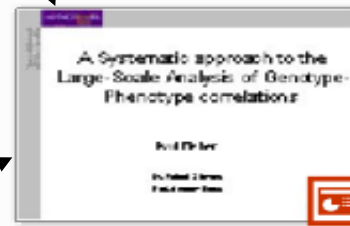
affymetrix | african
trypanosomiasis | cattle |
data-driven | disease | entrez
| genotype | Kegg
Pathways | KeggID |
link-integration | **microarray**
| mouse | pathway |
pathway-driven |
phenotype | sleeping
sickness | swissprot | uniprot |
web services

[edit]

Add Tags

Feeds into

Included in



complete.ppt

Slides

Published in



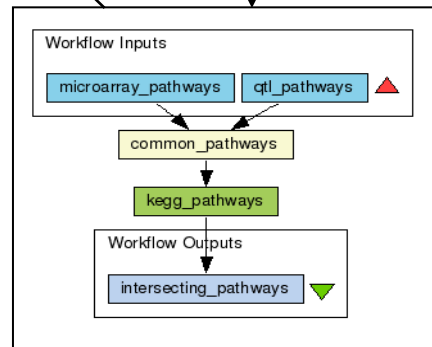
Paper

Included in

Published in

Common pathways

Workflow 13



produces

Primer Name	Left Flank Nucleotide sequence	Right Flank Nucleotide sequence
DAXX_1274_1812	CAGGAGGAATGGCGAGTG	AGCTTAGTCTTCCCAAGCC
DAXX_140754_456_1070	CTTGTAGGATTTGGACTGGG	TCTCCTCTCTTCTCCTCC
DAXX_2270_2720	TGGCAGGAGAGATGGTTC	ATGGTTCAAGGGAAGGGAAA
DAXX_2644_3187	TGTGTGATTGGCTGTGTT	GCAAATACGAGGAGTCTGGG
DAXX_exon5	TCCTCCTCTACCAATCAAA	AGCAGAACTAACCCACAAGG
Daxx_Upst_479_1104	CAGGCTTCTCATCAACACC	TGTCCTATGGCTGTGCAGG

Results

myGrid



Reusable. The key tenet of Research Objects is to support the sharing and reuse of data, methods and processes.

Repurposeable. Reuse may also involve the reuse of constituent parts of the Research Object.

Repeatable. There should be sufficient information in a Research Object to be able to repeat the study, perhaps years later.

Reproducible. A third party can start with the same inputs and methods and see if a prior result can be confirmed.

Replayable. Studies might involve single investigations that happen in milliseconds or protracted processes that take years.

Referenceable. If research objects are to augment or replace traditional publication methods, then they must be referenceable or citeable.

Revealable. Third parties must be able to audit the steps performed in the research in order to be convinced of the validity of results.

Respectful. Explicit representations of the provenance, lineage and flow of intellectual property.

- **Workflow** – pack contains a number of workflows
 - **Presentation** - encapsulation of a single presentation
 - **Collection** - a number of things: workflows, presentations, papers
 - **Heterogeneous** - where the workflows do not appear to have a clear common purpose
 - **Homogeneous** - workflows appear to be designed to work together
- **Paper** - source for a paper
 - **Tutorial** - tutorial material
 - **Data** - collection of data files
 - **Derived data** - results of workflow
 - **Benchmark** - benchmarking data
 - **Supplementary** - stuff associated with a paper
 - **Noise** - tests, tryouts, rubbish
 - **Oddity** - none of the above

Analysis by Sean Bechhofer



- Workflow Preservation
 - Research Objects
 - Provenance
 - Recommendation
- Astronomy and Genomics



Project ID card

Funded under: 7th FWP (Seventh Framework Programme)

Area: Digital Libraries and Digital Preservation. (ICT-2009.4.1)

Project reference: 270192

Total cost: 3.86 million euro

EU contribution: 2.94 million euro

Execution: From 2010-12-01 to 2013-11-30

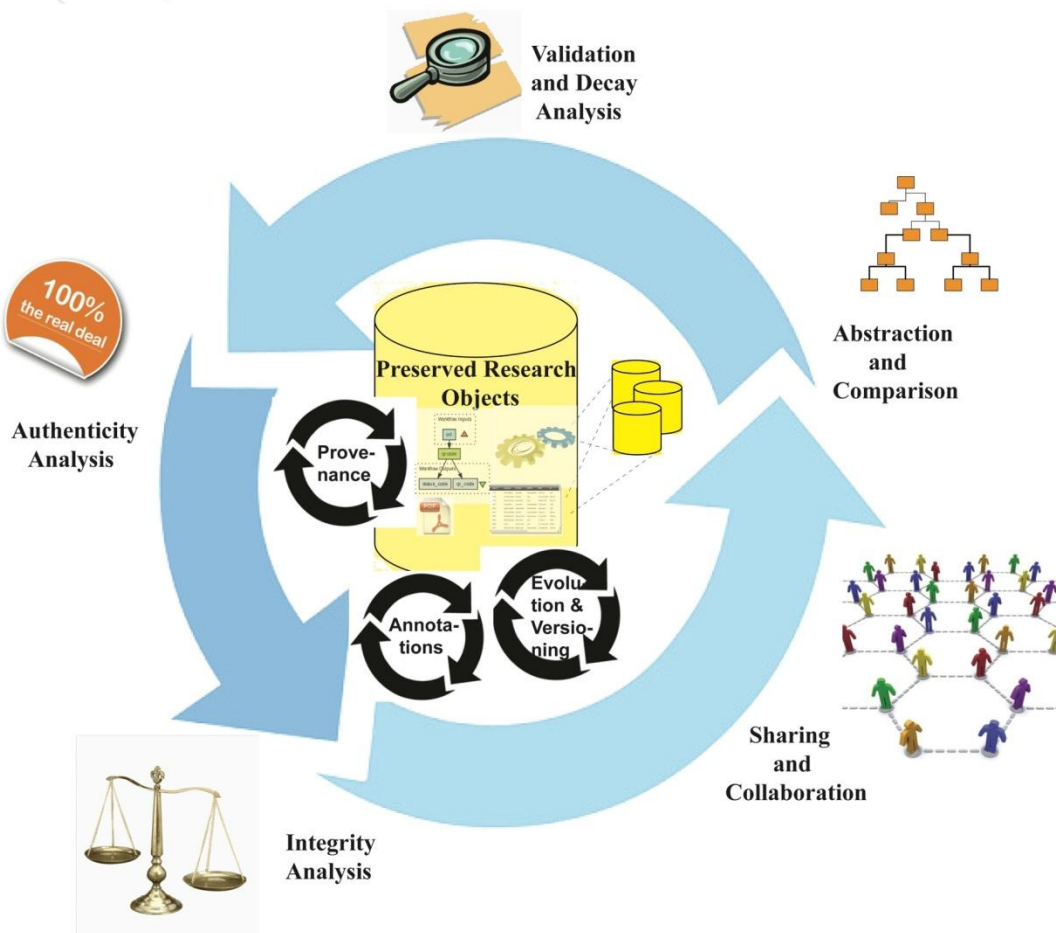
Duration: 36 months

<http://www.wf4ever-project.org/>



Preservation of scientific workflows in data-intensive science

- » Scientific workflows aim at the heart of experimental science
 - › Enable **automation** of scientific methods
 - › Encourage **best practices**
- » Need to be **preserved**
 - › **Reuse** is fundamental for incremental scientific development
 - › **Method reproducibility** is key for credit and publication
- » ...but workflow preservation is complex
 - › Heterogeneous types of information need to be **aggregated**, including workflows and related resources into **research objects**
 - › Research objects need to be **trusted** and **understandable** n years from now
 - › **Social aspects** need to be addressed in order to support reuse in scientific communities



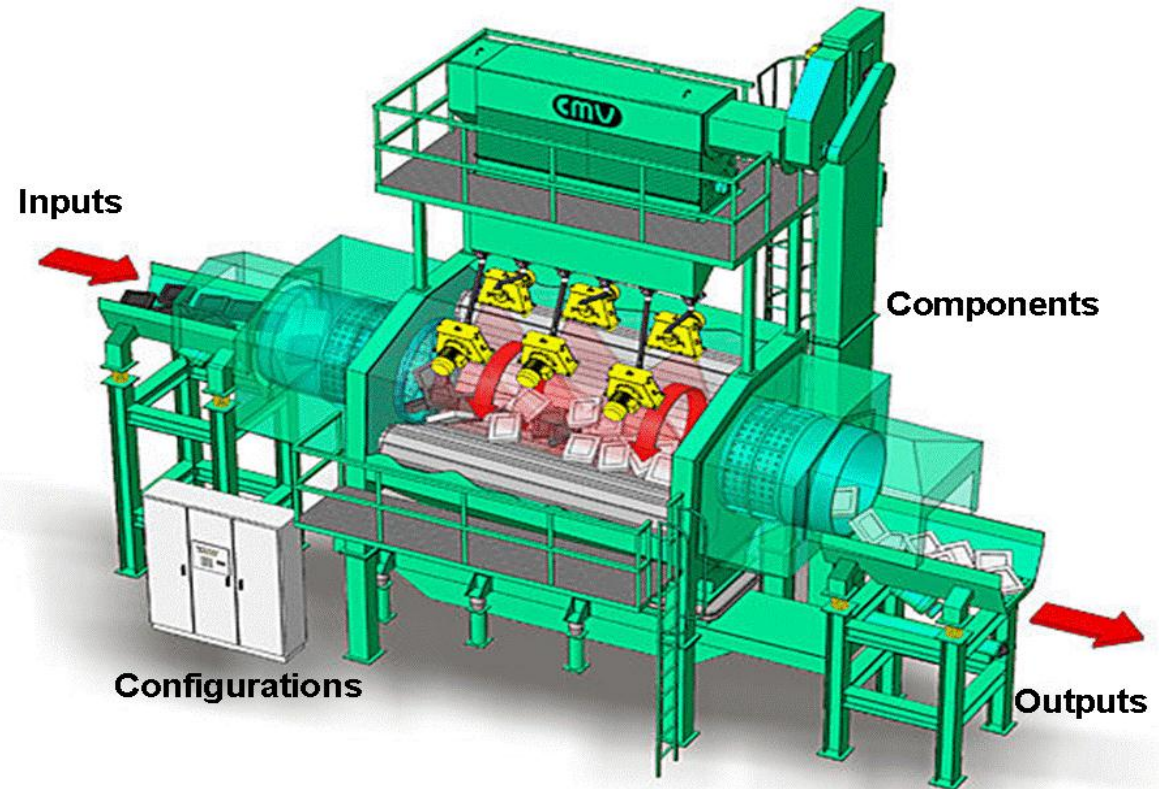
Stability, Completeness, Integrity, Authenticity, Quality

Workflow Decay

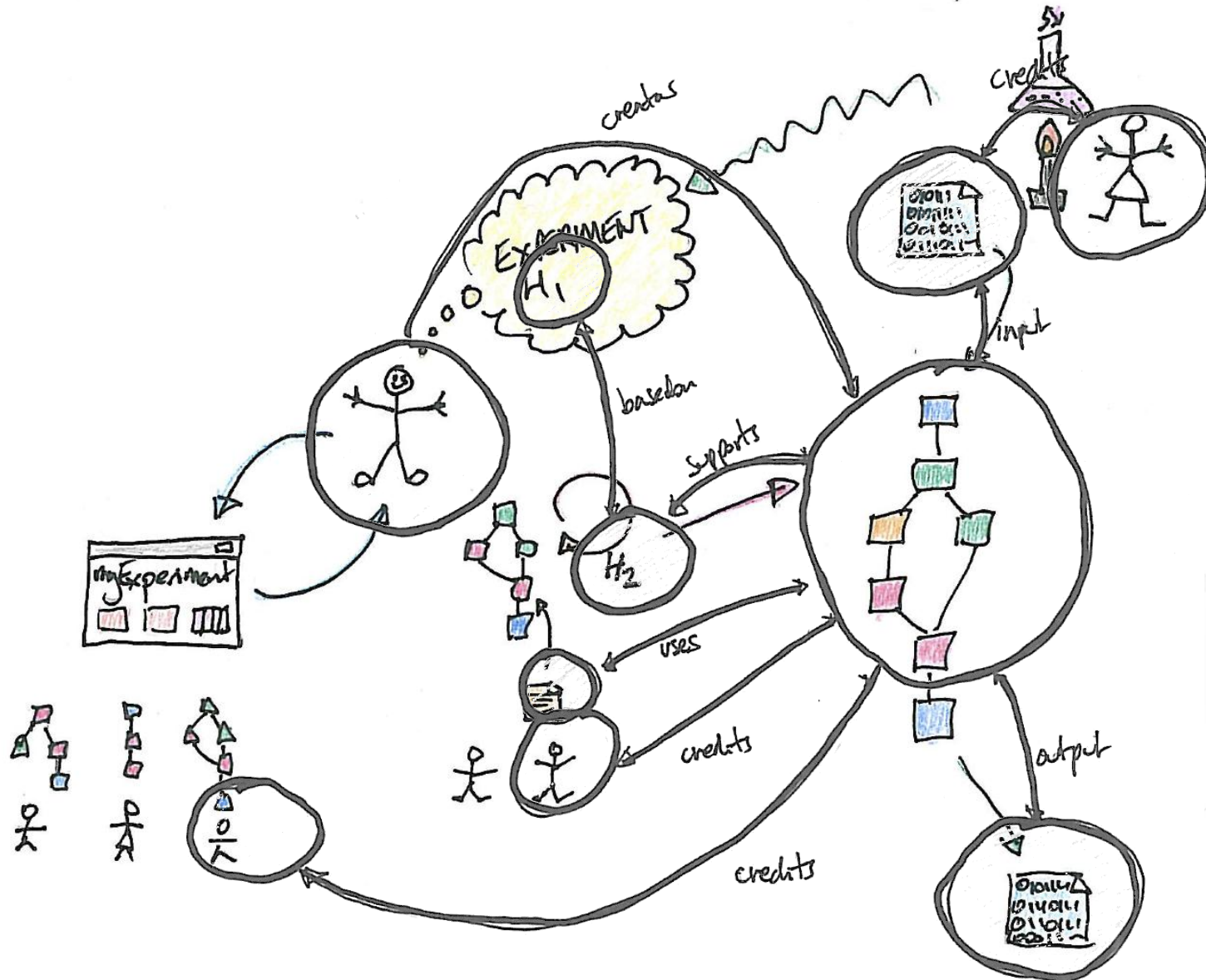
- Component level
- flux/decay/unavailability
- Data level
 - formats/ids/standards
- Infrastructure level
 - platform/resources

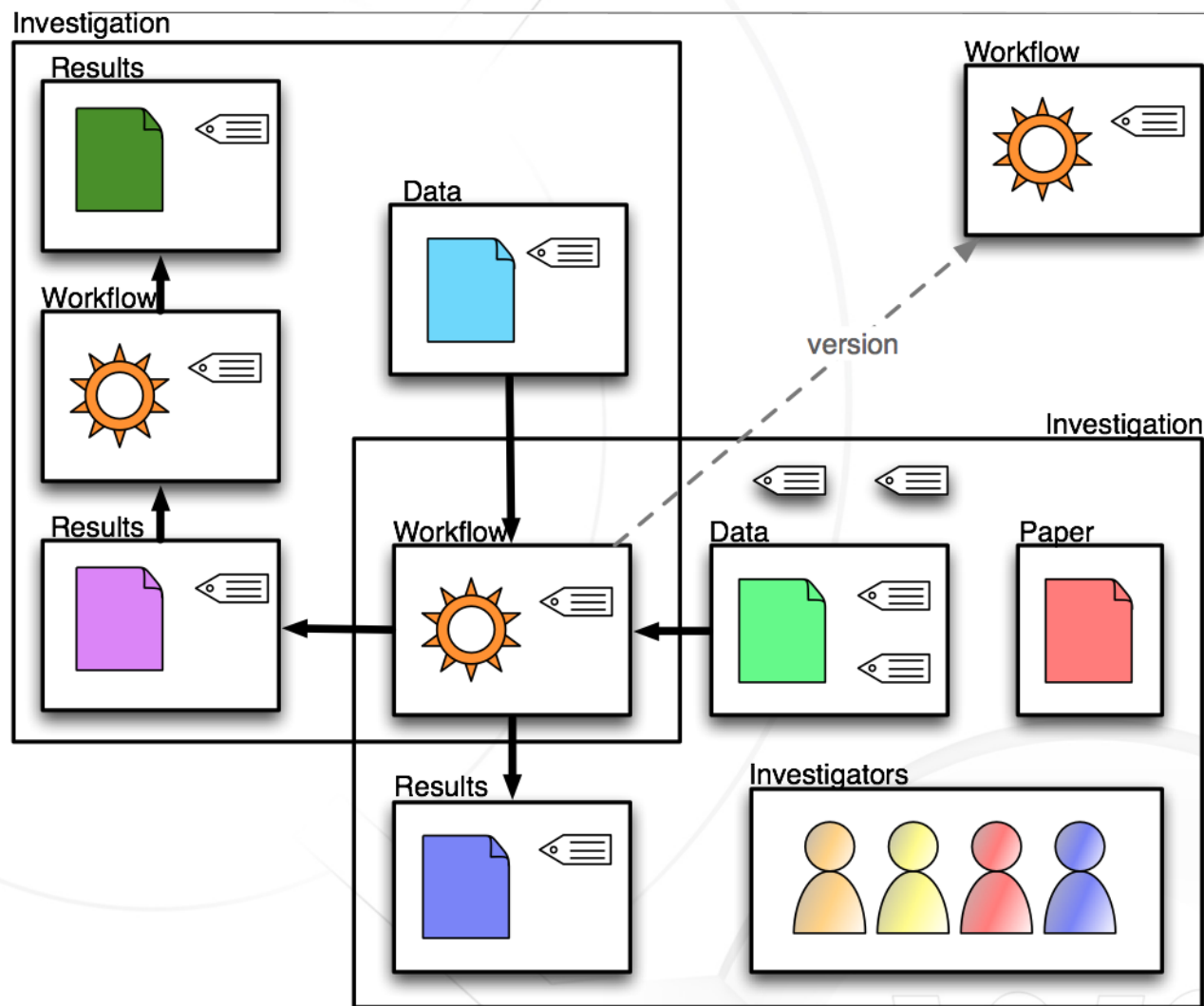
Experiment Decay

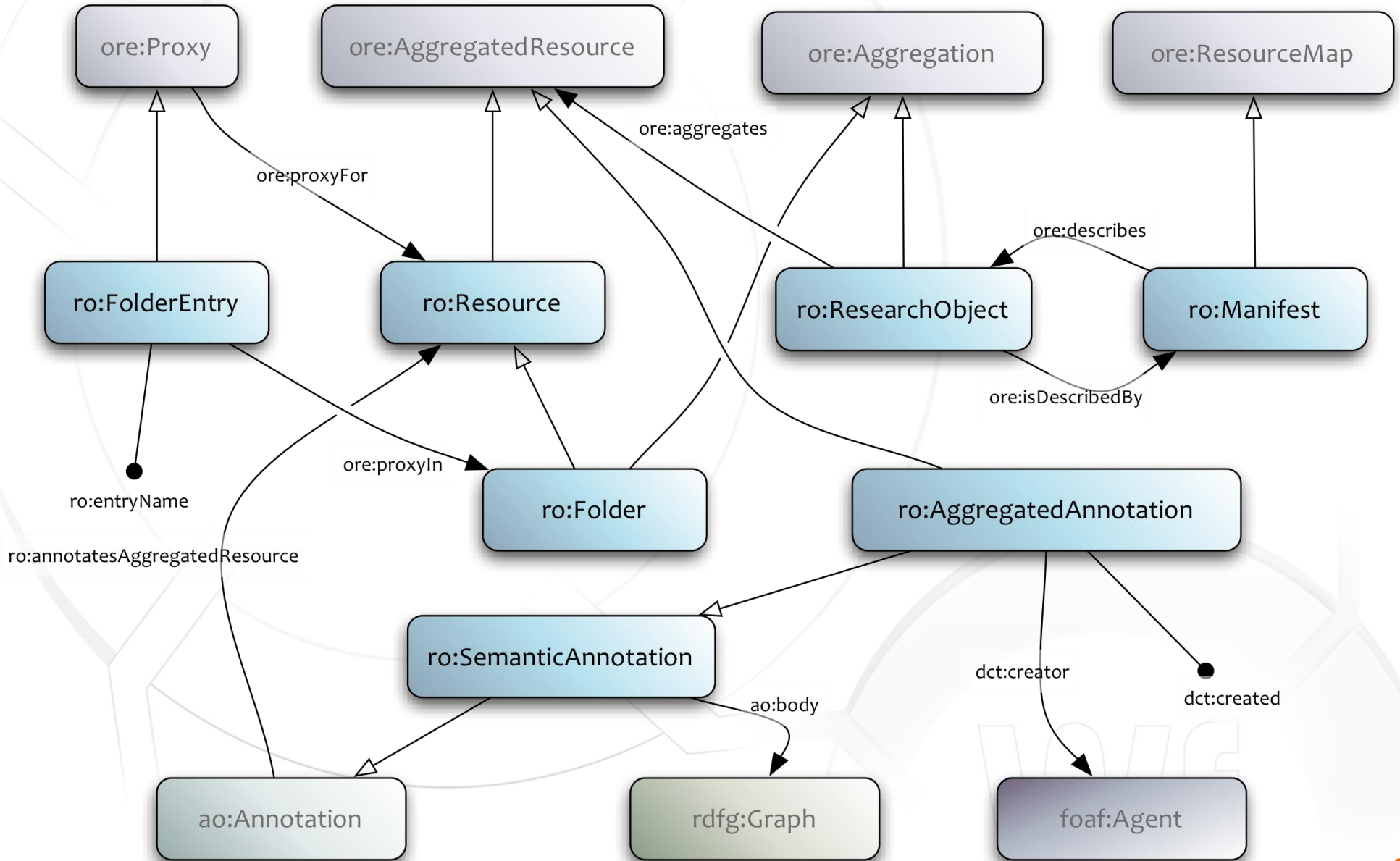
- Methodological changes
- New technologies
- New resources/components
- New data



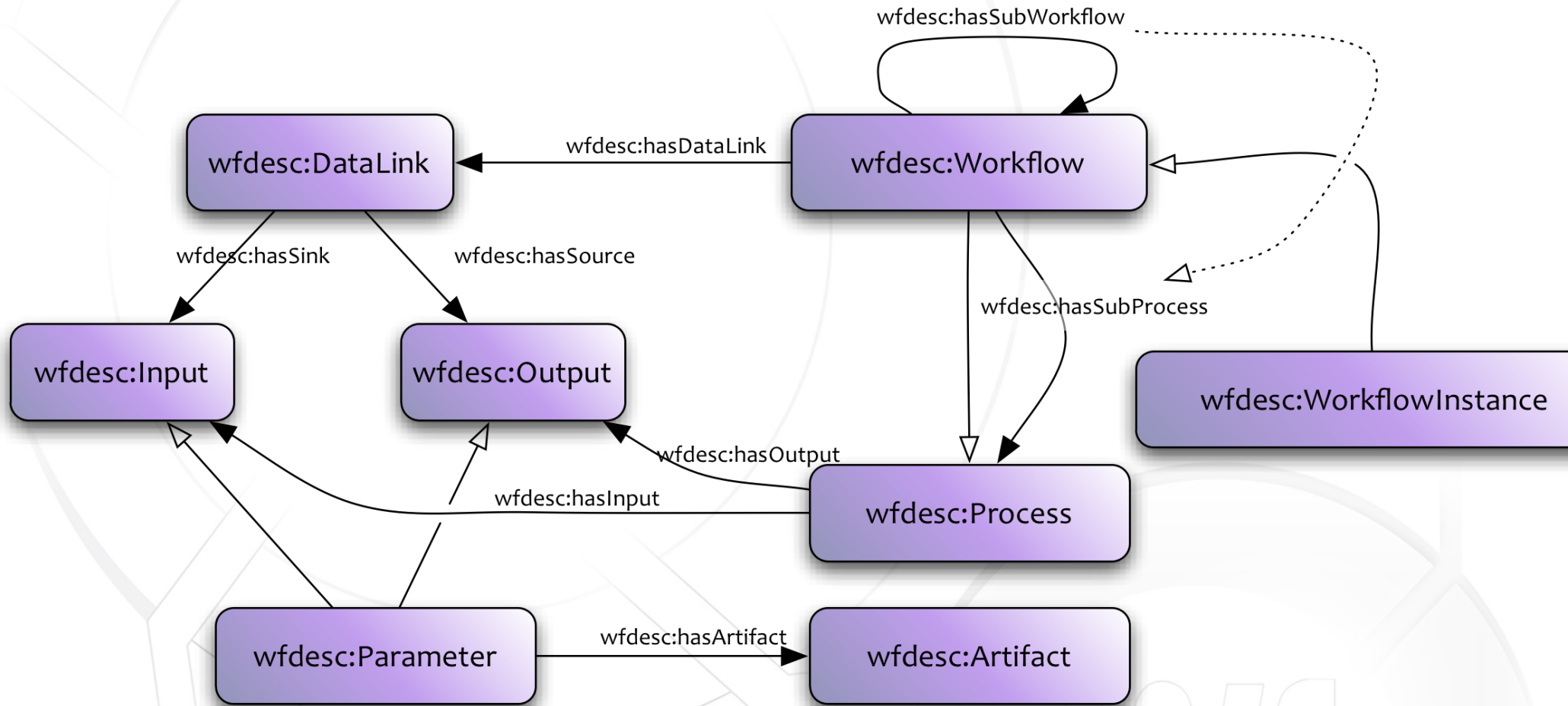
Research Objects as Social Objects



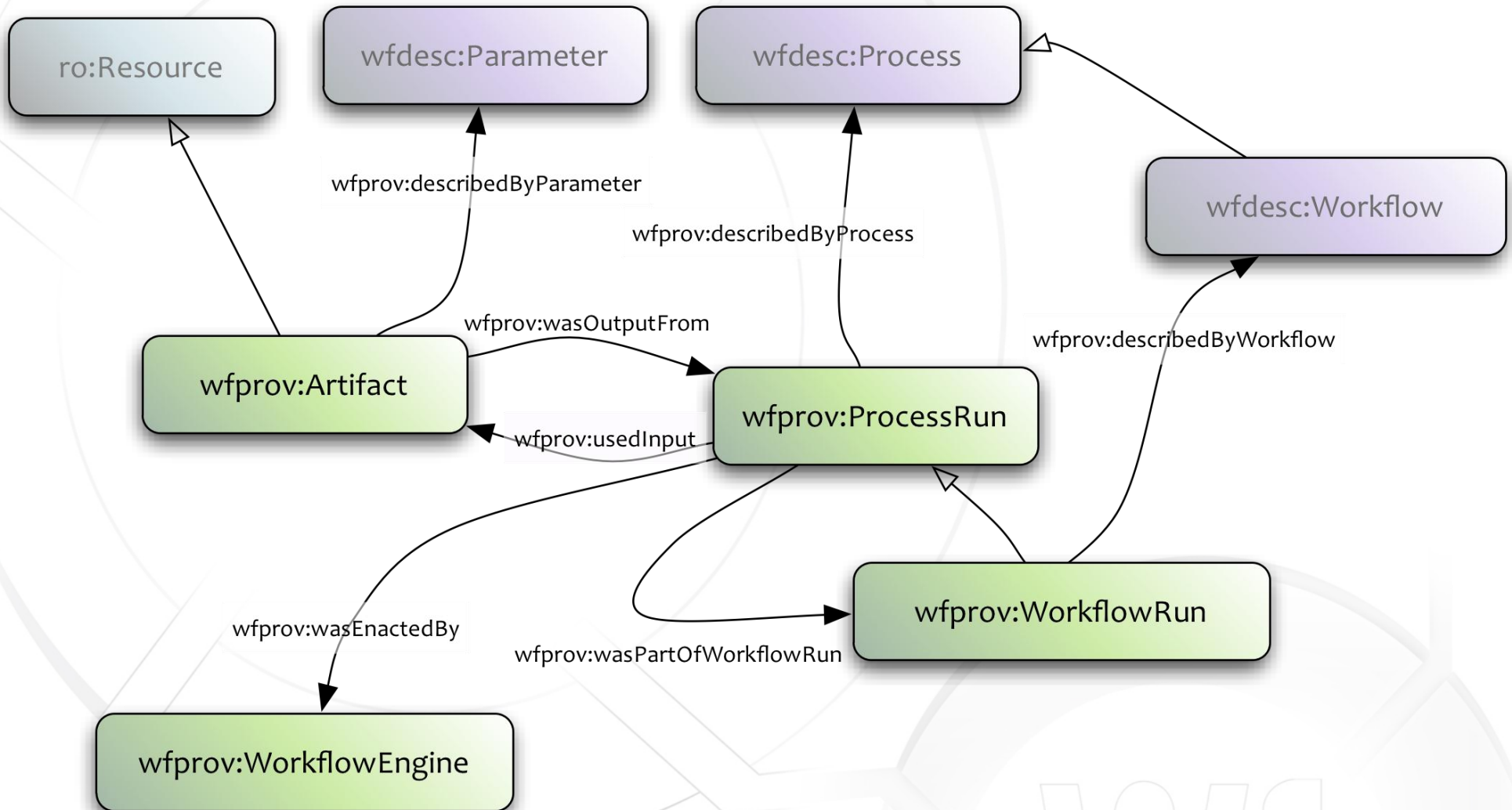




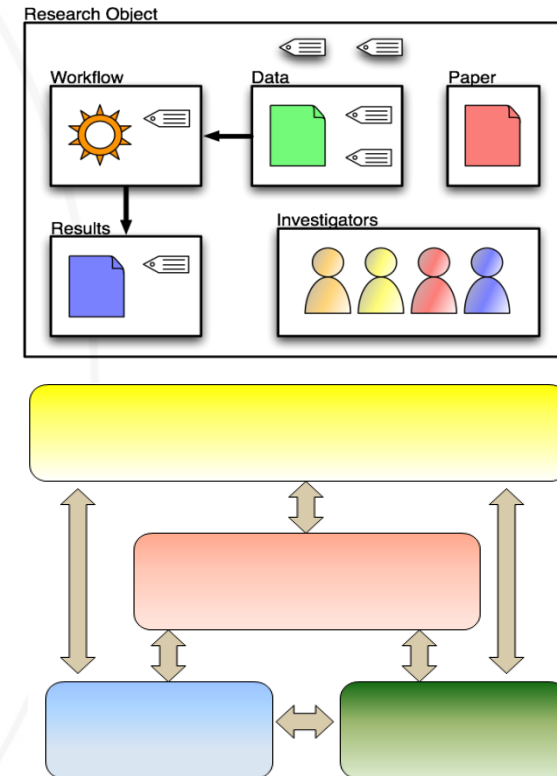
Workflow Description (wfdesc)

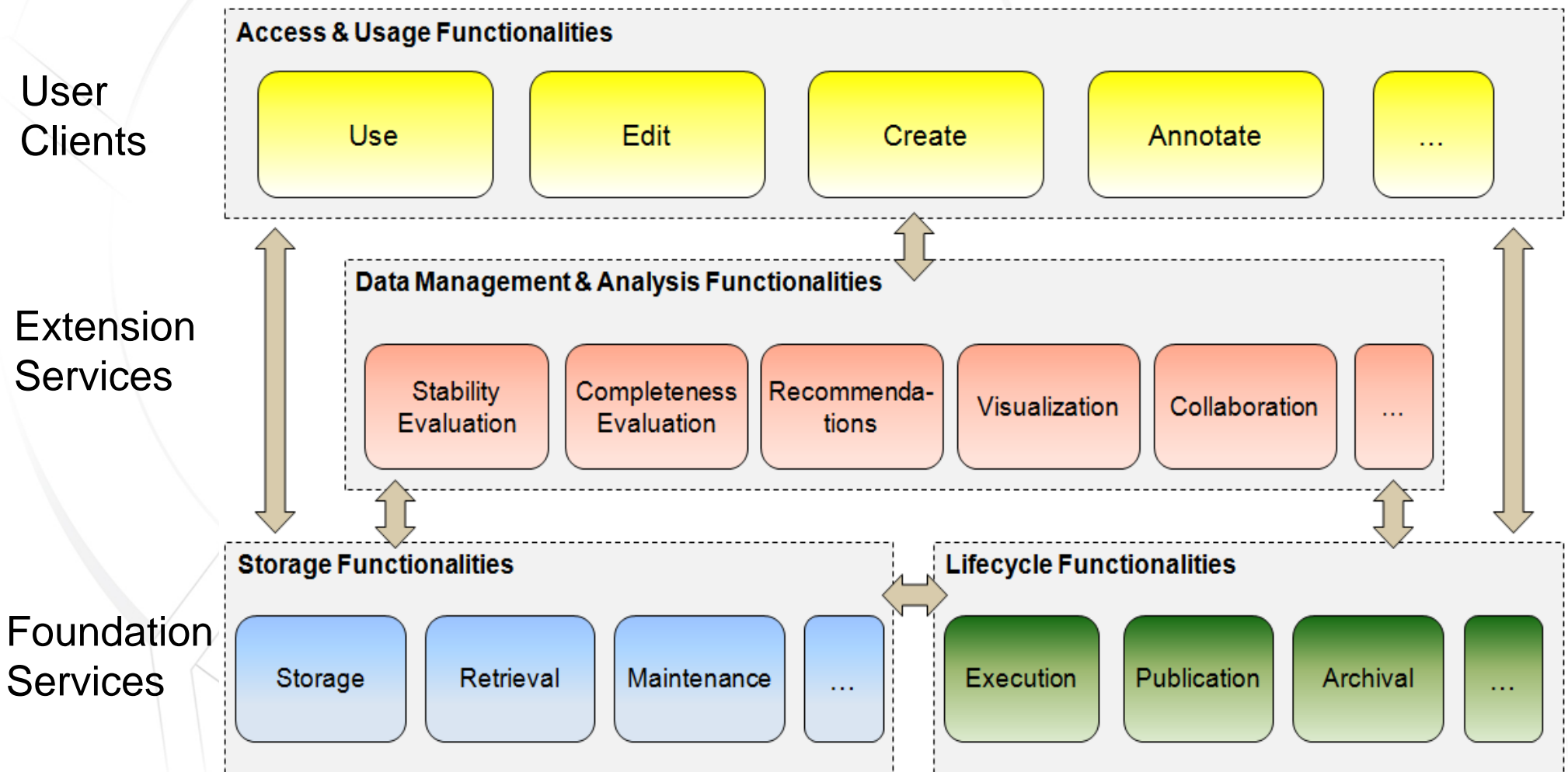


Workflow Provenance (wfprov)



- **Models**
 - Research Object
 - Annotation
 - Provenance
 - Evolution and Versioning
 - *Semantic Web Encoding*
- **Services**
 - Foundational, Extension, User
 - APIs, Architecture
 - *Web protocols/services*
- **Principles**
 - Map into standards
 - Adopt standards
 - Lightweight components
- **Ecosystem**
 - Command line
 - Third party systems





- » Analyse decay within all myExperiment workflows
 - › Estimate: Roughly 50% don't run correctly
 - › But why? Services gone? Did they ever work?
 - › How have the community evolved those workflows?
- » Service and wf substitutions; recommendations
- » Provenance analysis and use
 - › e.g. verifiability, replayability
- » SHIWA approach – can handle “What if the workflow system stops working”

Any Questions?

<http://www.myexperiment.org/>

<http://www.wf4ever-project.org/>

<http://www.mygrid.org.uk/>

<http://www.taverna.org.uk/>



*This work is licensed under the **Creative Commons Attribution 3.0 Unported License**. To view a copy of this license, visit <http://creativecommons.org/licenses/by/3.0/> or send a letter to Creative Commons, 444 Castro Street, Suite 900, Mountain View, California, 94041, USA.*