# Reproducibility and reuse in data-driven sciences: from provenance to summaries

## Franck Michel
*Univ. Côte d'Azur, CNRS, Inria, I3S*

## Alban Gaignard
*Institut du thorax, CNRS, Univ. de Nantes*

# Repeat > Replicate > Reproduce > Reuse

*S. Cohen-Boulakia, K. Belhajjame, O. Collin, J. Chopard, C. Froidevaux, A. Gaignard, K. Hinsen, P. Larmande, Y. Le Bras, F. Lemoine, F. Mareuil, H. Ménager, C. Pradal, C. Blanchet,* **Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities**, *Future Generation Computer Systems, Volume 75, 2017,* https://doi.org/10.1016/j.future.2017.01.012 *.*

# Repeat > Replicate > Reproduce > Reuse

Same experiment

Same setup

Same lab

*S. Cohen-Boulakia, K. Belhajjame, O. Collin, J. Chopard, C. Froidevaux, A. Gaignard, K. Hinsen, P. Larmande, Y. Le Bras, F. Lemoine, F. Mareuil, H. Ménager, C. Pradal, C. Blanchet, **Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities**, Future Generation Computer Systems, Volume 75, 2017, https://doi.org/10.1016/j.future.2017.01.012 .*

# Repeat > Replicate > Reproduce > Reuse

Same experiment     Same experiment

Same setup          Same setup

Same lab            ~~Same lab~~

# Repeat > Replicate > Reproduce > Reuse

| Same experiment | Same experiment | Same experiment |
| Same setup | Same setup | ~~Same setup~~ |
| Same lab | ~~Same lab~~ | ~~Same lab~~ |

# Repeat > Replicate > Reproduce > Reuse

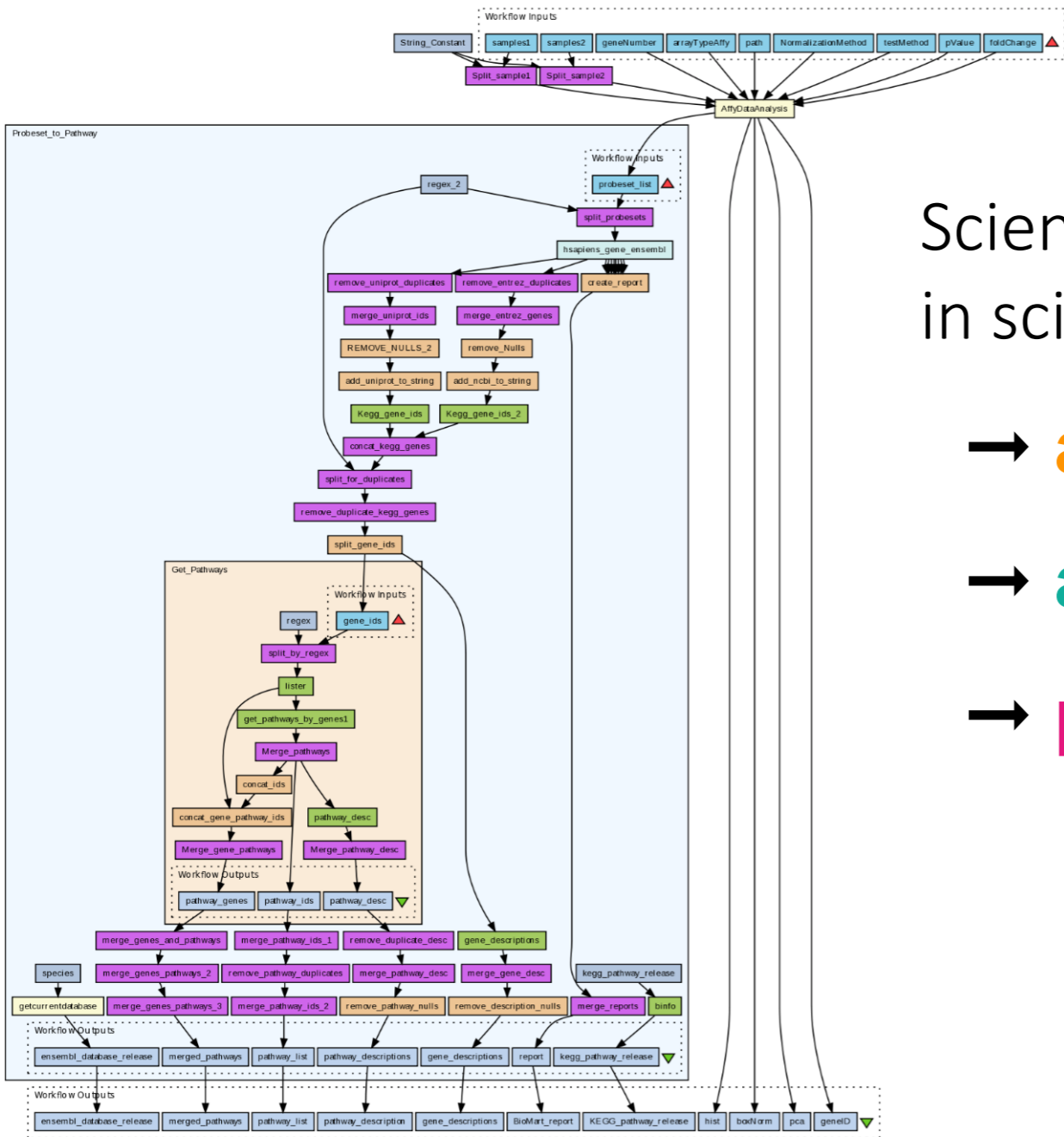| | | |
|---|---|---|
| Same experiment | Same experiment | Same experiment |
| Same setup | Same setup | ~~Same setup~~ |
| Same lab | ~~Same lab~~ | ~~Same lab~~ |

new ideas,
new experiment,
some commonalities

*S. Cohen-Boulakia, K. Belhajjame, O. Collin, J. Chopard, C. Froidevaux, A. Gaignard, K. Hinsen, P. Larmande, Y. Le Bras, F. Lemoine, F. Mareuil, H. Ménager, C. Pradal, C. Blanchet,* **Scientific workflows for computational reproducibility in the life sciences: Status, challenges and opportunities,** *Future Generation Computer Systems, Volume 75, 2017, https://doi.org/10.1016/j.future.2017.01.012 .*

# Scientific <span style="color:red">workflows</span> to the rescue
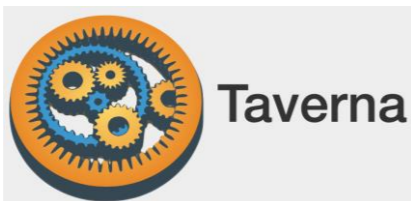
# What is a workflow ?

« a systematic way of **describing the methods** needed and provide the **interface** between **domain specialists** and **computing infrastructures**. »

*Malcolm Atkinson, Sandra Gesing, Johan Montagnat, Ian Taylor. **Scientific workflows: Past, present and future**. Future Generation Computer Systems, Elsevier, 2017, 75, pp.216 - 227. <10.1016/j.future.2017.05.041>*

Scientific workflows to enhance **trust** in scientific results:

➔ **automation** of data analysis (at scale)

➔ **abstraction** (describe/share methods)

➔ **provenance** (~tracability, trust, transparency)

# Provenance

## Definition in Computer Science

« Provenance information describes the **origins** and the **history of data in its life cycle**. »

« Today, (…) data is constantly being created, copied, moved around, and combined indiscriminately. Because information sources (…) vary widely in terms of quality, it is essential to provide **provenance and other context information** which can **help end users judge** whether query results are **trustworthy**. »

*James Cheney, Laura Chiticariu, and Wang-Chiew Tan. 2009. **Provenance in Databases: Why, How, and Where**. Found. Trends databases 1, 4 (April 2009), 379-474. DOI=http://dx.doi.org/10.1561/1900000006*

# Representing provenance

**W3C**

## PROV-O: The PROV Ontology

### W3C Recommendation 30 April 2013

**This version:**
  http://www.w3.org/TR/2013/REC-prov-o-20130430/
**Latest published version:**
  http://www.w3.org/TR/prov-o/
**Implementation report:**
  http://www.w3.org/TR/2013/NOTE-prov-implementations-20130430/
**Previous version:**
  http://www.w3.org/TR/2013/PR-prov-o-20130312/
**Editors:**
  Timothy Lebo, Rensselaer Polytechnic Institute, USA
  Satya Sahoo, Case Western Reserve University, USA
  Deborah McGuinness, Rensselaer Polytechnic Institute, USA
**Contributors:**
  (In alphabetical order)
  Khalid Belhajjame, University of Manchester, UK
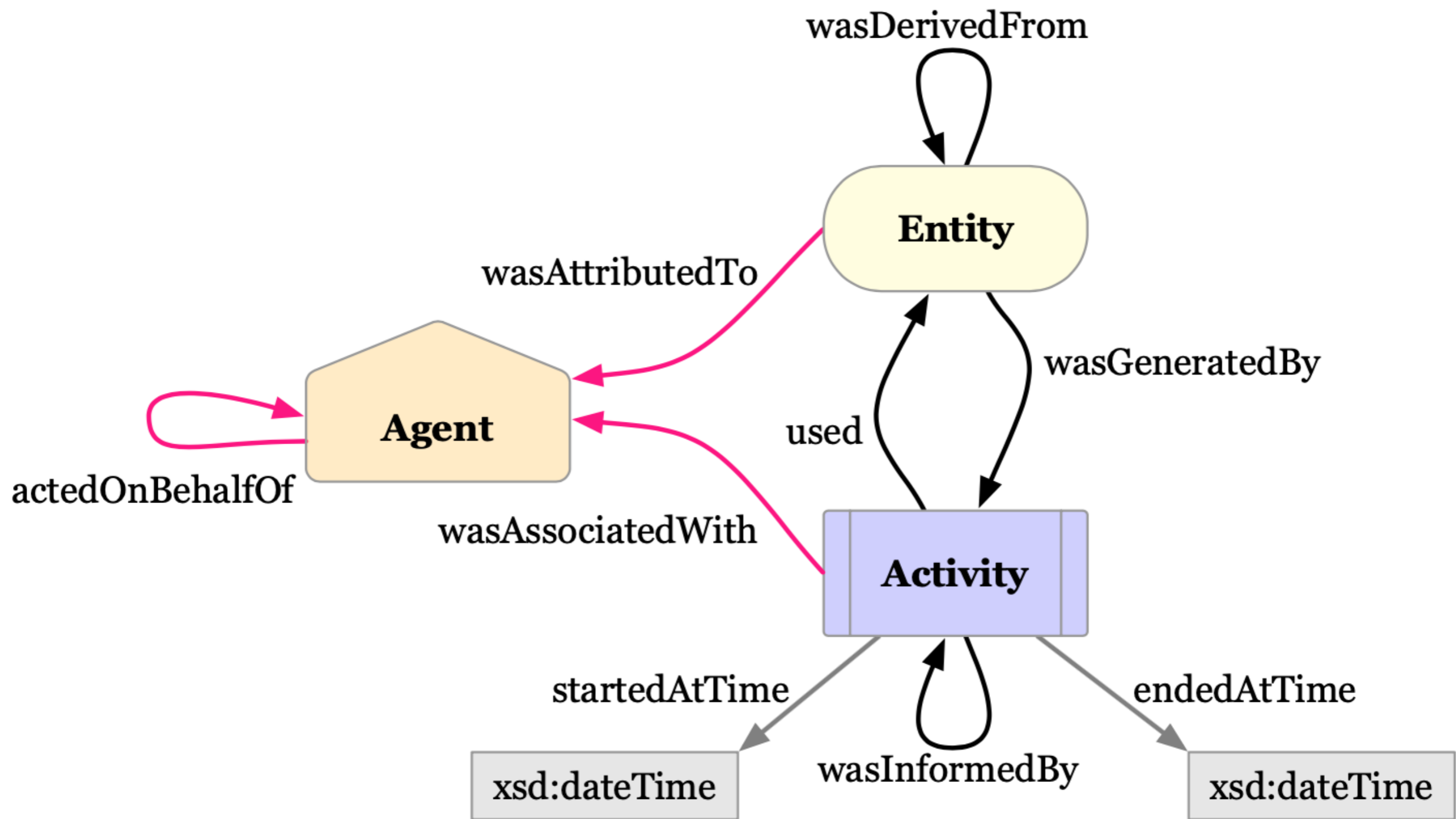  James Cheney, University of Edinburgh, UK
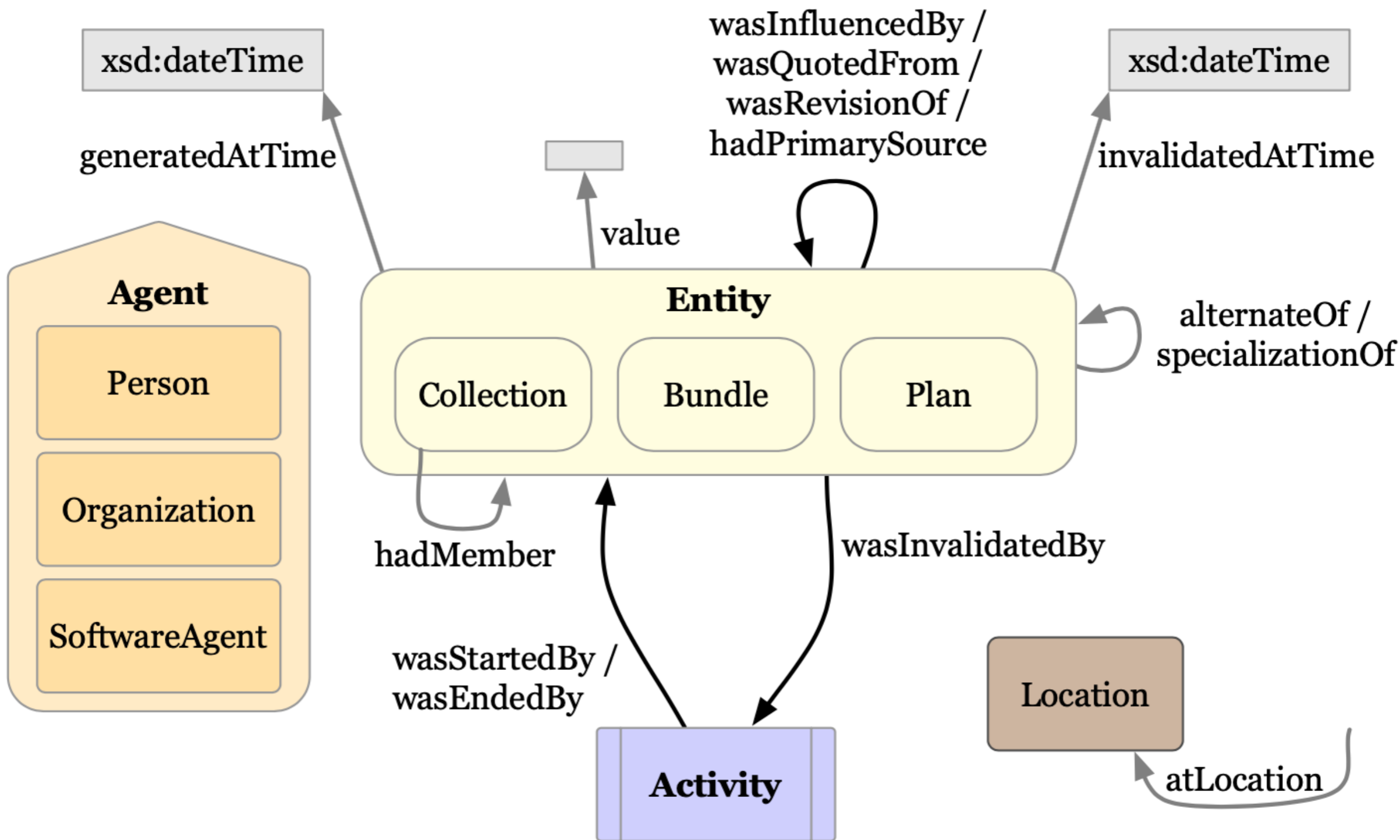  David Corsar, University of Aberdeen, UK
  Daniel Garijo, Ontology Engineering Group, Universidad Politécnica de Madrid, Spain
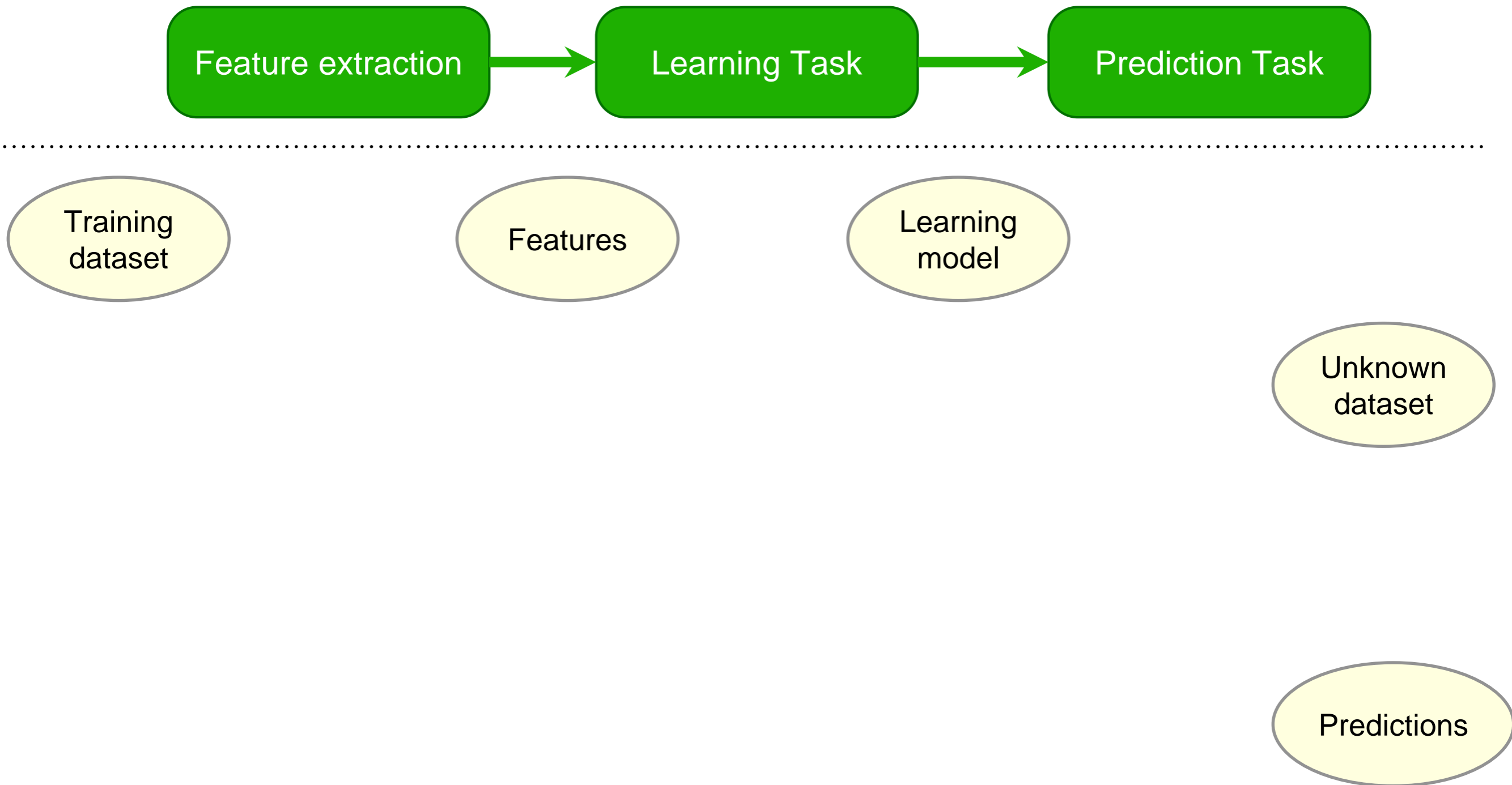  Stian Soiland-Reyes, University of Manchester, UK
  Stephan Zednik, Rensselaer Polytechnic Institute, USA
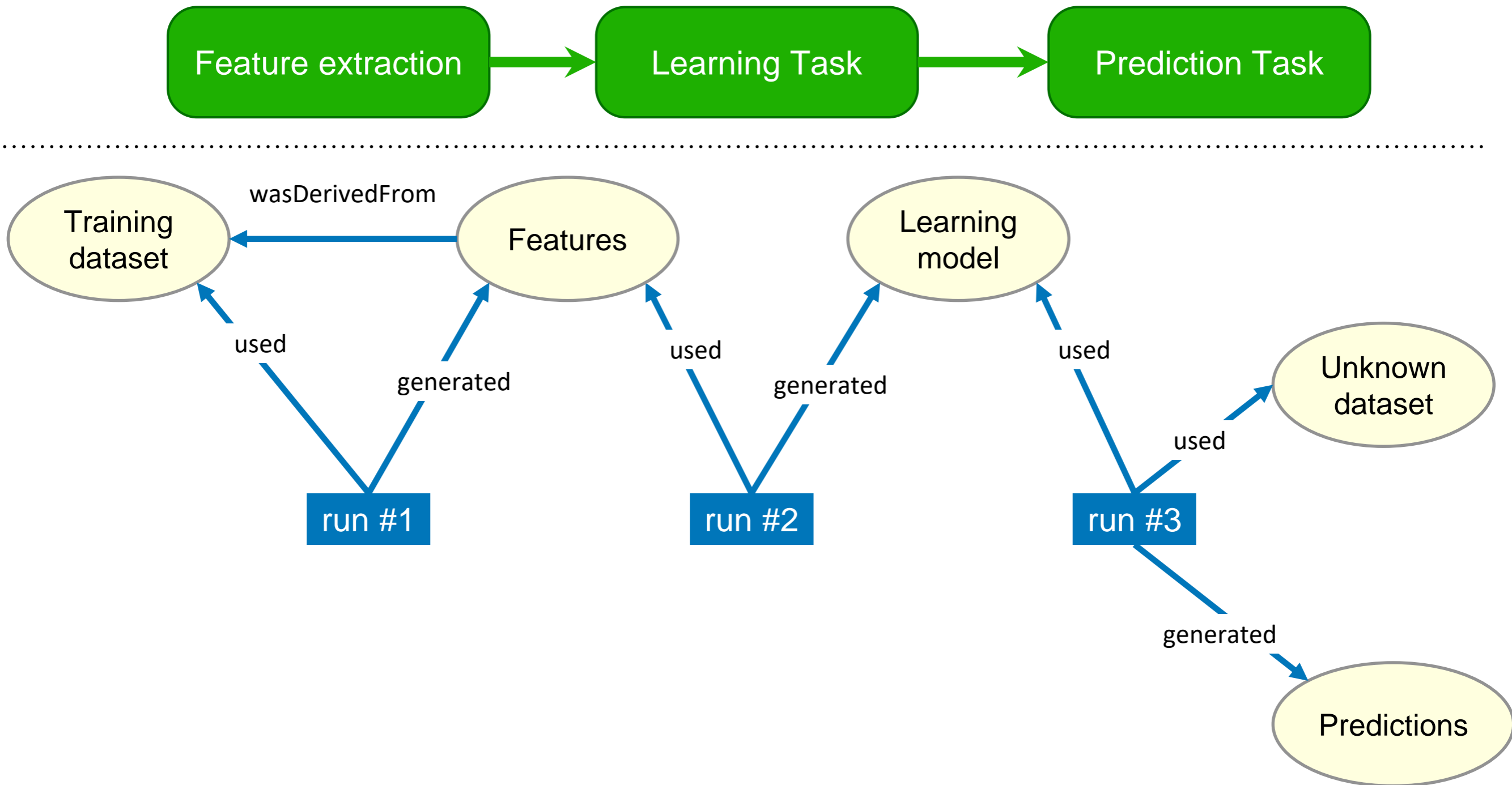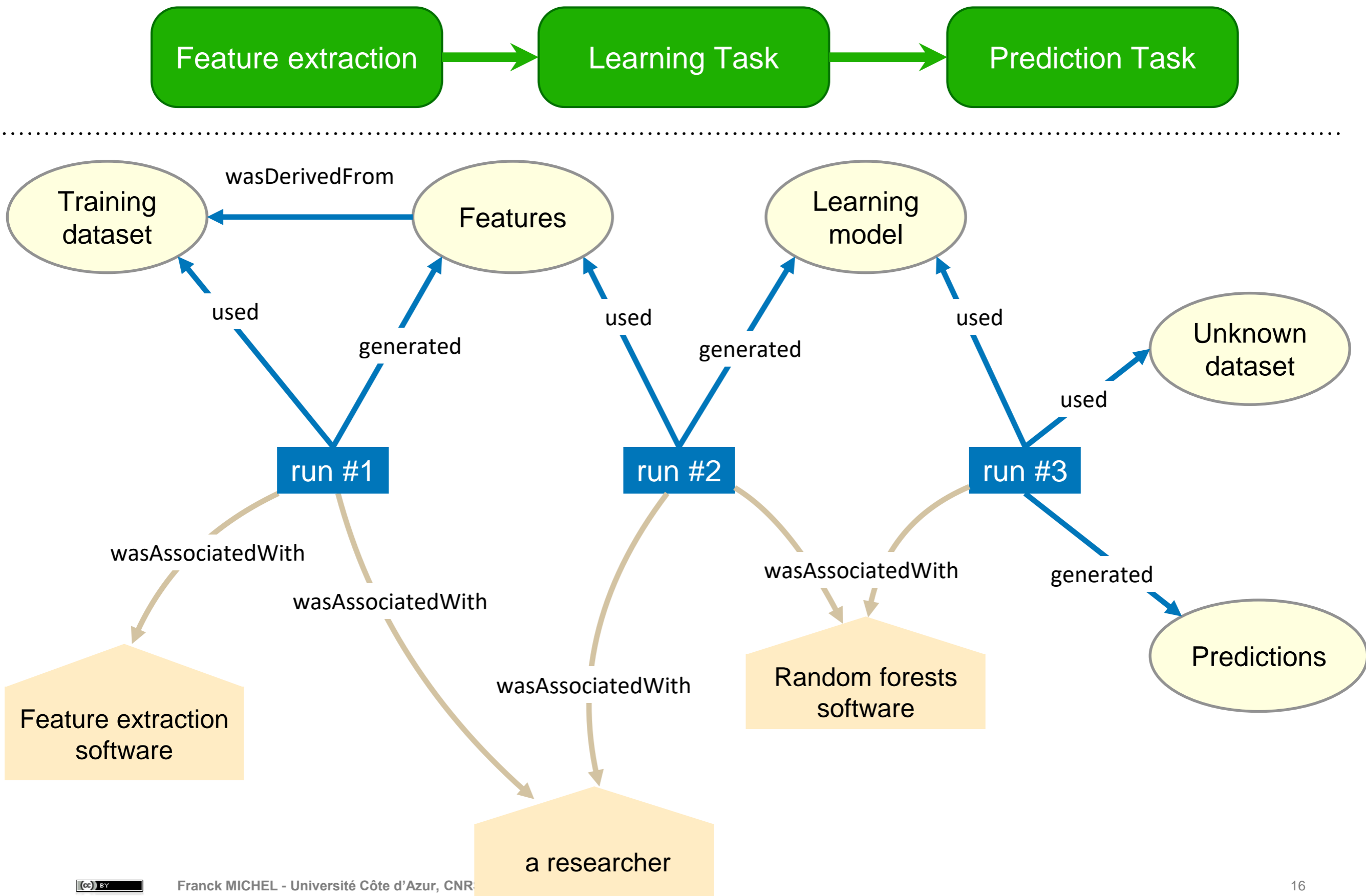  Jun Zhao, University of Oxford, UK

# Example: Provenance for a ML workflow
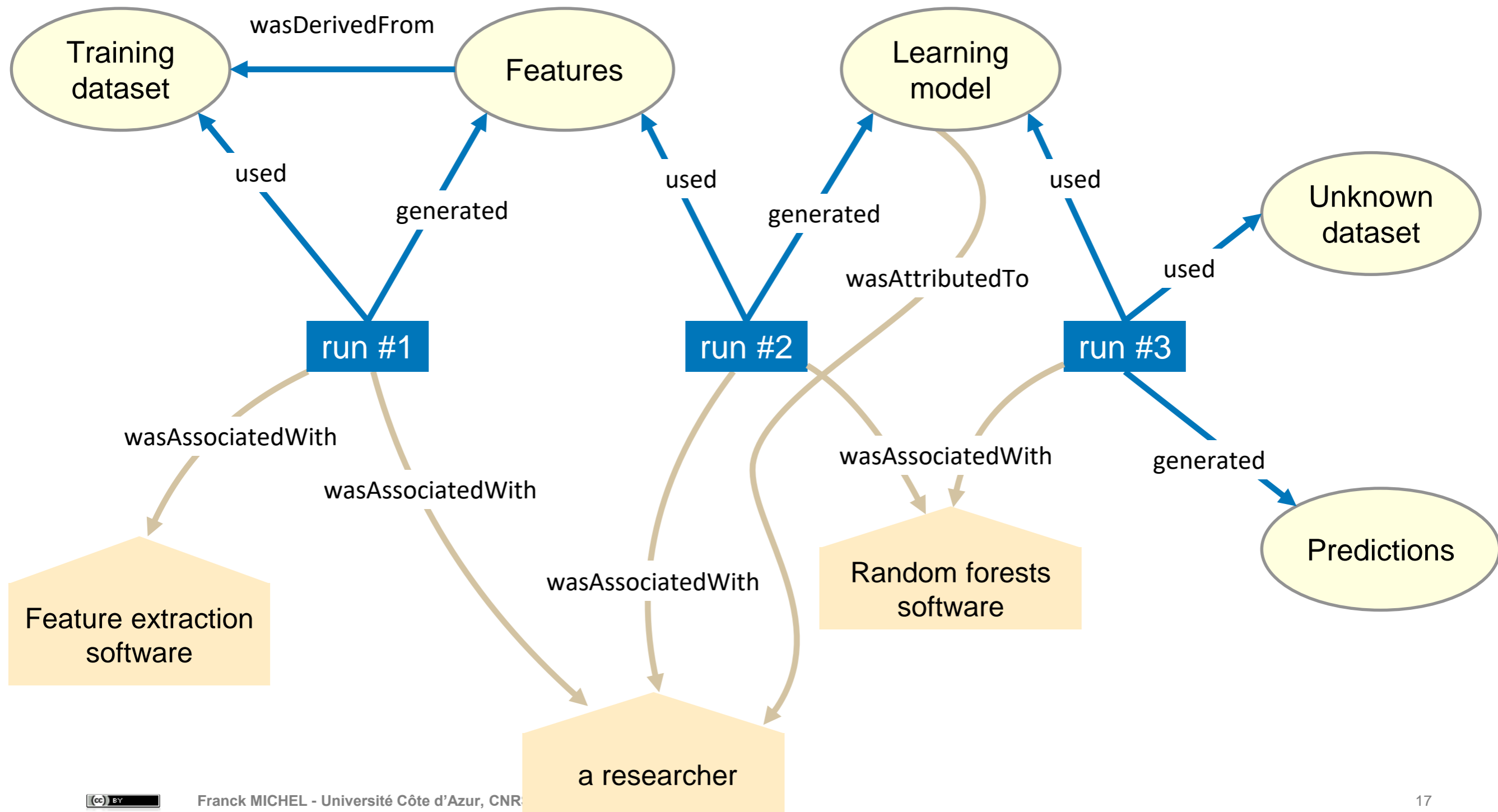
# Example: Provenance for a ML workflow

# Example: Provenance for a ML workflow

# Example: Provenance for a ML workflow

# Reasoning with provenance

W3C

## Constraints of the PROV Data Model

### W3C Recommendation 30 April 2013

**This version:**
http://www.w3.org/TR/2013/REC-prov-constraints-20130430/
**Latest published version:**
http://www.w3.org/TR/prov-constraints/
**Test suite:**
http://dvcs.w3.org/hg/prov/raw-file/default/testcases/process.html
**Implementation report:**
http://www.w3.org/TR/2013/NOTE-prov-implementations-20130430/
**Previous version:**
http://www.w3.org/TR/2013/PR-prov-constraints-20130312/ (color-coded diff)
**Editors:**
James Cheney, University of Edinburgh
Paolo Missier, Newcastle University
Luc Moreau, University of Southampton
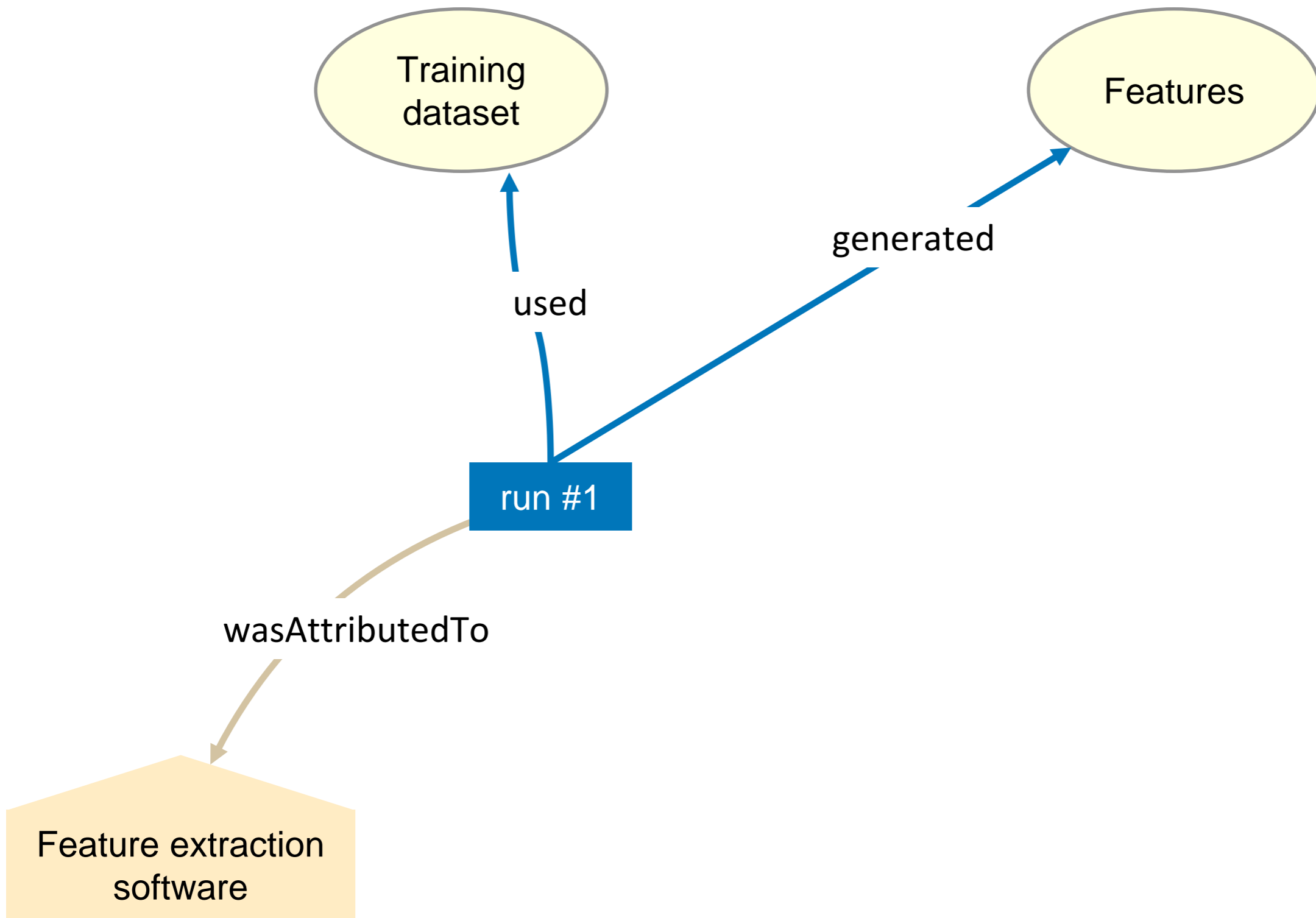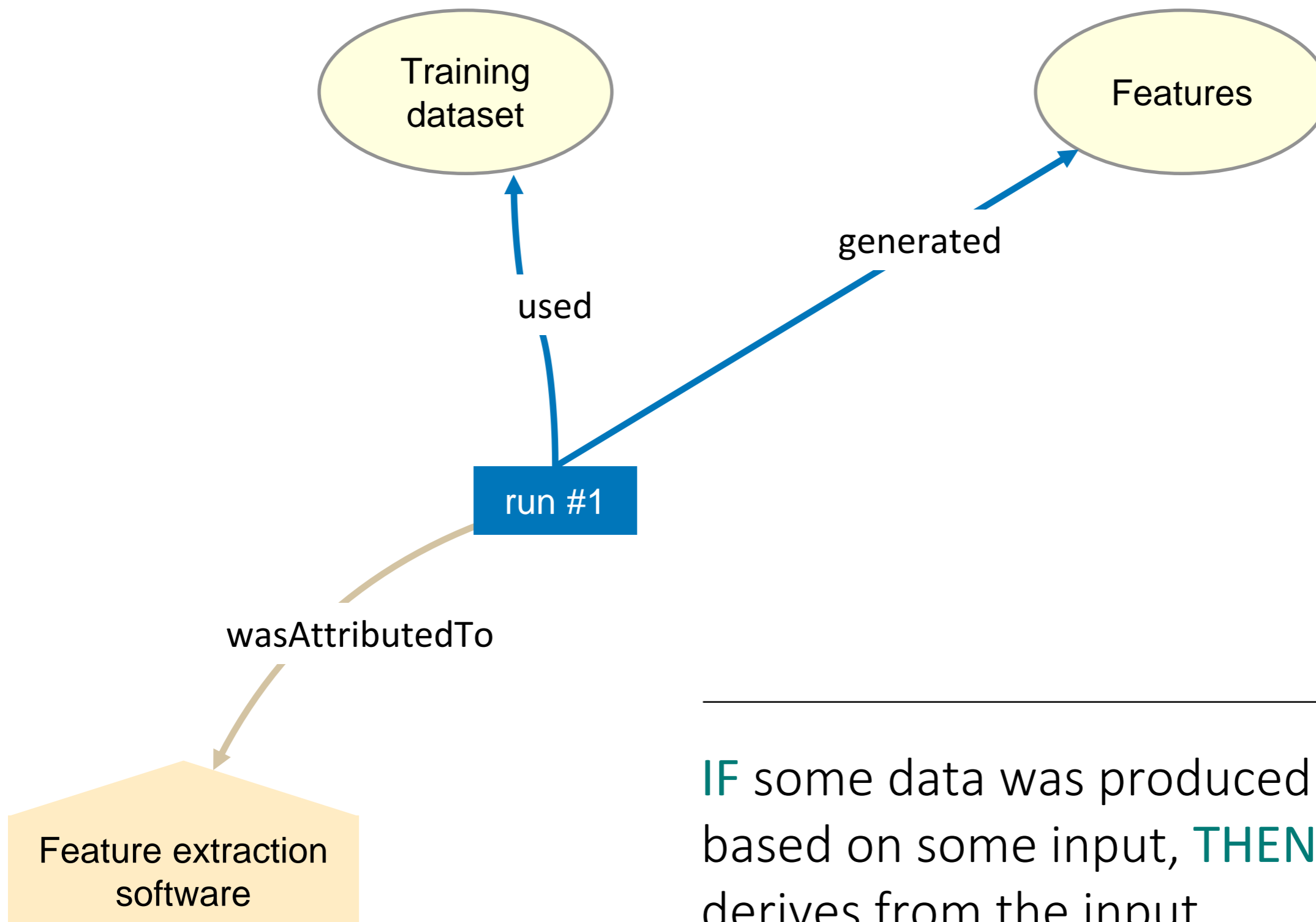**Author:**
Tom De Nies, iMinds - Ghent University

Please refer to the **errata** for this document, which may include some normative corrections.

The English version of this specification is the only normative version. Non-normative translations may also be available.
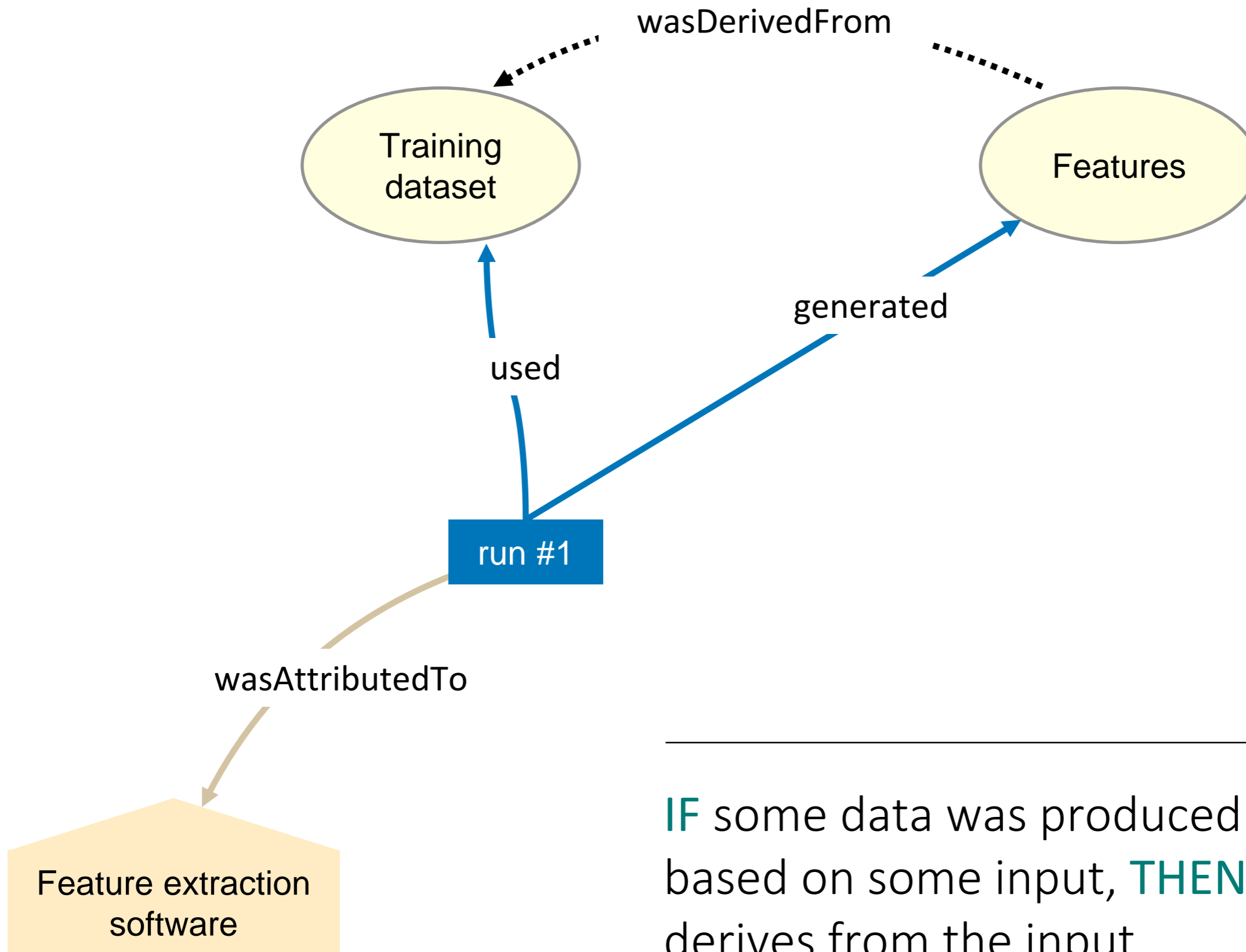
**Training dataset**

**Features**

used

generated

**run #1**

wasAttributedTo

**Feature extraction software**

---

IF some data was produced by a tool based on some input, THEN this data derives from the input

IF some data was produced by a tool based on some input, THEN this data derives from the input

# Is provenance enough for reproducibility?

"An activity used an entity". What to denote inputs' roles?

E.g. image registration: 2 inputs image, image to register + atlas

Problem of hidden parameters

Configuration parameters passed inline in a script

Hyper-parameters of a DNN

Hence the need for **PROV-O extensions**

ProvONE: PROV Extension for Scientific Workflow Provenance

Sensor Data Provenance: SSNO and PROV-O

PAV: Provenance, Authoring and Versioning ontology

…

# Is provenance enough for reuse?

```
11    a prov:Bundle, prov:Entity;
12    prov:wasAttributedTo <#galaxy2prov>;
13    prov:generatedAtTime "2016-04-14T18:18:37.000409"^^xsd:dateTime;
14  .
15
16  <#72486b583fe152f0>
17    a prov:Activity ;
18    prov:wasAssociatedWith <#cat1> ;
19    prov:startedAtTime "2015-12-15T12:54:50.749845"^^xsd:dateTime;
20    prov:endedAtTime "2015-12-15T12:55:57.016799"^^xsd:dateTime;
```

Too fine-grained
No domain concepts

**Visualise**

# Semantic tools catalogs

# Domain-centric provenance summary



Knowledge Graph

**2**

# What about ML tools?

# ML-Schema

**W3C** Community Group

Goal: *Improve **interoperability**, **reproducibility** and **interpretability** of DM/ML experiments*

How:

- Define a schema to represent/share information on DM/ML **algorithms**, **datasets** and **experiments**

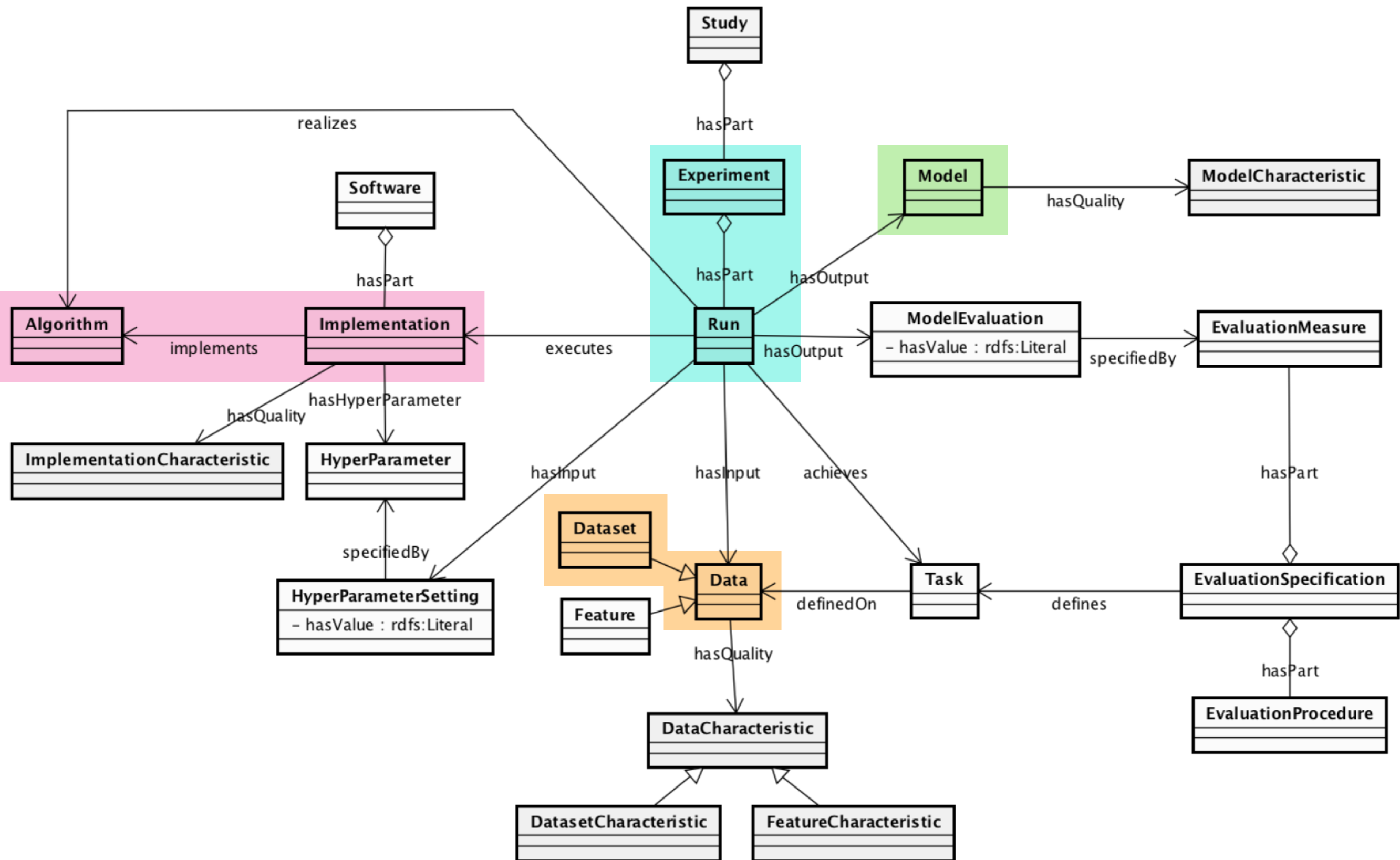- Align existing DM/ML ontologies to this schema, develop ontologies for specific purposes/applications

- Turn DM/ML algorithms and results into **Linked Open Data**

# ML-Schema

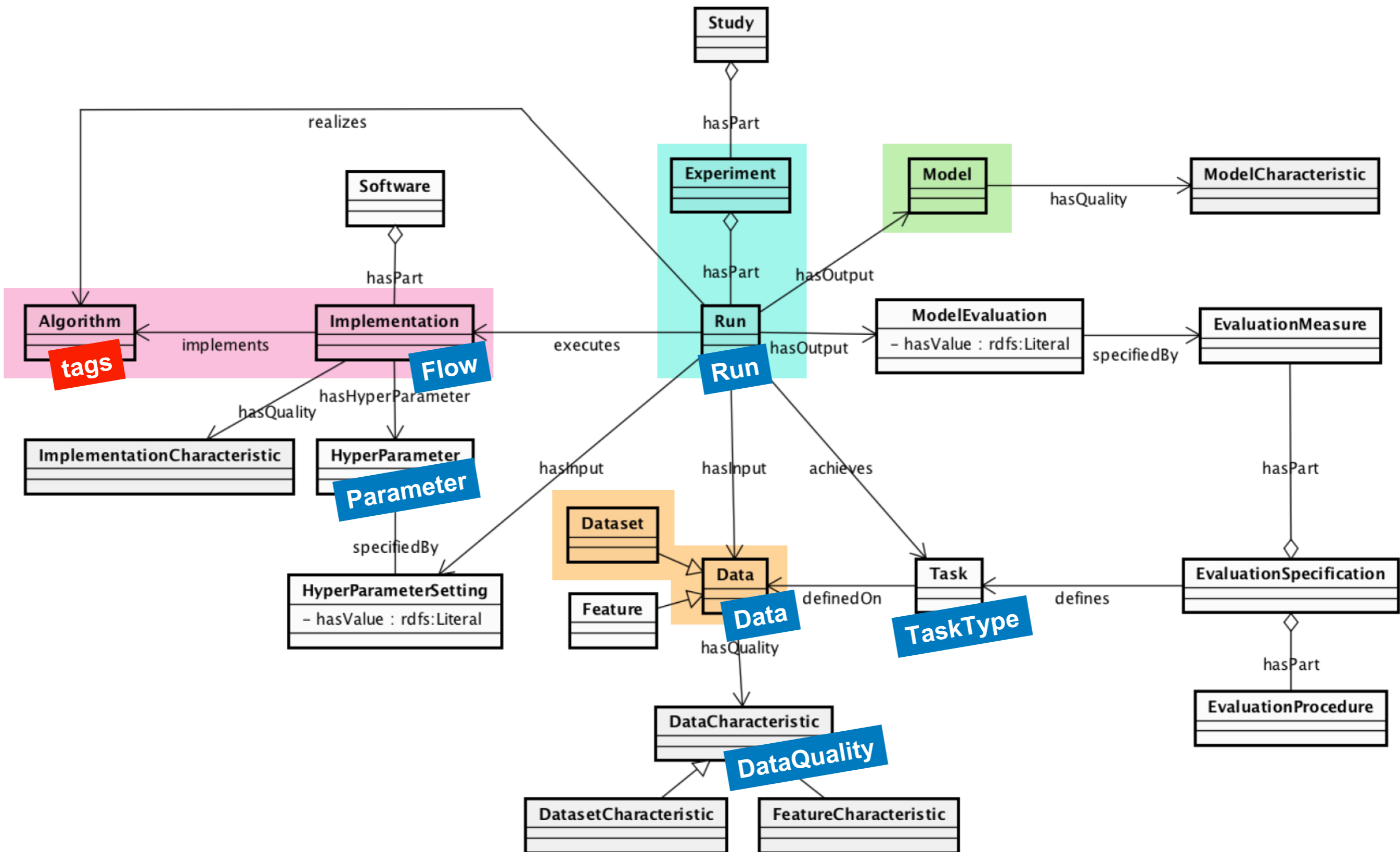## Inspired by previous works on ML/DM vocabularies

- **OntoDM**-core ontology: DM, based on BFO

- **Exposé** ontology: ML experiments, on top of OntoDM. Used in OpenML

- **D**ata **M**ining **OP**timization ontology: taxonomy of DM algos and ML models

- **MEX**: lightweight vocabulary to exchange basic ML metadata

Publio G. C., Esteves D., Ławrynowicz A., Panov P., Soldatova L., Soru T., Vanschoren J. & Zafar H. (2018). **ML-Schema: Exposing the Semantics of Machine Learning with Schemas and Ontologies**. In *Proc. of the 2nd Reproducibility in Machine Learning*, p. 5. Stockholm, Sweeden.

# ML-Schema

Publio G. C., Esteves D., Ławrynowicz A., Panov P., Soldatova L., Soru T., Vanschoren J. & Zafar H. (2018). **ML-Schema: Exposing the Semantics of Machine Learning with Schemas and Ontologies**. In *Proc. of the 2nd Reproducibility in Machine Learning*, p. 5. Stockholm, Sweeden.
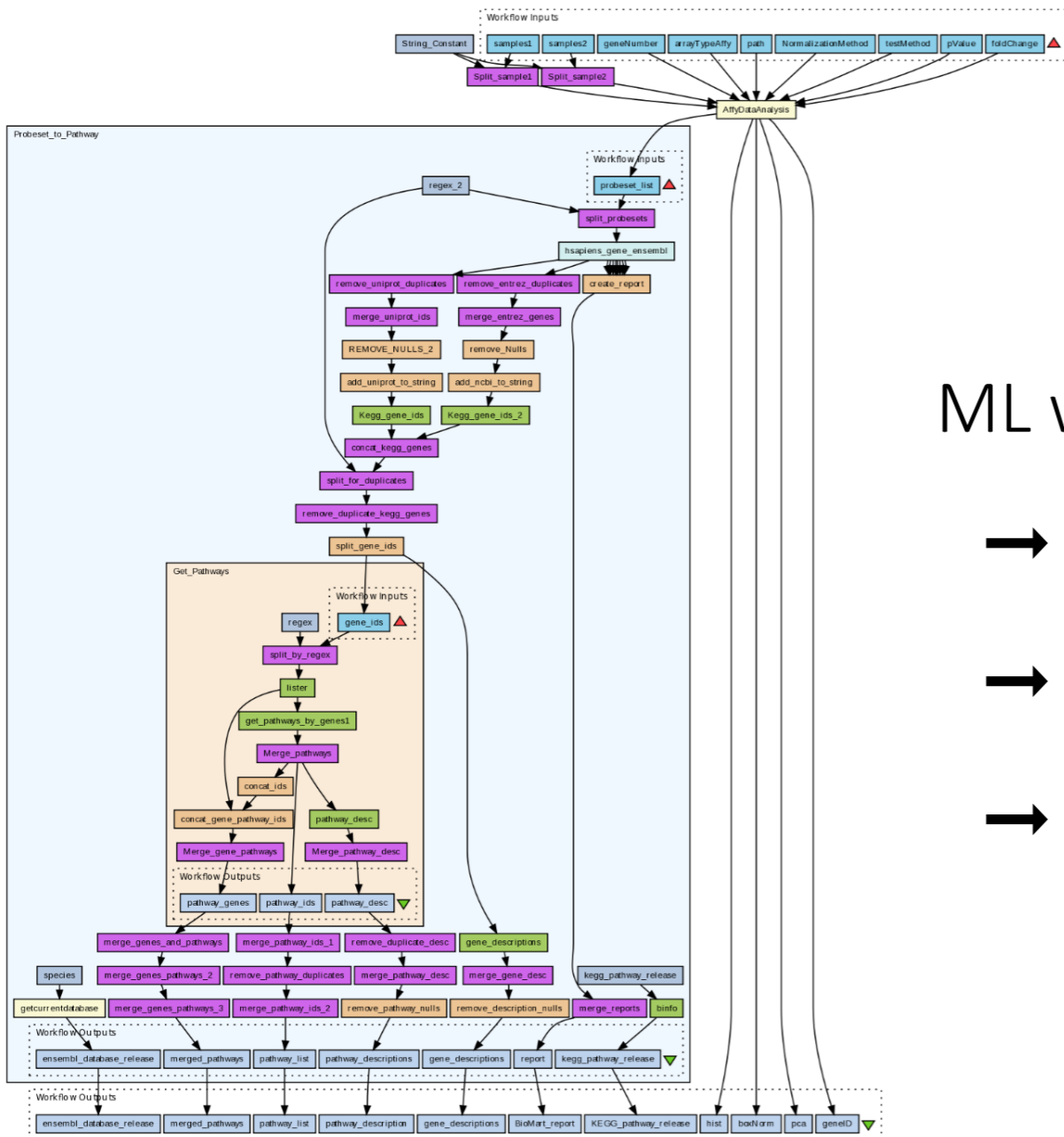
# Imperfect mapping of ML-Schema to OpenML

Source: https://github.com/ML-Schema/core/wiki/Vocabulary-(mappings)

# ML-Schema

## Current status (personal feeling)

- Still a preliminary work: a simple Turtle file

- Need for more complete recommendations
  (e.g. wrt. use of 3rd party vocabularies: DC Terms, DCAT etc.)

- Need for richer mapping descriptions

- Need for real interoperability tests

But a key to instrument ML algos/datasets/experiments together with PROV-O-instrumented workflow engines

ML workflows:

→ **automation** of data analysis

→ **abstraction**: ML-Schema + extensions

→ **provenance**: PROV-O + extensions

# Take-aways & open questions

**Scientific Workflows** ➡ automation, abstraction, provenance

Standards for **provenance representation** and **reasoning**

Reproducibility/reuse requires **domain-specific tools' description** and **provenance-enabled workflow engines**
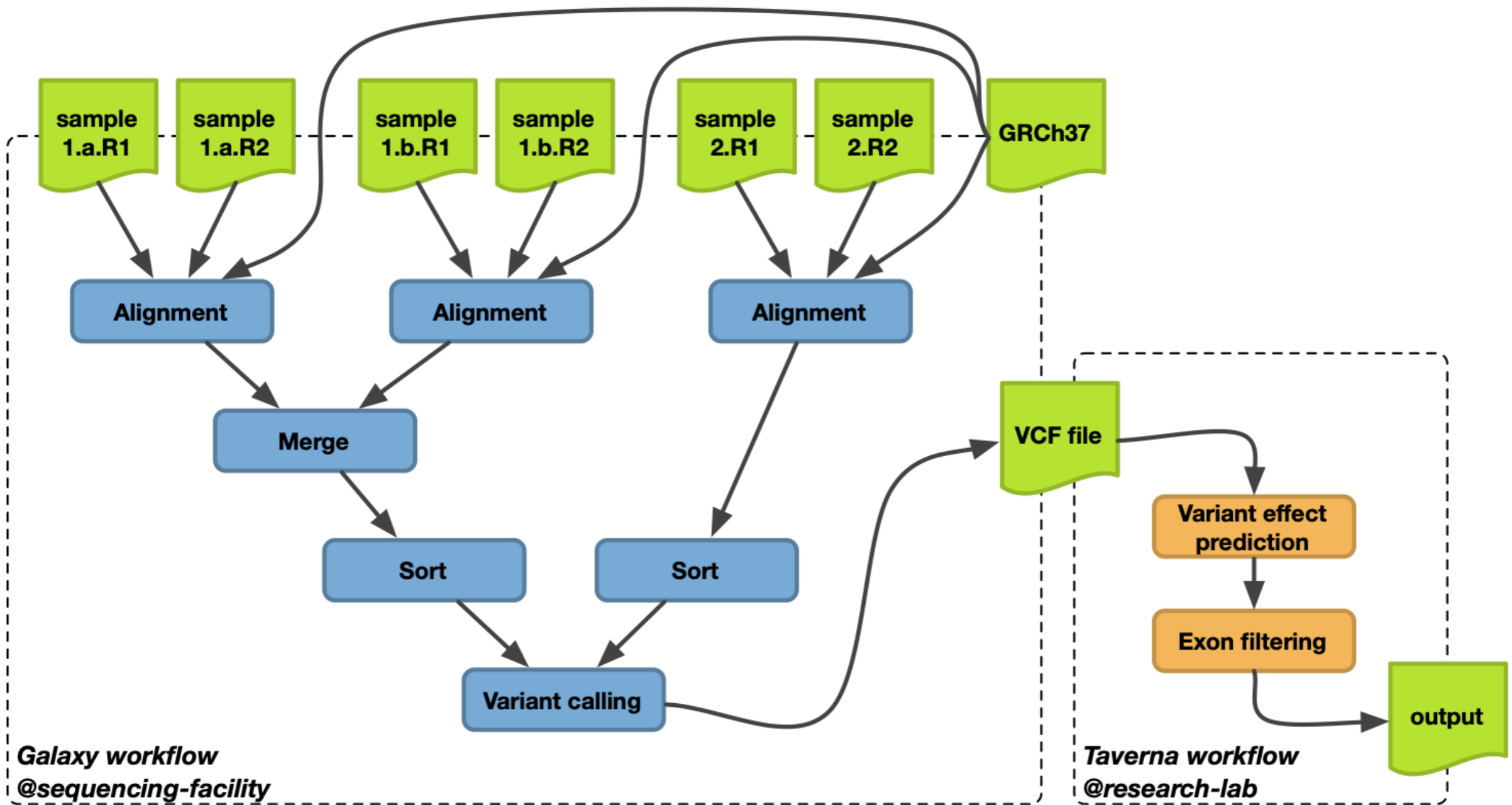
**ML-Schema + PROV-O** ➡ the future winning couple?

Distributed data analysis ➡ **Distributed provenance, reasoning**?
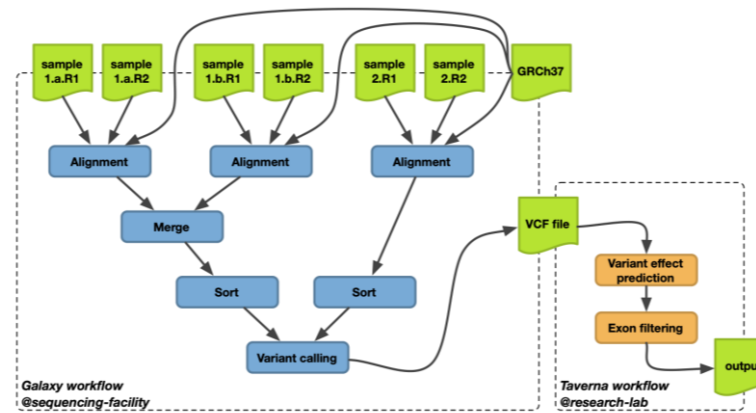
# Backup slides

# Provenance
# in **multi-site** studies ?

# Multi-site studies ➡ ≠ workflow engines !



Scattered provenance capture ?

# Provenance issues



« Which alignment algorithm was used
when predicting these pathogenic store? »

« A new version of a reference genome is available, which
genome was used when predicting these phenotypes ? »

Need for an overall tracking of provenance
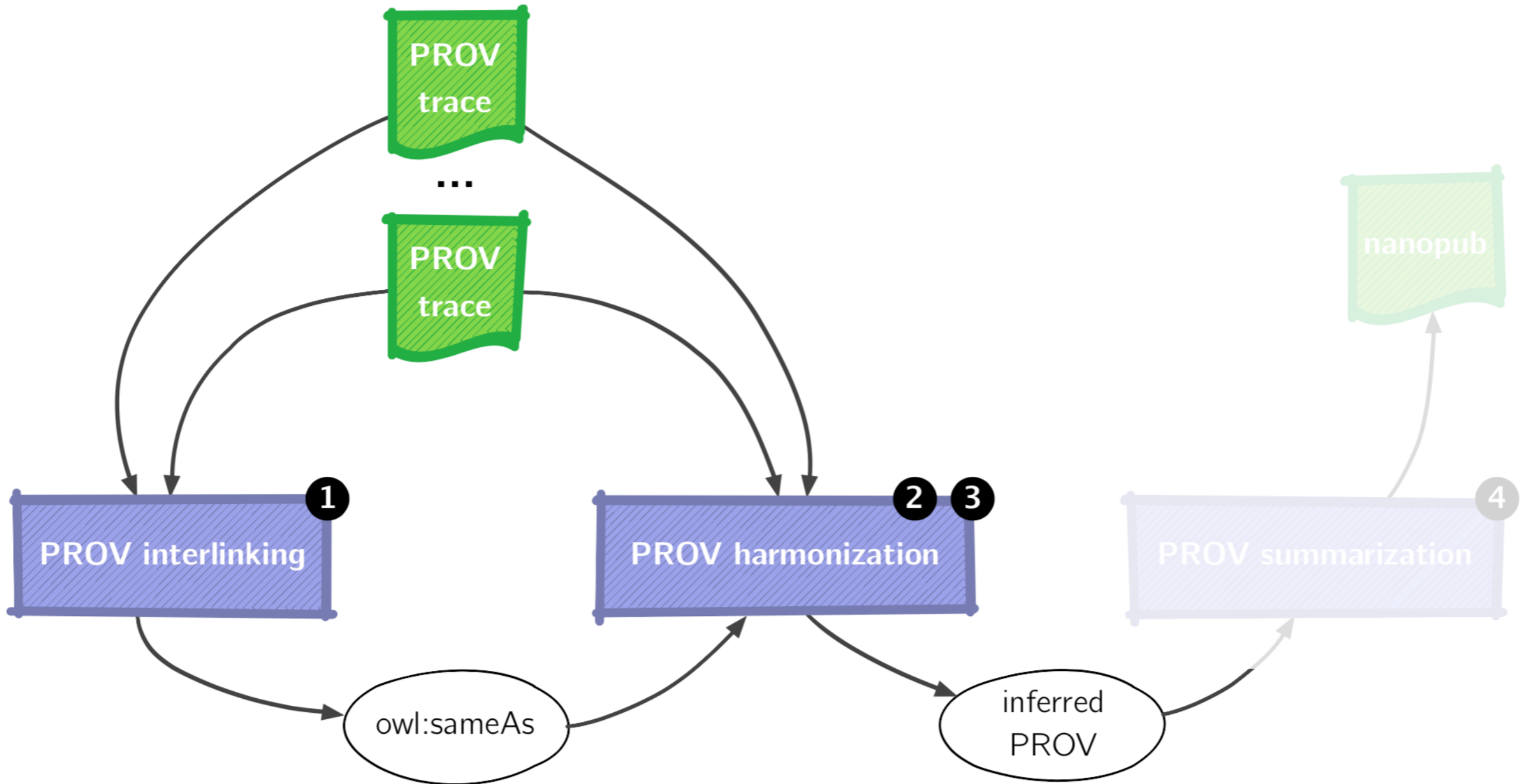over both Galaxy and Taverna workflows !

# Provenance « heterogeneity »

| Galaxy PROV predicates | counts |
|---|---|
| prov:wasDerivedFrom | 118 |
| rdf:type | 76 |
| rdfs:label | 62 |
| prov:used | 61 |
| prov:wasAttributedTo | 34 |
| prov:wasGeneratedBy | 33 |
| prov:endedAtTime | 26 |
| prov:startedAtTime | 26 |
| prov:wasAssociatedWith | 26 |
| prov:generatedAtTime | 1 |

| Taverna PROV predicates | counts |
|---|---|
| rdf:type | 54 |
| rdfs:label | 13 |
| prov:atTime | 8 |
| wfprov:describedByParameter | 6 |
| rdfs:comment | 6 |
| prov:hadRole | 6 |
| prov:activity | 5 |
| dcterms:hasPart | 4 |
| prov:agent | 4 |
| prov:endedAtTime | 4 |
| prov:hadPlan | 4 |
| prov:qualifiedAssociation | 4 |
| prov:qualifiedEnd | 4 |
| prov:qualifiedStart | 4 |
| prov:startedAtTime | 4 |
| prov:wasAssociatedWith | 4 |
| tavernaprov:content | 3 |
| wfprov:usedInput | 3 |
| wfprov:wasEnactedBy | 3 |
| wfprov:wasOutputFrom | 3 |

How to reconcile these provenance traces?
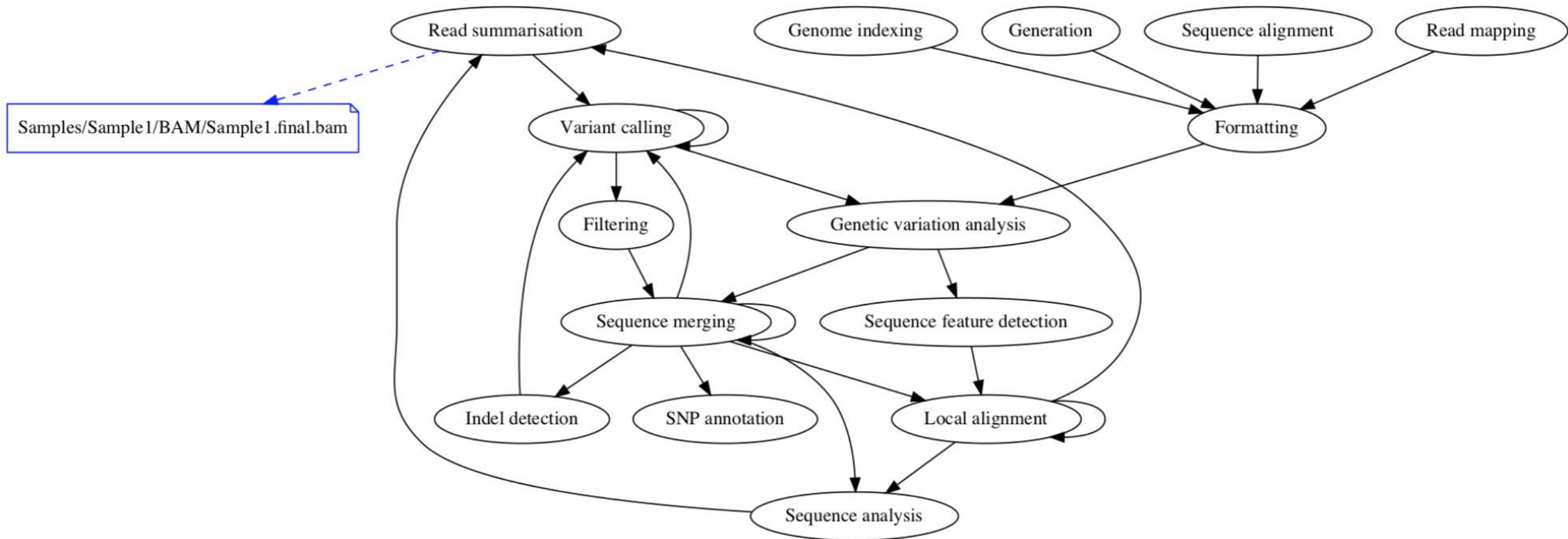
# Approach



A. Gaignard, K. Belhajjame, H. Skaf-Molli. **SHARP: Harmonizing and Bridging Cross-Workflow Provenance**. *The Semantic Web: ESWC 2017 Satellite Events Portorož, Slovenia, May 28 – June 1, 2017, Revised Selected Papers, 2017*

# Results

Reconciled provenance as an « influence graph »

https://github.com/albangaignard/sharp-prov-toolbox

# Provenance summary



```
...
The file Samples/Sample1/BAM/Sample1.realign.bai results from
tool gatk2_indel_realigner-IP which Locally align two or more molecular
sequences.

It was produced in the context of Rare Coding Variants in ANGPTL6 Are
Associated with Familial Forms of Intracranial Aneurysm
...
```

**4**