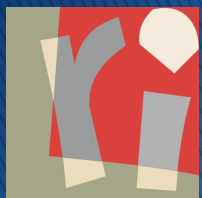# Computational reproducibility in the Life Sciences and Research in Computer Science: round trip

**Sarah Cohen-Boulakia**

Université Paris-Saclay, Laboratoire de Recherche en Informatique

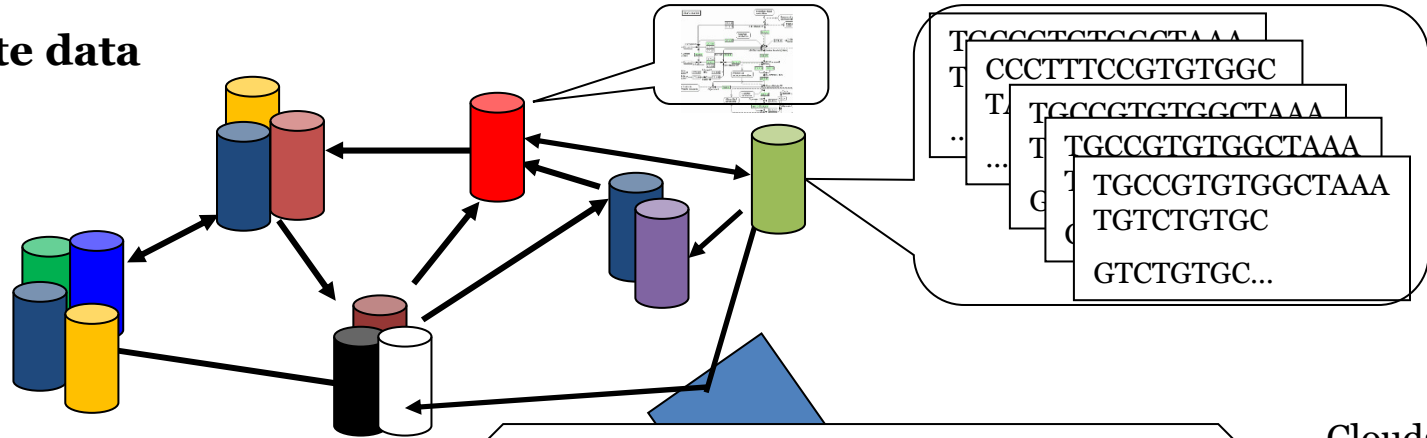CNRS UMR 8623, Université Paris-Saclay, Orsay, France

# Bioinformatics analysis

**Public and private data sources**

Distributed

Heterogeneous

**> 1,500**

TGCCGTCTCCCTAAA
CCCTTTCCGTGTGGC
TGCCGTGTGGCTAAA
TGCCGTGTGGCTAAA
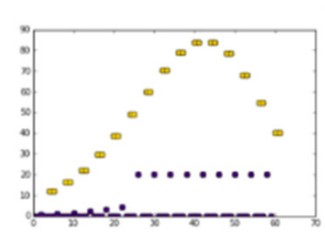TGCCGTGTGGCTAAA
TGTCTGTGC
GTCTGTGC...

How has this plot been generated?
With which input data?
With which tools?
Parameters?

**→ Reproducibility**

Binarization    Water Use Efficiency
Segmentation    **Java**
**Python**       Web services

Clouds

Grids

Clusters

Desktop

**Tools**

Distributed **> 13,000**

Heterogeneous

To be chained

**Biologist's workspace**

MaDICS

# Studies on reproducibility

- Nekrutenko & Taylor, Nature Genetics (2012)
  - 50 papers published in 2011 using the Burrows-Wheeler Aligner for Mapping Illumina reads.
  - 31/50 (62%) provide no information
    - no version of the tool + no parameters used + no exact genomic reference sequence
  - 7/50 (14%) provide all the necessary details
- Alsheikh-Ali et al, PLoS one (2011)
  - 10 papers in the top-50 IF journals → 500 papers (publishers)
    - 149 (30%) were not subject to any data availability policy (0% made their data available)
    - Of the remaining 351 papers
      - 208 papers (59%) did not adhere to the data availability instructions
      - 143 make a statement of *willingness* to share
      - 47 papers (9%) deposited full primary raw data online

# Context, Challenges

*Computational reproducibility crisis*

## Increasing number of irreproducible results

- Even published in high IF venues
- Not (always) deliberately
- Computational irreproducibility increases
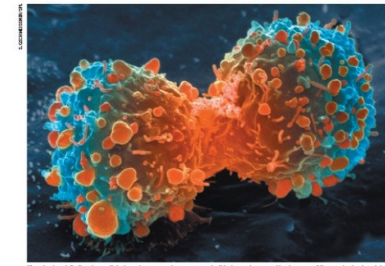
## Various scientific domains

- Consequences may be huge (preclinical studies...)

## Major challenge

- The cost of irreproducible preclinical studies have been evaluated to >$10 Billions per year (USA)

## Becoming mandatory

- NSF projects, editors, ANR...



**Must try harder**
*Too many sloppy mistakes are creeping into scientific papers. at the data — and at themselves.*

**Error prone**
*Biologists must realize the pitfalls massive amounts of data.*

**If a job is worth doing, it is worth doing twice**
*Researchers and funding agencies need to put a premium on ensuring that results are reproducible, argues Jonathan F. Russell.*

The case for open computer programs

**Six red flags for suspect work**
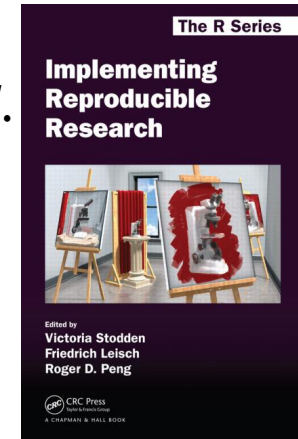C. Glenn Begley explains how to recognize the preclinical papers in which the data won't stand up.

Know when your numbers are significant

Raise standards for preclinical cancer research
C. Glenn Begley and Lee M. Ellis propose how methods, publications and incentives must change if patients are to benefit.

47/53 "landmark" publications could not be replicated

[Begley, Ellis Nature, 483, 2012]

# Reproducibility

V. Stodden *et al.*



## *Empirical reproducibility*
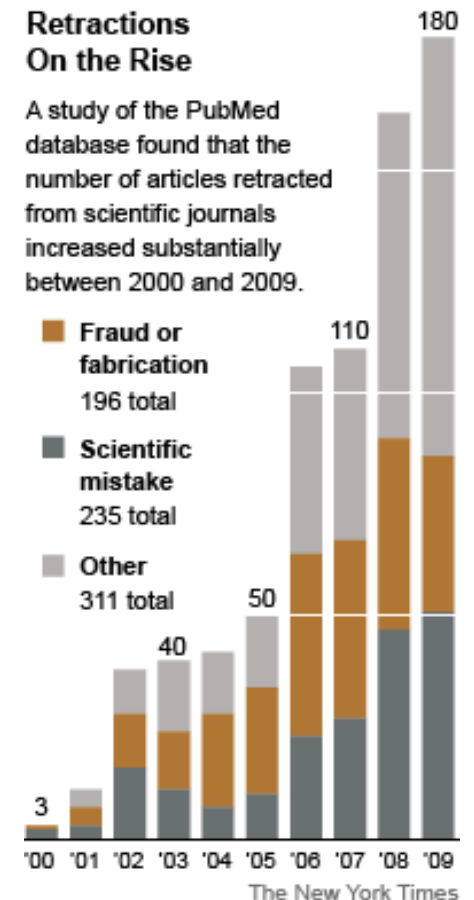
- ◦ detailed information about non-computational empirical scientific experiments and observations
- ◦ In practice this is enabled by making data freely available, as well as details of how the data was collected.

## *Statistical reproducibility*

- ◦ detailed information about the choice of statistical tests, model parameters, threshold values, etc.
- ◦ This relates to pre-registration of study design to prevent p-value hacking and other manipulations.

## *Computational reproducibility*

- ◦ detailed information about code, software, hardware and implementation details
  - → Goal: document how data has been produced



Retractions
On the Rise

A study of the PubMed database found that the number of articles retracted from scientific journals increased substantially between 2000 and 2009.

- Fraud or fabrication 196 total
- Scientific mistake 235 total
- Other 311 total

The New York Times

# Scripts and reproducibility?
## Good practices

Providing scripts is an excellent first step

+ Using git/github for versioning, collaborative development

But scripts do not allow to

Distinguish between steps of the analysis
- ◦ piece of codes, methods/functions

… and execution of the analysis
- ◦ data sets used as inputs and then produced

Emphasize the major steps of the analysis

Provide solution for data management
- ◦ Naming convention for produced files, storage…

→ Scripts are difficult to share, exchange and reuse (repurpose)

# Outline

Context

## Systems and tools to enhance reproducibility
- ◦ Scientific workflow systems
- ◦ Companion tools

Lessons learnt on using such systems and tools
- ◦ Reprohackathons
- ◦ Levels of reproducibility with scientific workflows
- ◦ Reproducibility-friendly features

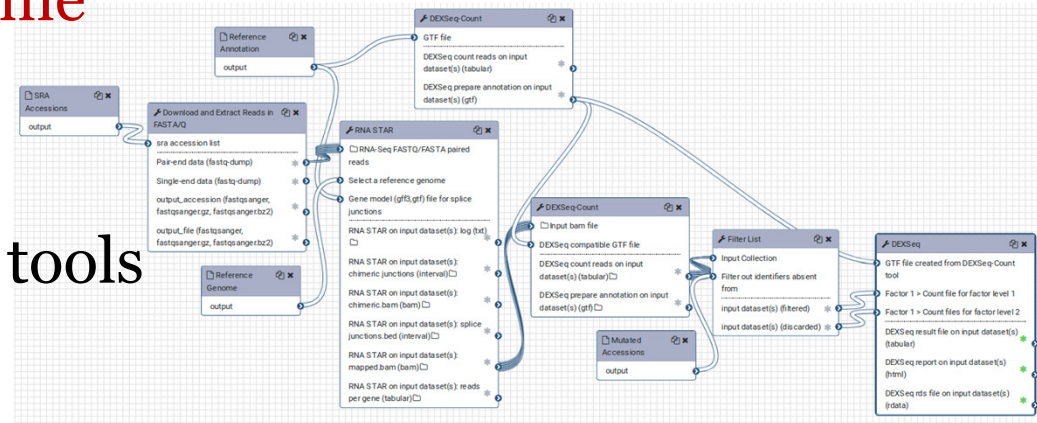Open Computer Science research problems

Conclusion

# Scientific workflow systems

SWFS = "Data analysis pipeline"

Data flow driven

Encapsulation of scripts

WF specification: connected tools
*steps of the analysis*



WF execution: data consumed/produced
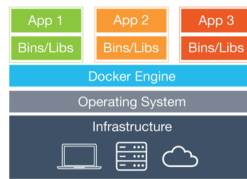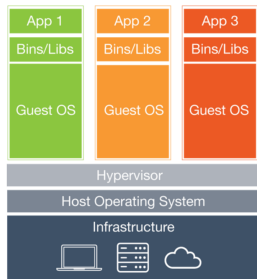
Provenance modules
*data management*

SWFS scheduling, logging,
May be equipped with GUI
Galaxy, NextFlow, SnakeMake…

# Capturing the programming environment

Ensuring your workflow has everything it needs to run
Libraries, dependencies...
Virtual machines capture the programming environment
Container solutions

○ package an application
  - with all of its dependencies
  - into a standardized unit for software development
include the application and its dependencies
○ but share the kernel with other containers
○ They
  - are not tied to any specific infrastructure;
  - run on any computer, on any infrastructure and in any cloud

Lighter solution than classical VM

➜ BioContainers: a registry of containers!

# Outline

Context

Systems and tools to enhance reproducibility

- Scientific workflow systems
- Companion tools

## Lessons learnt on using such systems and tools

- Reprohackathons
- Levels of reproducibility with scientific workflows
- Reproducibility-friendly features

Open Computer Science research problems

Conclusion

# Our new concept: ReproHackathon

## ReproHackathon

- A hackathon where
  - Given a scientific publication + input data (+ possibly contacts with authors)
  - Several (groups of) developers reimplement the methods to try to get the same result on the Cloud@IFB
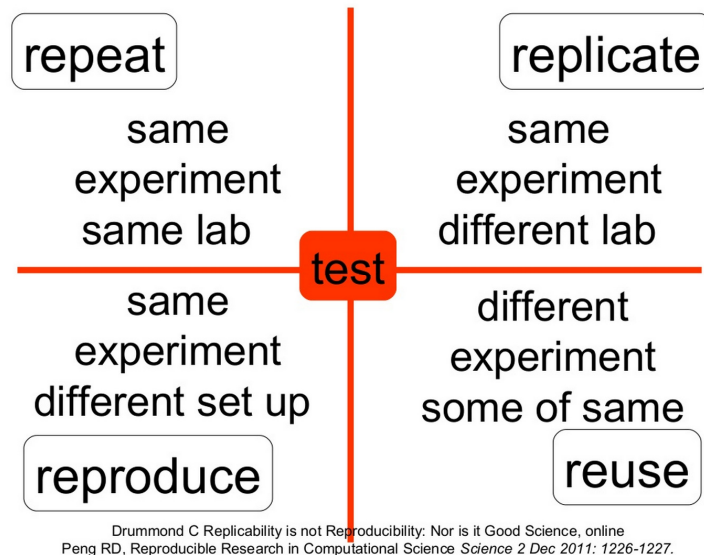- Aim: Ability of current workflow systems and companion tools to reproduce a scientific result





First edition
- RNA-Seq data from patients with uveal melanoma: genes involved
- Divergent published results…
- 25 participants (IGRoussy, Curie, Pasteur, Saclay,…)

Reprohackathon 2 Lyon, July 2018
Phylogenetics

Reprohackathon 3 Montpellier Nov 2019
Plant phenotyping

# Levels of computational reproducibility



Drummond C Replicability is not Reproducibility: Nor is it Good Science, online
Peng RD, Reproducible Research in Computational Science *Science 2 Dec 2011: 1226-1227.*

## 3 ingredients

Workflow Specification
Chained Tools
Workflow Execution
Input data and parameters
Environment
OS/librairies …

## Repeat

- ◦ *Redo*: exact same context
- ◦ Same workflow, execution setting, environement
- ◦ Identical *output*
→Aim = proof for reviewers ☺

## Replicate

- ◦ Variation allowed in the workflows, execution setting, environement
- ◦ Similar *output*
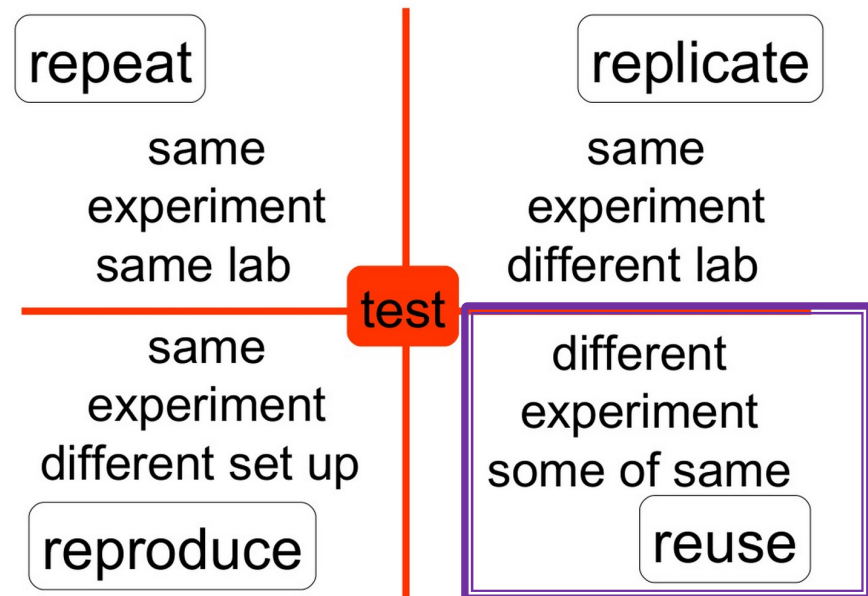→ Aim = robustness

# A continuum of possibilities

Reproduce

- Same *scientific result*
- But the means used may be changed
- Different workflows, execution setting, environment
- Different output but in accordance with the result

Reuse

- Different scientific result
- Use of tools/... designed in another context



repeat — same experiment same lab

replicate — same experiment different lab

test

reproduce — same experiment different set up

reuse — different experiment some of same

Drummond C Replicability is not Reproducibility: Nor is it Good Science, online
Peng RD, Reproducible Research in Computational Science *Science 2 Dec 2011: 1226-1227.*

# Reproducibility-friendly features

6 Systems: Galaxy, Nextflow, SnakeMake, VisTrails, OpenAlea, Taverna

Future Generation Computer Systems
Volume 75, October 2017, Pages 284-298

Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities
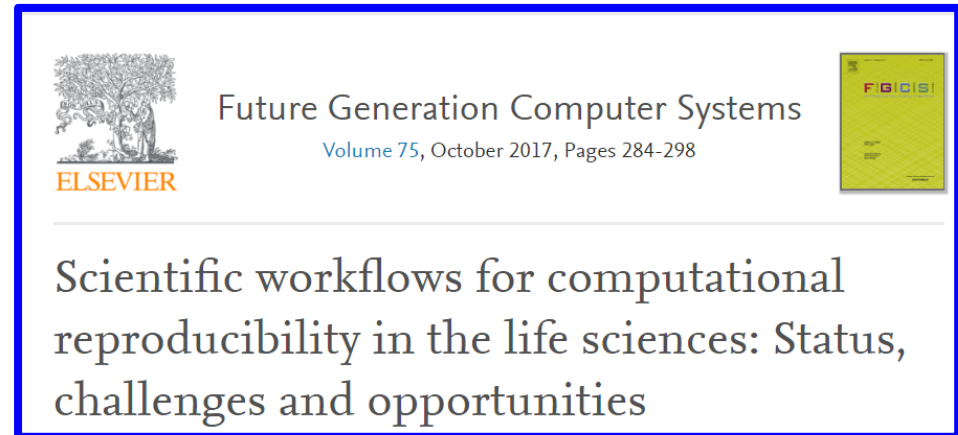
## Specification

Language (XML, Python…)

Interoperability (CWL…)

Description of steps
- Remote services
- Command line
- Access to source code

Modularity (nested workflows?)

Annotation (tags, ontologies, myexperiment…)

## Execution

Language and standard (PROV…,) → repeat … reuse

Presentation (interactivity with the results/provenance, notebooks) → replicate … reuse

Annotations → reuse

## Environment

Ability to run workflows within a given environment

Virtual machines
  ◦ VMWare, KVM, VirtualBox, Vagran,…

Lighter solutions (containers)
  ◦ Docker, Rocket, OpenVZ, LXC, Conda

Capturing the command-line history, input/output, specification: CDE, ReproZip

# Outline

Context

Systems and tools to enhance reproducibility

Lessons learnt on using such systems and tools

Open Computer Science research problems
Research topics
Focus on 3 graph-oriented kind of problems

Conclusion

Sarah Cohen-Boulakia, Université Paris-Saclay

# Developing workflows

<span style="color:red">Bridge the gap between scripts and workflows</span>

Supporting <span style="color:blue">several programming languages</span> in the same environment of development

<span style="color:blue">Tests</span> in workflows
- Unit tests, integration tests...
- Providing samples may be an issue (privacy...)

Workflow Maintenance: set of compatible libraries?
- Docker (containers), VM allows to freeze the environment

→ <span style="color:blue">Need to liquefy!</span>
- Given a program P that can be repeated in an environment E... ... Find an environment E' (E' uses more recent versions of libraries than E) where P still *works*

# Discovering workflows [Reuse]

Query languages for repositories?

*Given a workflow – find similar workflows*

Detecting patterns within workflows

Indexing workflows

Reconstruct their histories

## Core of the problem:

Workflow similarity

State-of-the-art [SCB+14]

Based on the graph structures or annotations (ontologies)

Need to design hybrid and efficient solutions

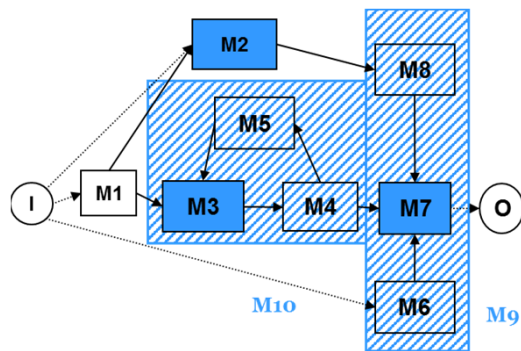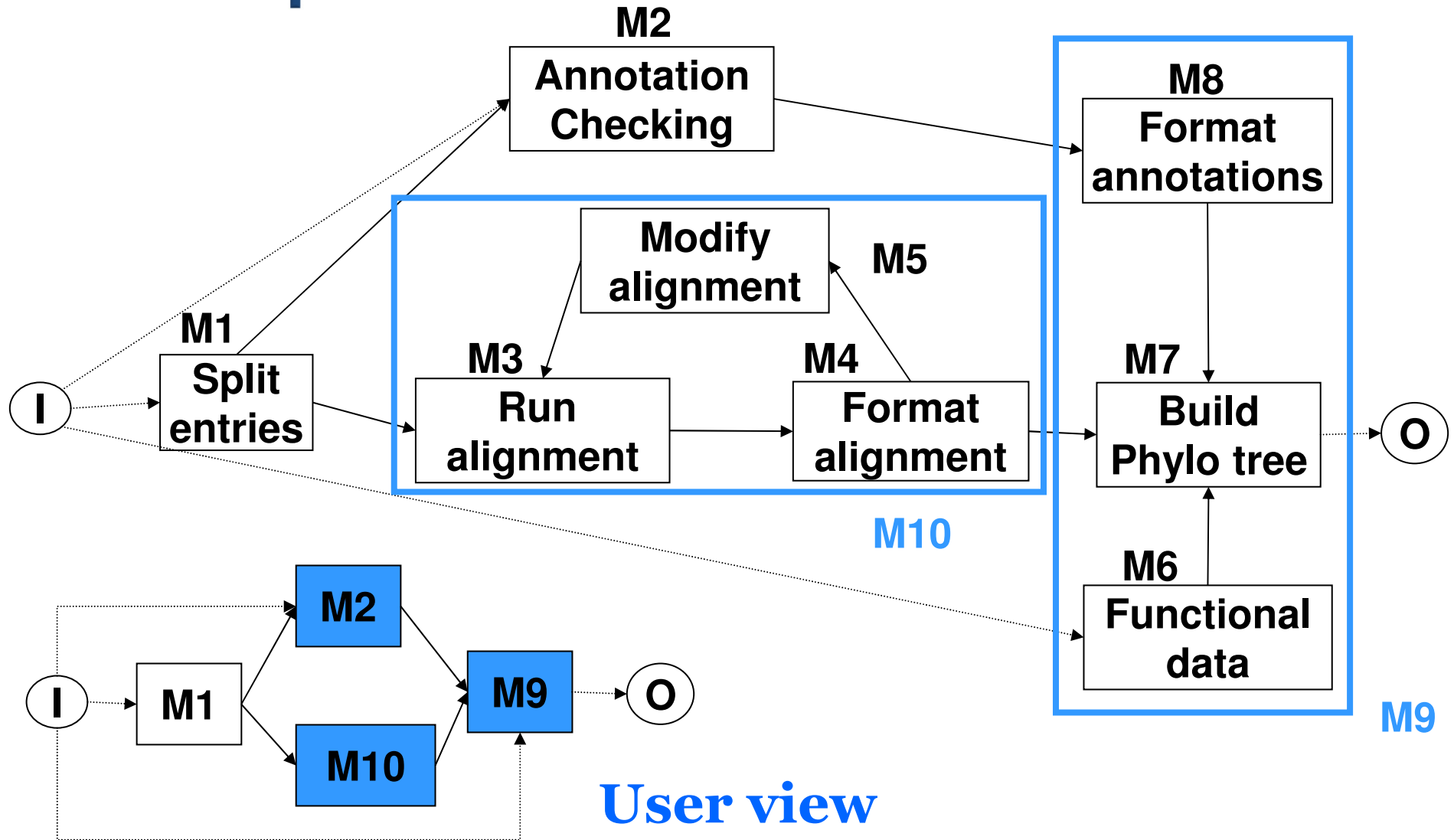NB : Reusing (and searching for)

Notebooks is another open point

# Simplifying workflows [Reuse]

Designing more coarse-grained workflows

◦ Automatic Design of subworkflows (graph-based)

◦ Abstraction of provenance traces

◦ Summarization (Web Semantics)

Refactoring workflows

◦ Remove redundancies in workflows

◦ Rewritting, Anti-patterns

# Outline

Context

Systems and tools to enhance reproducibility

Lessons learnt on using such systems and tools

Open Computer Science research problems
Research topics
Focus on 3 graph-oriented kinds of problems

Conclusion

# Problem 1: Abstracting workflows (composition)

Reducing the complexity of workflows
making them easier to share

# Composite modules

# Composite modules



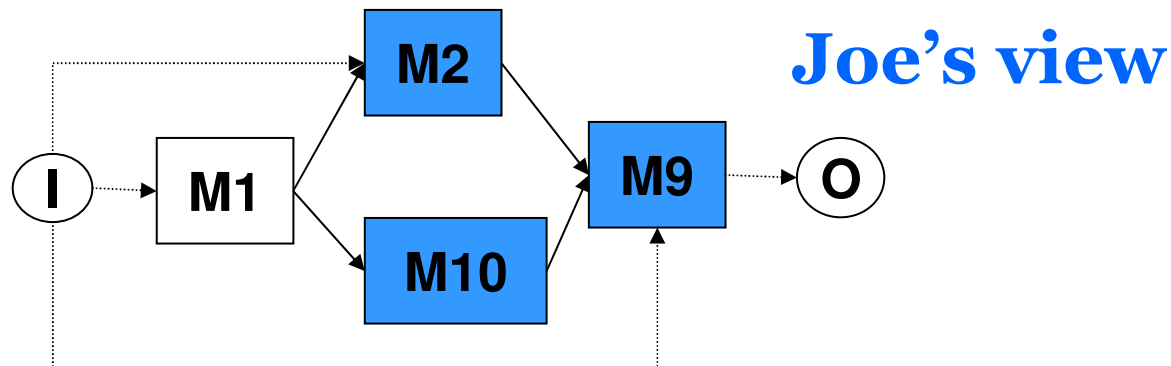Composition simplifies provenance

(and the wf specification)

# Relevant user view
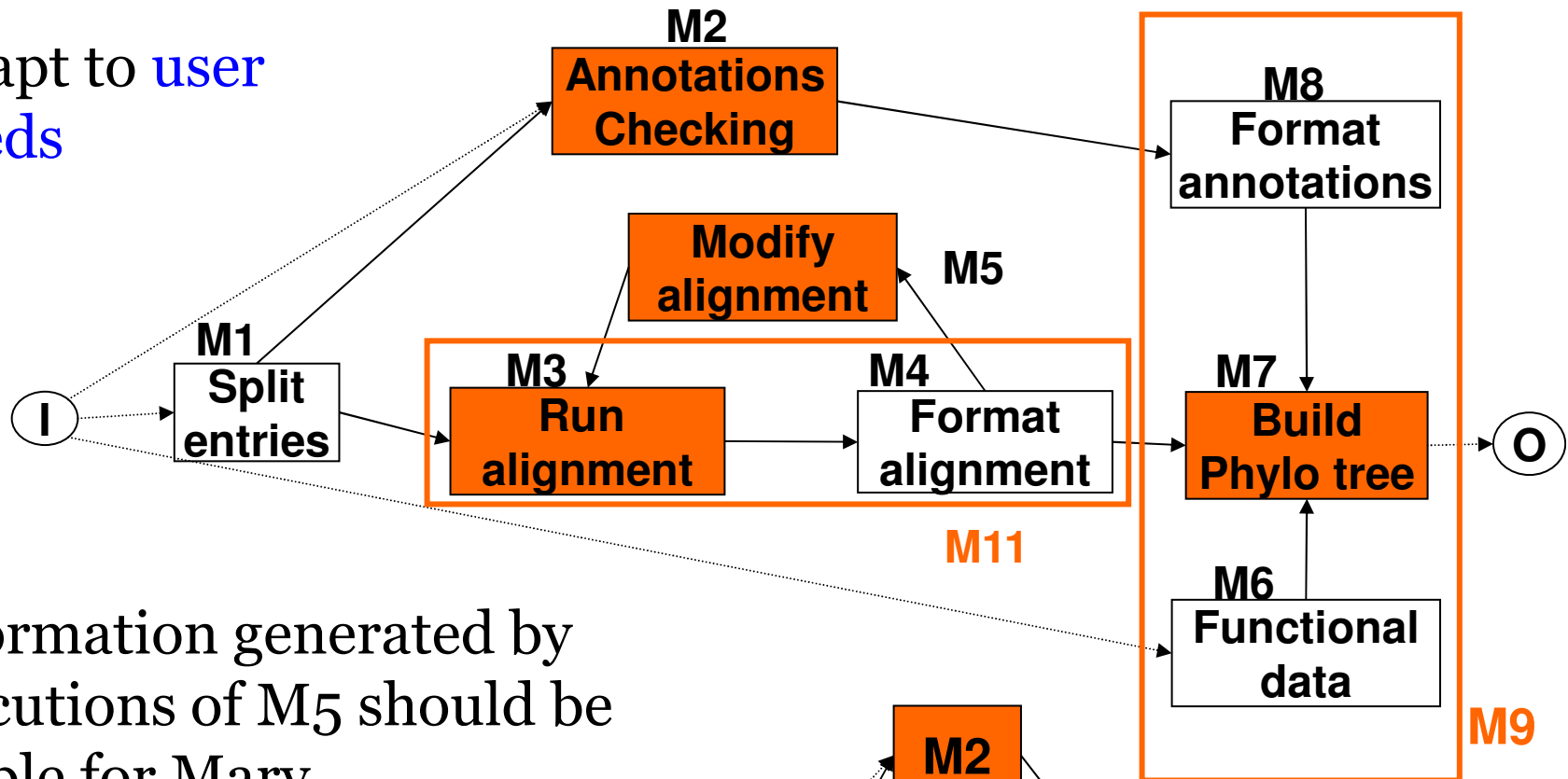
Modules Joe
considers
relevant

Each composite module takes the
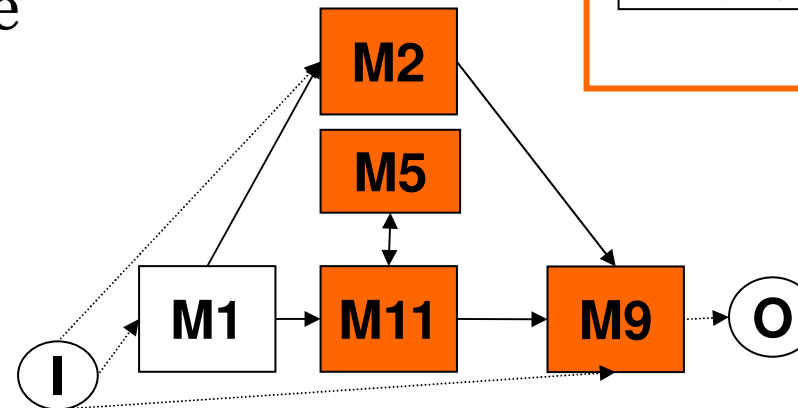meaning of the relevant module it
contains



**Joe's view**
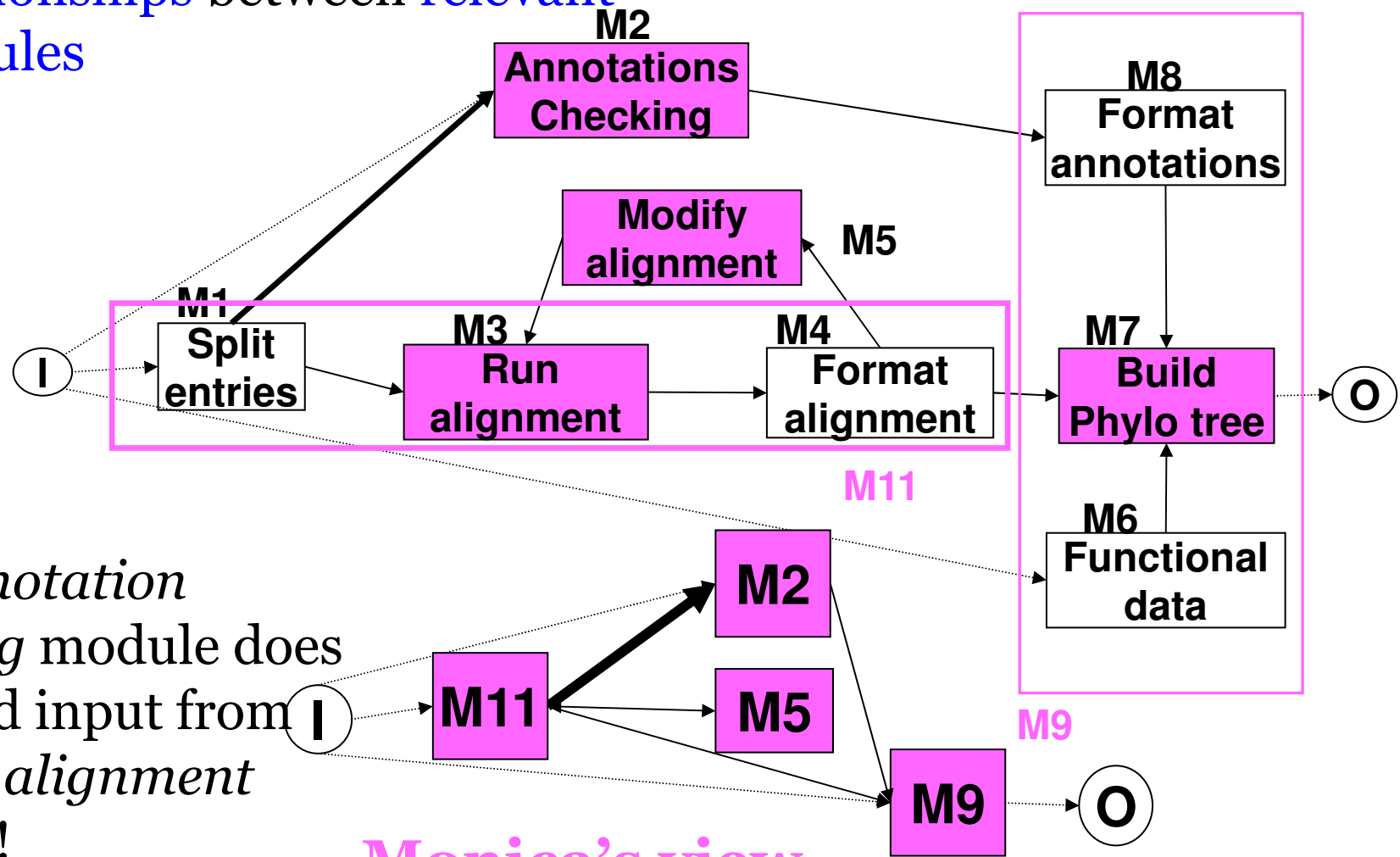
# User views may differ

Adapt to user needs

Information generated by executions of M5 should be visible for Mary

**Mary's view**

# Grouping may be error-prone!

Grouping should preserve the relationships between relevant modules



**M2**
Annotations Checking

**M8** Format annotations

Modify alignment

**M5**

**M1** Split entries

**M3** Run alignment

**M4** Format alignment

**M7** Build Phylo tree

**M11**

**M6** Functional data

**M9**

The *annotation checking* module does not need input from the *run alignment* module!

**M2**

**M11**  **M5**

**M9**

**Monica's view**

# Hope and next challenges

▸ Hope
- *ZOOM* provides a polynomial-time algorithm to automatically construct user views
  - which preserve the dataflow (no missing path and no new path between two relevant tasks)
  - and produces a minimal user view

In collaboration with  Penn
UNIVERSITY of PENNSYLVANIA

▸ Next challenges

Repairing user views badly designed
Providing such functionalities in real systems
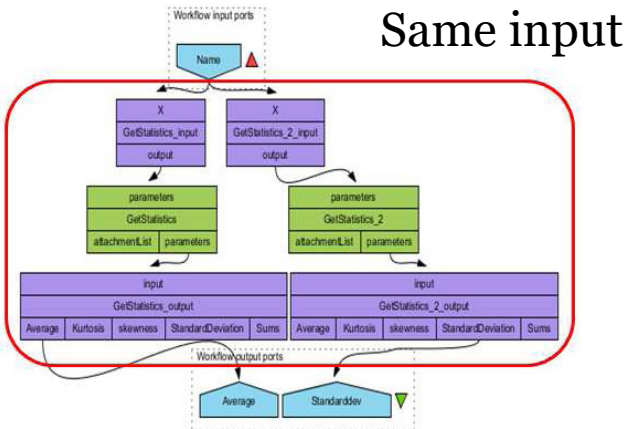Using it on provenance information in provenance systems

# Problem 2: rewritting workflows

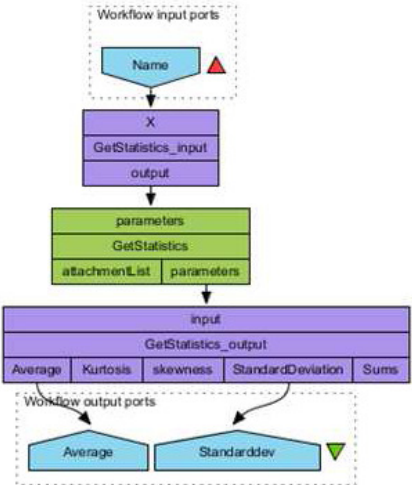Repairing workflows,
making them easier to share

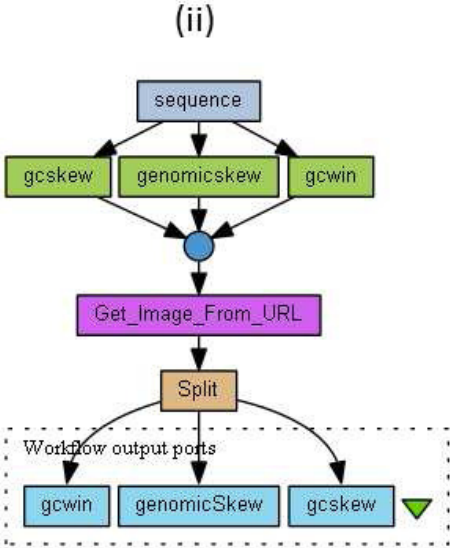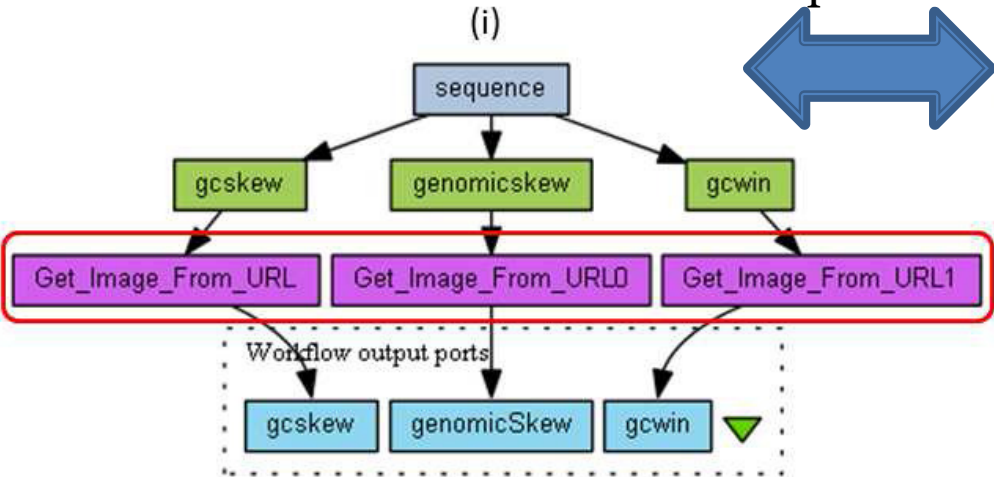# Redundancy in Workflows

3 processors duplicated!

No redundancy
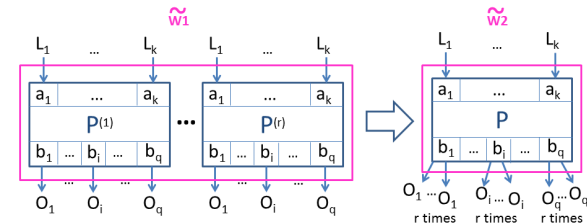
Same input

Equivalent

Equivalent

# Rewriting workflows

▸ Assumptions before merging several copies of a processor

  ◦ Only copies with the exact same code
  ◦ Only copies that do not depend on each other
  ◦ Only deterministic processors (same input → same output)

▸ Need to understand the semantics of the system

  ◦ Determining the anti-patterns and designing their corresponding rewritings

# Hope and next challenges

- **Hope**
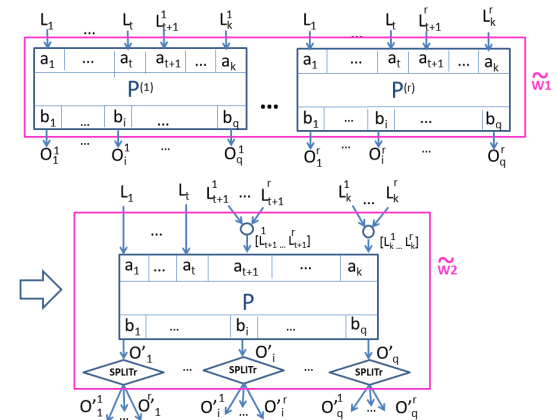  - *DistillFlow* detects anti-patterns and rewrites Taverna workflows

- **Next Challenges**
  - Larger sets of anti-patterns
  - Rewritting on in-use systems (Galaxy, NextFlow, SnakeMake)

In collaboration with



$L_i$ can be one single value or a list of values

Processor P applies cross product to values on ports $a_1$ to $a_t$ and dot product to values on ports $a_{t+1}$ to $a_k$
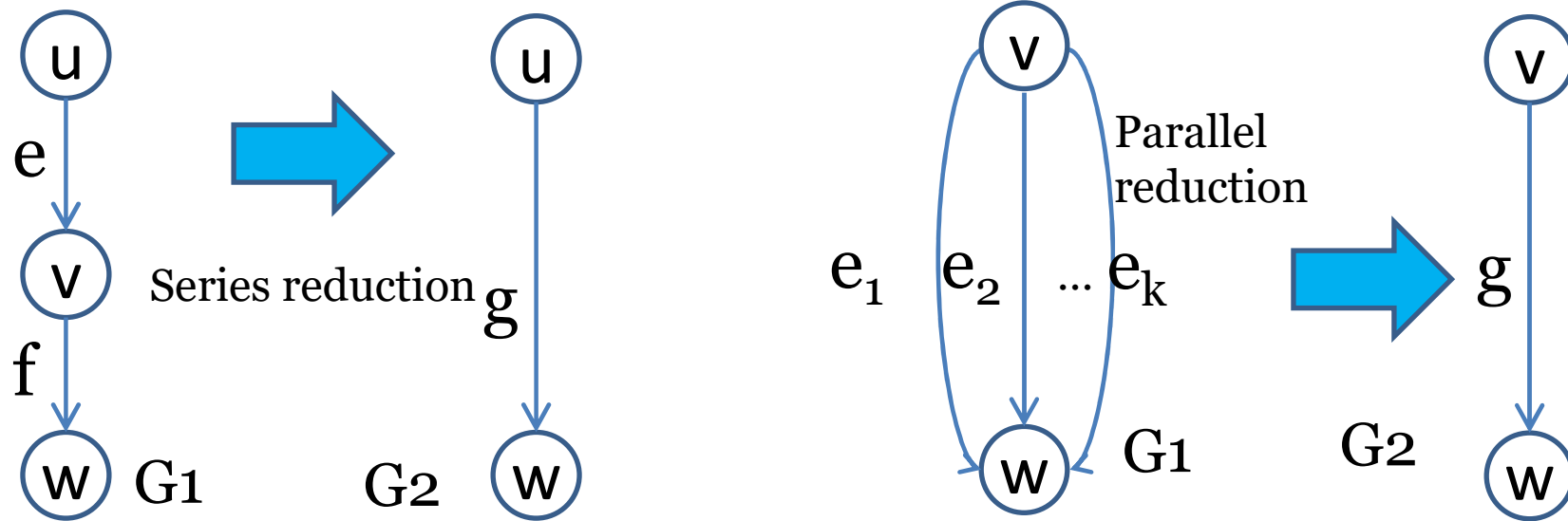
# Problem 3: exploiting specific graph structures

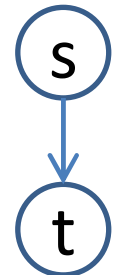Comparing graphs is a difficult problem that can be simplier on specific graph structures

# Definition of SP-graphs

**G is SP iff MaxRed(G) = BSP**

▸ **MaxRed(G):** iteratively performs series and parallel reductions on a given graph G



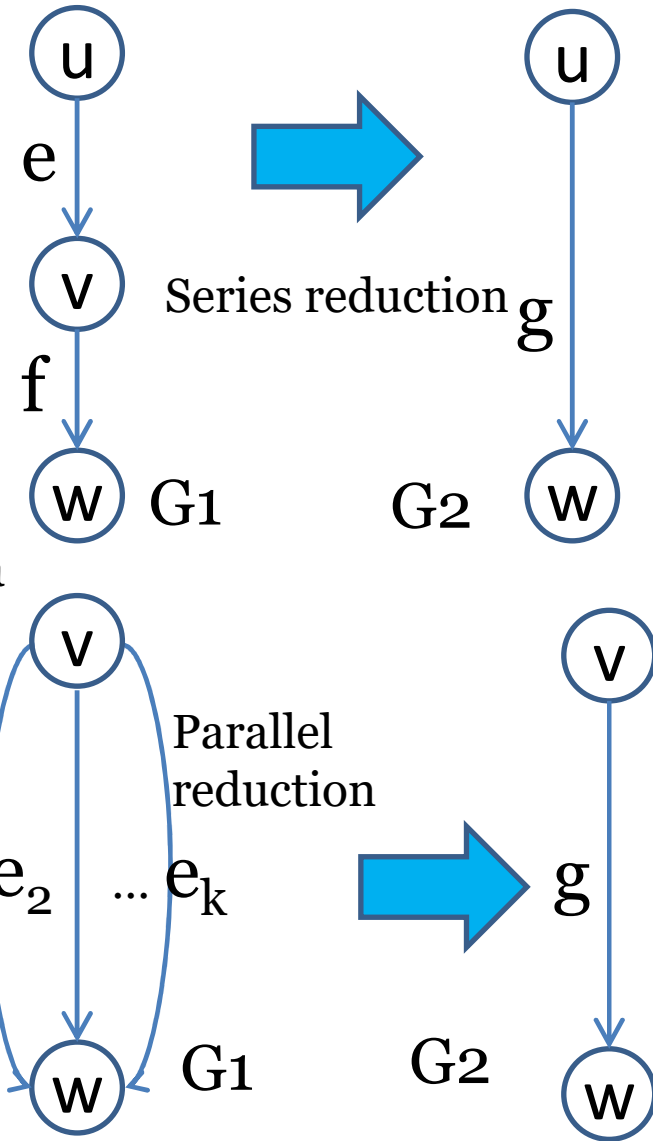Series reduction

Parallel reduction

▸ **BSP**: Basic Series-Parallel

# Is it Series-Parallel?

YES!

G is SP iff
$$MaxRed(G) = BSP$$

▸ MaxRed(G): iteratively performs series and parallel reductions on a given graph G

▸ BSP: Basic Series-Parallel

$(G_o)$

$(G_o)$
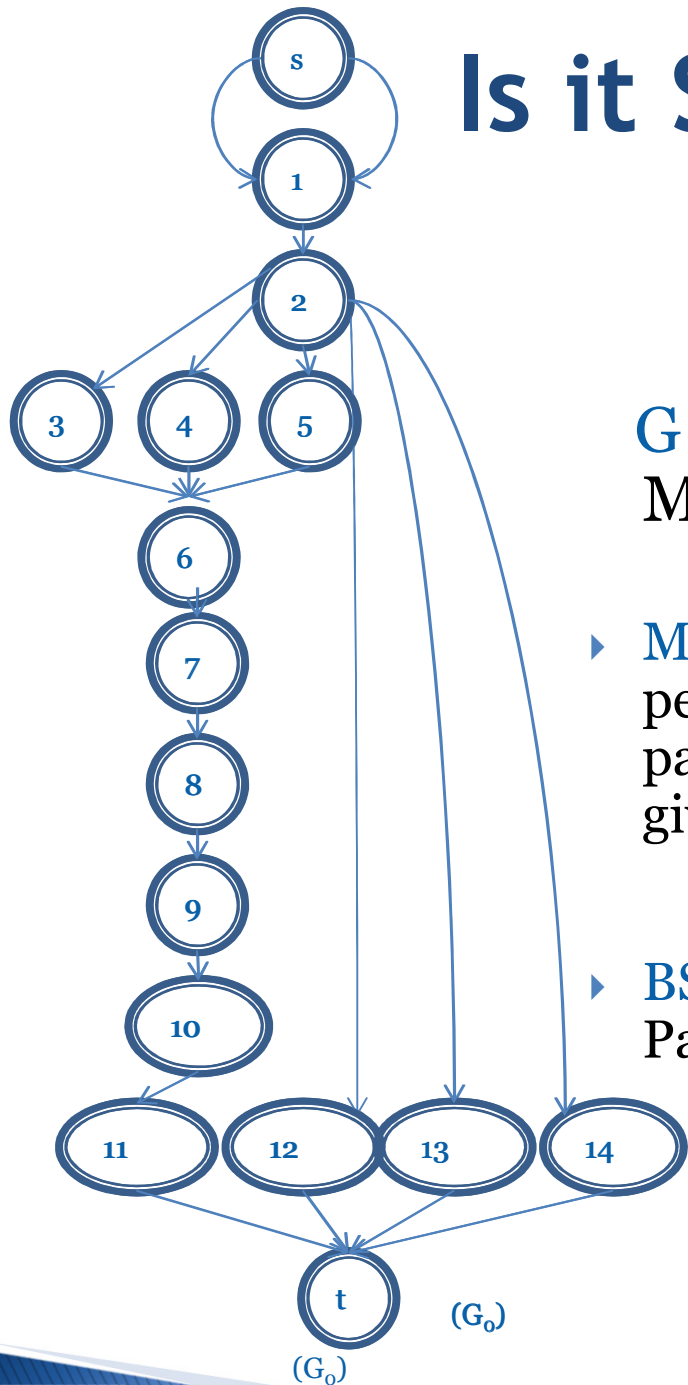
Series reduction

Parallel reduction

$e_1 \quad e_2 \quad ... \quad e_k$

$G_1$

$G_2$

# Is it Series-Parallel?

NO!
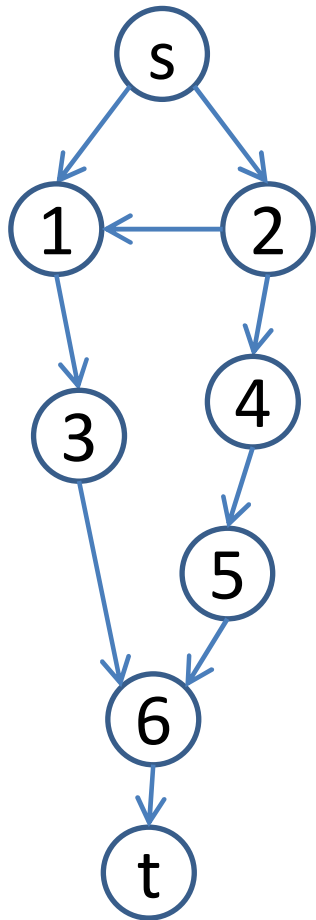
G is SP iff
MaxRed(G) = BSP

▸ MaxRed(G): iteratively performs series and parallel reductions on a given graph G

▸ BSP: Basic Series-Parallel

Series reduction

$G_1$ → $G_2$
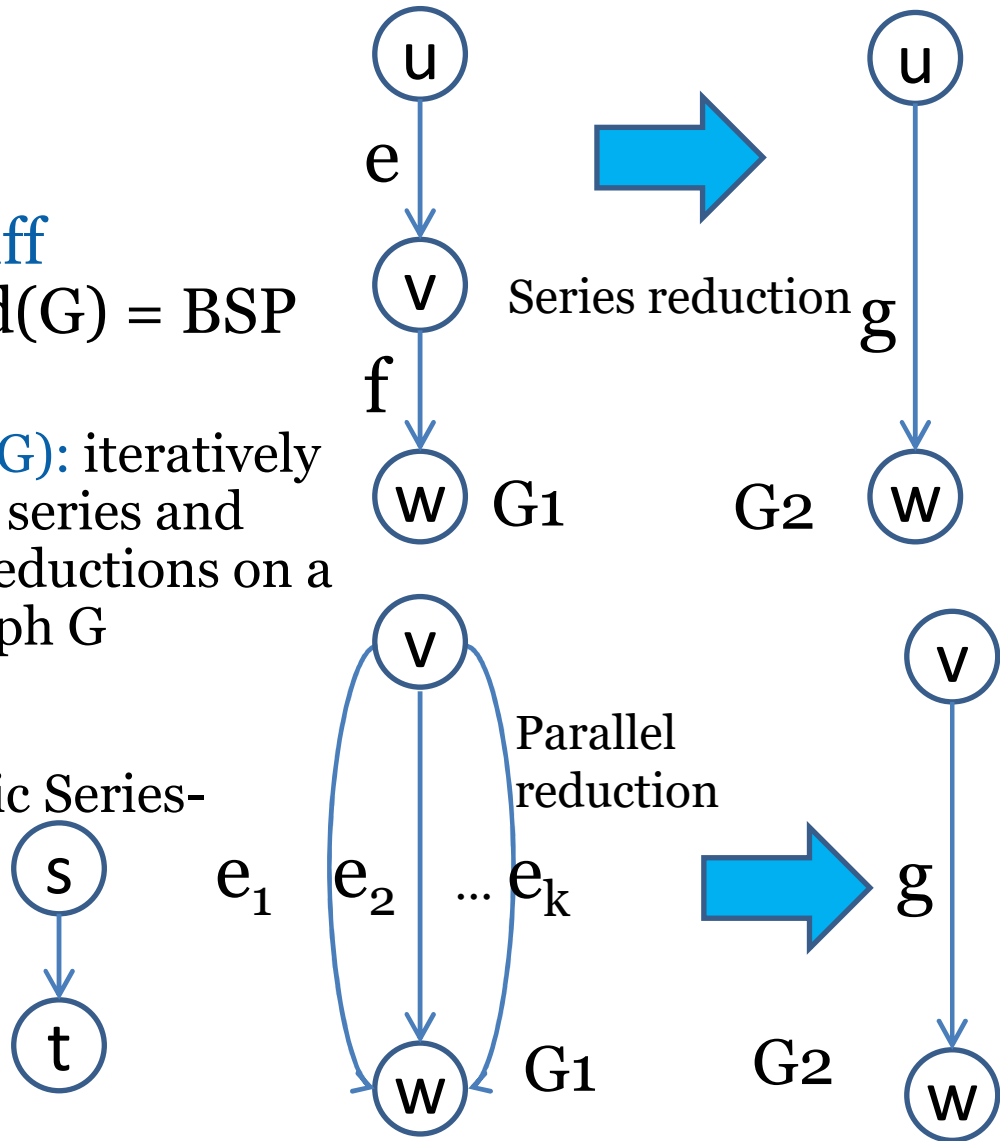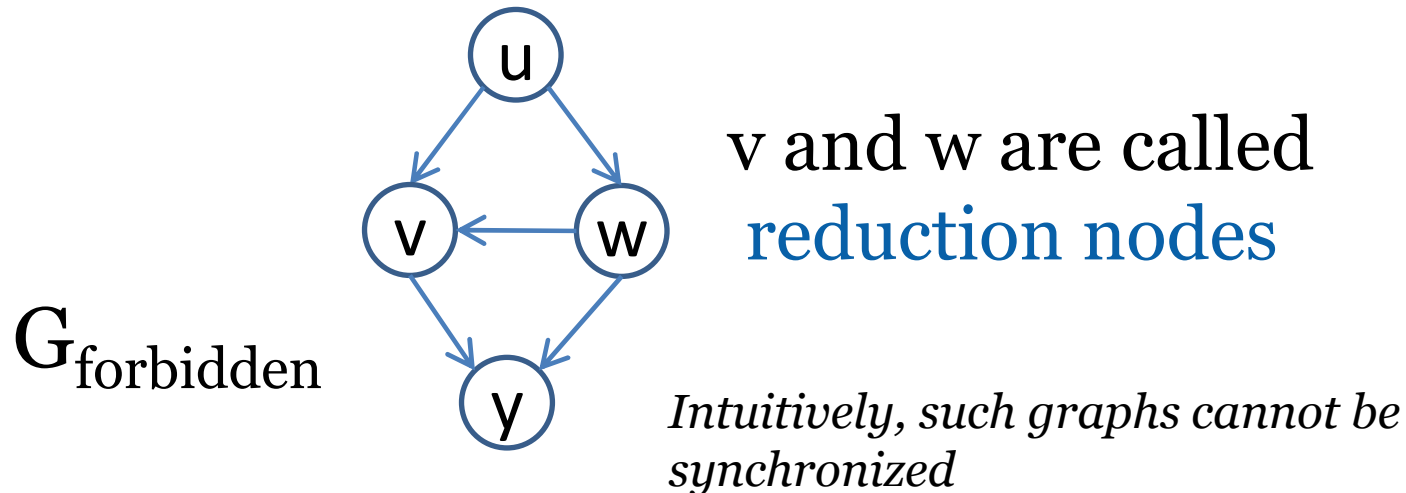
Parallel reduction

$e_1$ $e_2$ ... $e_k$ → $g$

$G_1$ → $G_2$

# Another definition (Non SP-graphs)

G is non-SP iff MaxRed(G) contains $G_{forbidden}$

$G_{forbidden}$



v and w are called reduction nodes

*Intuitively, such graphs cannot be synchronized*

Subgraph isomorphism is polynomial for SP graphs

# Some hope and next Challenges

▸ Hope

- ◦ In *PDiffView* SP-graphs have been used to solve the problem of providing a polynomial algorithm to compare two executions of the same workflow

In collaboration
with

▸ Next challenges

- ◦ Guiding users in designing (close to) SP-graph structures
- ◦ Extending SP structures to cover more expressive workflow structures while making sure that polynomial-time algorithms can be found to compare workflows

→ Impact on querying workflow repositories, reconstructing workflow history, …

# Conclusion

Many scientific results are not computationally reproducible
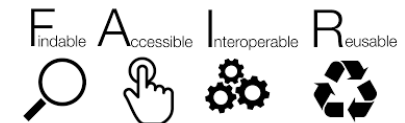
Providing scripts is an excellent start

Scientific workflows are increasingly mature solutions

- Tracking the exact connected tools used
- Track the exact data used, produced and tool parameters setting
  - →Provenance modules
- Coarse-grain version of the analysis to better capture the analysis steps

Several open challenges are directly related to improvement in research in computer science (graphs, algorithmics…)

Workflows play key role to produce FAIR data
FAIR metrics for workflows have to be defined too!

Findable  Accessible  Interoperable  Reusable

Join us!

cohen@lri.fr

**MaDICS CNRS GDR**

Sarah Cohen-Boulakia, Université Paris-Saclay