# Research Objects: Preserving Scientific Workflows and Provenance
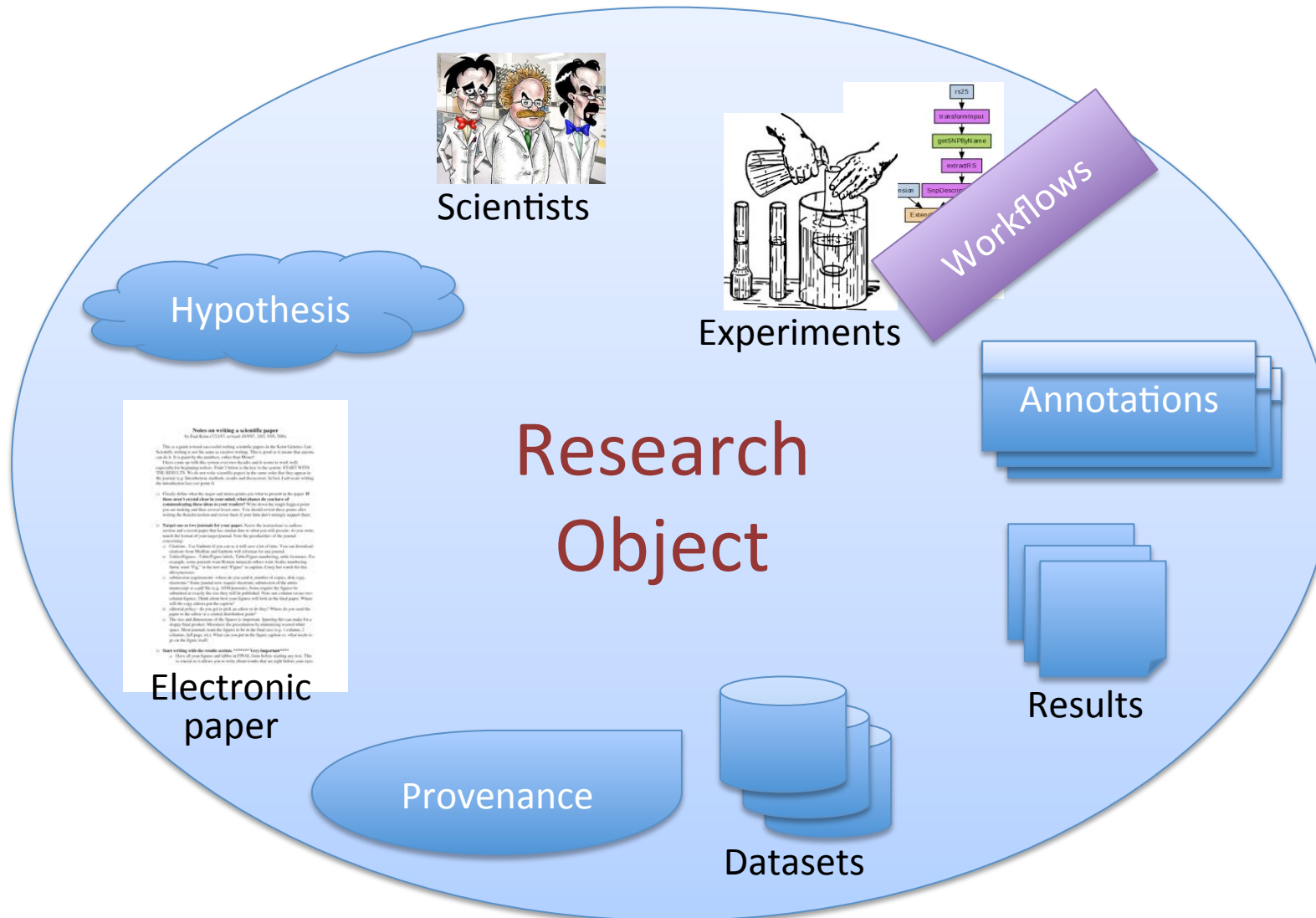
Khalid Belhajjame
Université Paris Dauphine

1

# Storyline

- Why Research Objects?

- Overview of the Research Objects

- Portfolio of Research Object Management Tools

- Provenance Distillation Through Workflow Summarization
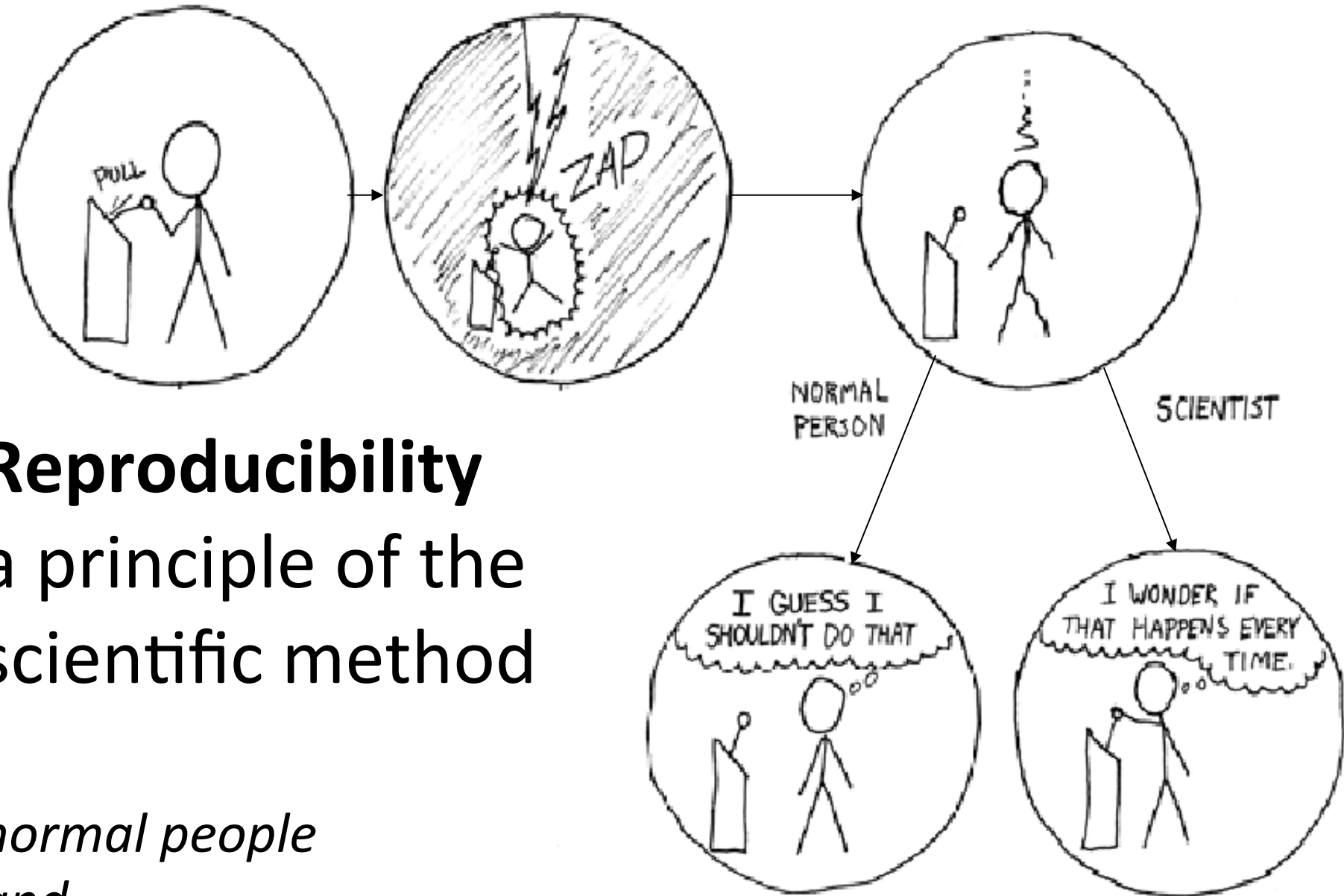
# Electronic papers are not enough



Electronic
paper

# Electronic papers are not enough

Scientists

Workflows

Experiments

Hypothesis

Annotations

Research Object

Electronic paper

Results

Provenance

Datasets

4

# Benefits Of Research Objects

- A research object aggregates all elements that are necessary to understand research investigations.

- Methods (experiments) are viewed as first class citizens

- Promote reuse

- Enable the verification of reproducibility of the results

**Reproducibility**
a principle of the
scientific method
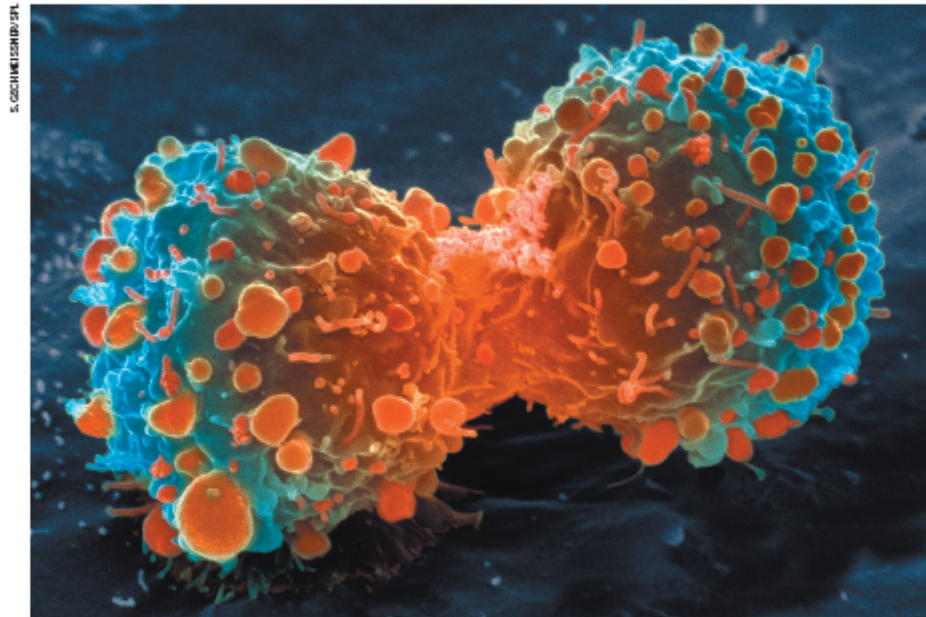
*normal people
and
scientist*

http://xkcd.com/242/

Many landmark findings in preclinical oncology research are not reproducible, in part because of inadequate cell lines and animal models.

47 of 53 "landmark" publications could not be replicated

Inadequate cell lines and animal models

Nature, 483, 2012

# Overview of the Research Object Model

# Research Object as an ORE Aggregation

# Scientific Workflows



- Data driven analysis pipelines

- Systematic gathering of data and analysis tools into computational solutions for scientific problem-solving

- Tools for automating frequently performed data intensive activities

**Provenance** for the resulting datasets
  - The method followed
  - The resources used
  - The datasets used

**WF Description
Prospective Provenance:**
Intended method for analysis

**WF Execution Trace
Retrospective Provenance:**
Actual data used, actual invocations, timestamps and data derivation trace

*PROV Primer, Gil et al*

# Specifying Workflows using WfDESC

# Specifying Workflow Provenance using WfPROV

# Portfolio of Research Object Tools



- Basic support for RO management
- Focus on developers
- Command-line tool

**RO Manager**

**RO management tools**

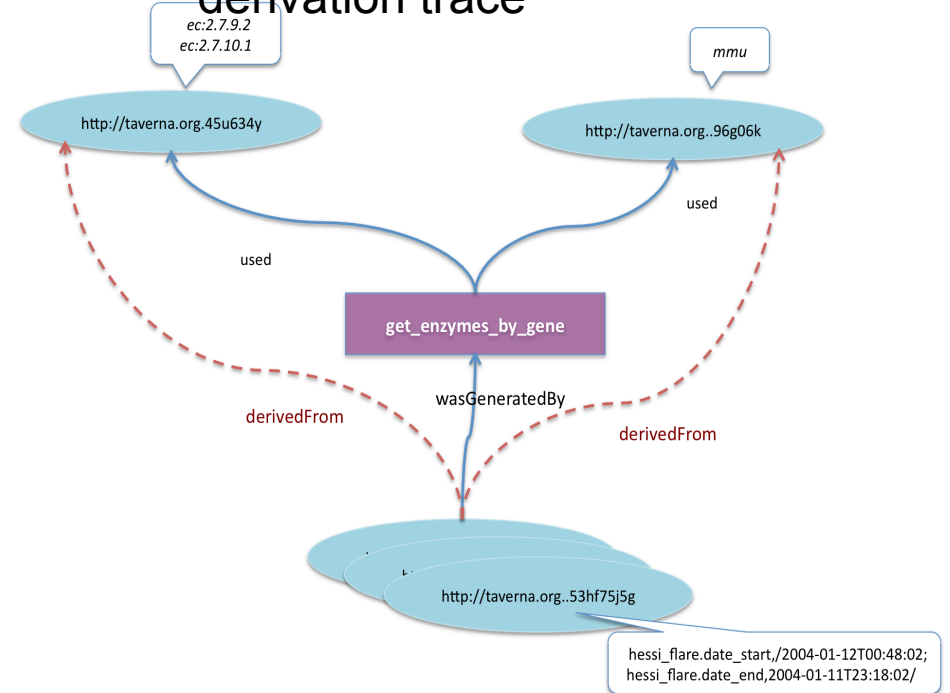**my experiment**

myExperiment
- Focus on scientists
- Online platform
- Workflow-centric
- Existing user community
- Incremental added value

**RODL**

- Holistic approach to content preservation
- Focus on scientists, librarians, others
- Online platform
- Research assets
- Non workflow-specific

14

# DEMO

# Workflows can get complex!

- Overwhelming for users who are not the developers
- Abstractions required for reporting
- Lineage queries result in very long trails

# Overall Approach



Workflow4Ever

Workflow Designer

Taverna Workbench

Motif Catalog

WF Description

Summarizer

WF Summary

Summarization Rules

SPARQL

OWL

# PART-1: Scientific Workflow Motifs

- Domain Independent categorization
  - Data-Oriented Nature
  - Resource/Implementation-Oriented Nature

- Captured In a lightweight OWL Ontology

http://purl.org/net/wf-motifs

# Motif annotations over operations



motifs(color_pathway_by_objects) = {m1:DataRetrieval}

motifs(Get_Image_From_URL_2) = {m2:DataMoving}

# PART-2: Workflow reduction primitives

- Collapse (Up/Down)
- Compose
- Eliminate

# Eliminate

# Two sample strategies

- ## By-Elimination

  – Minimal annotation effort

  – Single rule

  If $\exists$ path $p$ in $W$
  where $p = op_A$
  and $motif(op_A)$ contains $< m1 : DataPreparation >$,
  then $eliminate\_op(W, op_A)$

- ## By Collapse

  – More specific annotation

  – Multiple rules

  If $\exists$ path $p$ in $W$
  where $p = op_A$
  and $motif(op_A)$ contains $< m1 : Augmentation >$,
  then $collapse\_op\_downstream(W, op_A)$

  If $\exists$ path $p$ in $W$
  where $p = op_A$
  and $motif(op_A)$ contains $< m1 : Merging >$,
  then $collapse\_op\_upstream(W, op_A)$

By-Collapse

By-Elimination
24

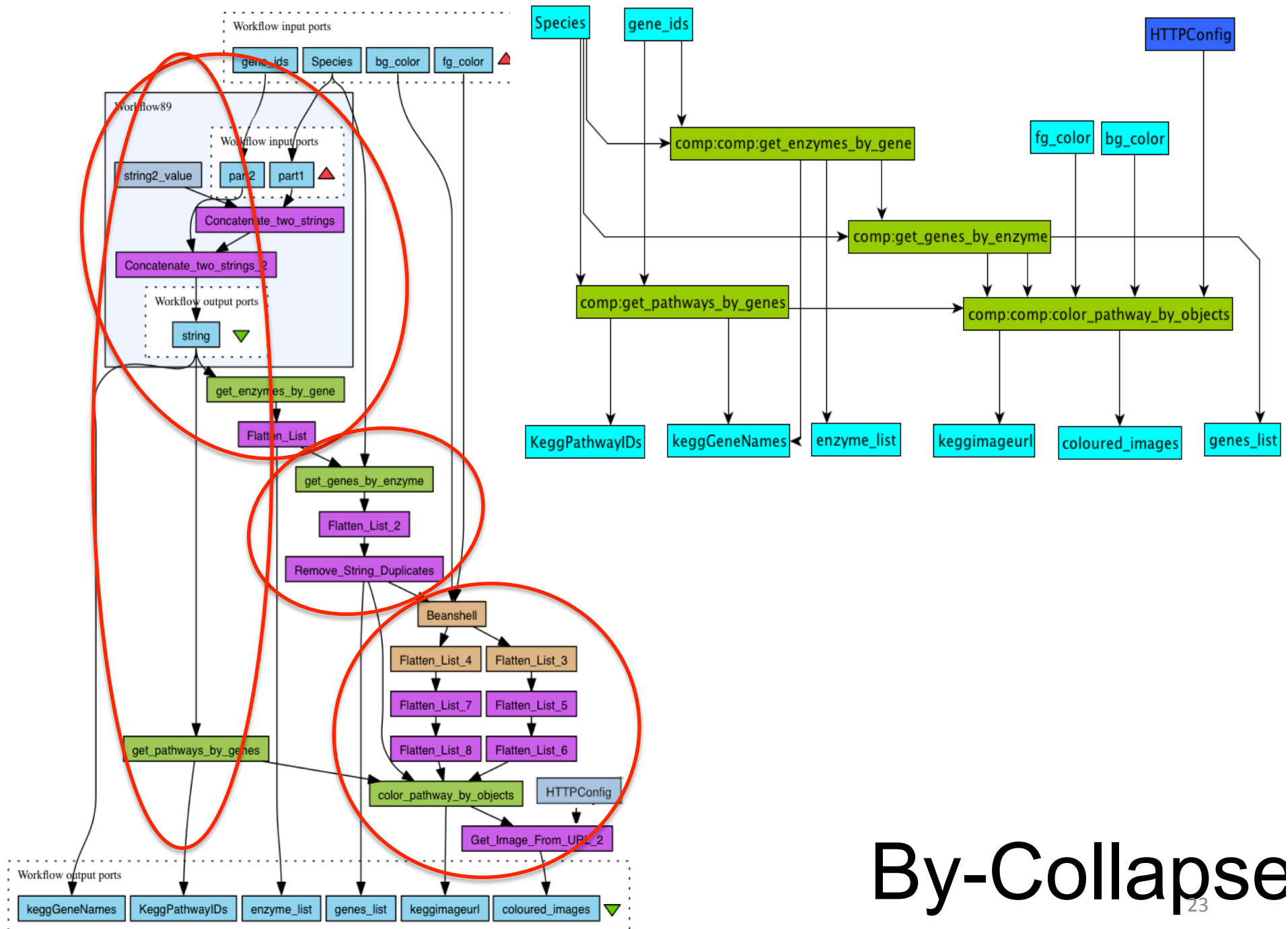# Analysis Data Set

- 30 Workflows from the Taverna system

- Entire dataset & queries accessible from
  http://www.myexperiment.org/packs/467.html

- Manual Annotation using Motif Vocabulary

# Mechanistic Effect of Summarization

|  | By-Elm. | By-Col. |
|---|---|---|
| Avg. Decrease in (Process-Wise)# of Operations | 68% | 63% |
| Avg. Decrease in # of Links (Process-Wise) | 31% | 37% |
| Avg. Decrease in Workflow Depth (Process-Wise) | 62% | 57% |
| Avg. Decrease in Workflow Depth (Data-Wise) | 33% | 57% |

# User Summaries vs. Summary Graphs



"Find pathways in which all the genes in the list are involved. For each pathway draw the pathway diagram. Color all enzyme boxes with colors specified. This workflow still has one problem. The list of colors have to be specified. I would like ideally to only except one background and one foreground color and expand that to a list with length equivalent to the number of enzymes found - just duplicating the specified colors."

| | By-Elm. Precision | By-Elm. Recall | By-Col. Precision | By-Col. Recall |
|---|---|---|---|---|
| Process-Wise | 0.74 | 0.92 | 0.65 | 0.93 |
| Data Wise | 0.14 | 0.55 | 0.33 | 0.43 |

# Highlights

- Research Object model and associated management tools

- Annotations of Workflow Using Motifs

- Methods for Summarizing Workflow and distilling their provenance traces

- Algorithms for Repairing Workflows

# Ongoing Work

- Validation of the workflow summarization

- Querying of Workflow Execution Provenance using summaries.

# References

- P Alper, K Belhajjame, C Goble, P Karagoz. Small Is Beautiful: Summarizing Scientific Workflows Using Semantic Annotations. IEEE International Congress on Big Data, 2013

- Pinar Alper, Khalid Belhajjame, Carole A. Goble, Pinar Karagoz: Enhancing and abstracting scientific workflow provenance for data publishing. EDBT/ICDT Workshops 2013.

- Belhajjame K, Corcho O, Garijo D, et al. Workflow-Centric Research Objects: A First Class Citizen in the Scholarly Discourse. In proceedings of the ESWC2012 Workshop on the Future of Scholarly Communication in the Semantic Web (SePublica2012), Heraklion, Greece, May 2012

- Belhajjame K, Zhao J, Garijo D, et al. The Research Object Suite of Ontologies: Sharing and Exchanging Research Data and Methods on the Open Web, submitted to the Journal of Web Semanics.

- Daniel Garijo, Pinar Alper, Khalid Belhajjame, Óscar Corcho, Yolanda Gil, Carole A. Goble: Common motifs in scientific workflows: An empirical analysis. eScience 2012.

- Zhao J, Gómez-Pérez JM, Belhajjame K, Klyne G, et al. Why workflows break - Understanding and combating decay in Taverna workflows. IEEE eScience 2012.

# Acknowledgement

# Acknowledgement