

# Découvrabilité et réutilisation de données produites par des workflows (en génomique)

*IC 2021*

Alban Gaignard<sup>1</sup>, Hala Skaf-Molli<sup>2</sup>, Khalid Belhajjame<sup>3</sup>

02 juillet 2021

<sup>1</sup> l'institut du thorax, INSERM, CNRS, Université de Nantes

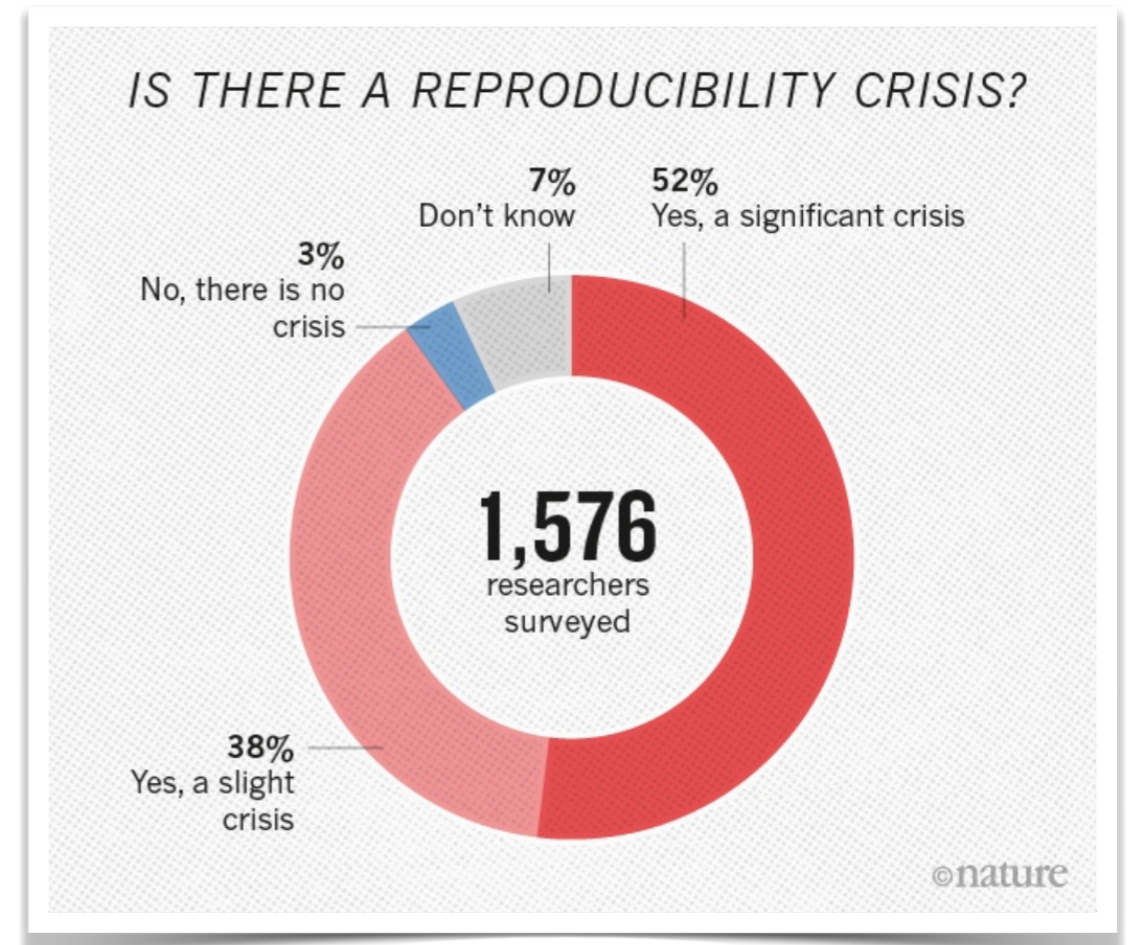
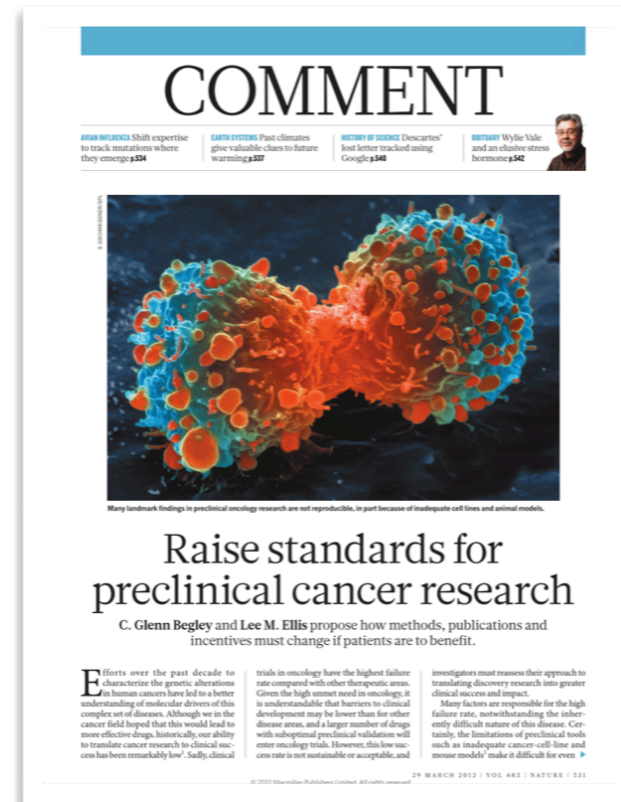
<sup>2</sup> LS2N - CNRS, Université de Nantes, France

<sup>3</sup> LAMSADE – PSL, Université Paris-Dauphine, France



Context

# Knowledge production



COMMENT 612 | NATURE | VOL 505 | 30 JANUARY 2014

## NIH plans to enhance reproducibility

Francis S. Collins and Lawrence A. Tabak discuss initiatives that the US National Institutes of Health is exploring to restore the self-correcting nature of preclinical research.

A growing chorus of concern, from scientists and laypeople, contends that the complex system for ensuring the reproducibility of biomedical research is failing and is in need of restructuring<sup>1,2</sup>. As leaders of the US National Institutes of Health (NIH), we share this concern and here explore some of the significant interventions that we are planning.

Science has long been regarded as 'self-correcting', given that it is founded on the shorter term, however, the checks and balances that once ensured scientific fidelity have been hobbled. This has compromised the ability of today's researchers to reproduce others' findings.

Let's be clear: with rare exceptions, we have no evidence to suggest that irreproducibility is about scientific misconduct. In 2011, the Office of Research Integrity of the US Department of Health and Human Services pursued only 12 such cases<sup>3</sup>.

« researchers made headlines when they declared that they had been **unable to reproduce the findings in 47 of 53 'landmark' cancer papers** »

# "Reusing" is challenging

## Repeat

Same experiment

Same setup

Same lab

## Replicate

Same experiment

Same setup

~~Same lab~~

## Reproduce

Same experiment

~~Same setup~~

~~Same lab~~

## Reuse

*new ideas,*

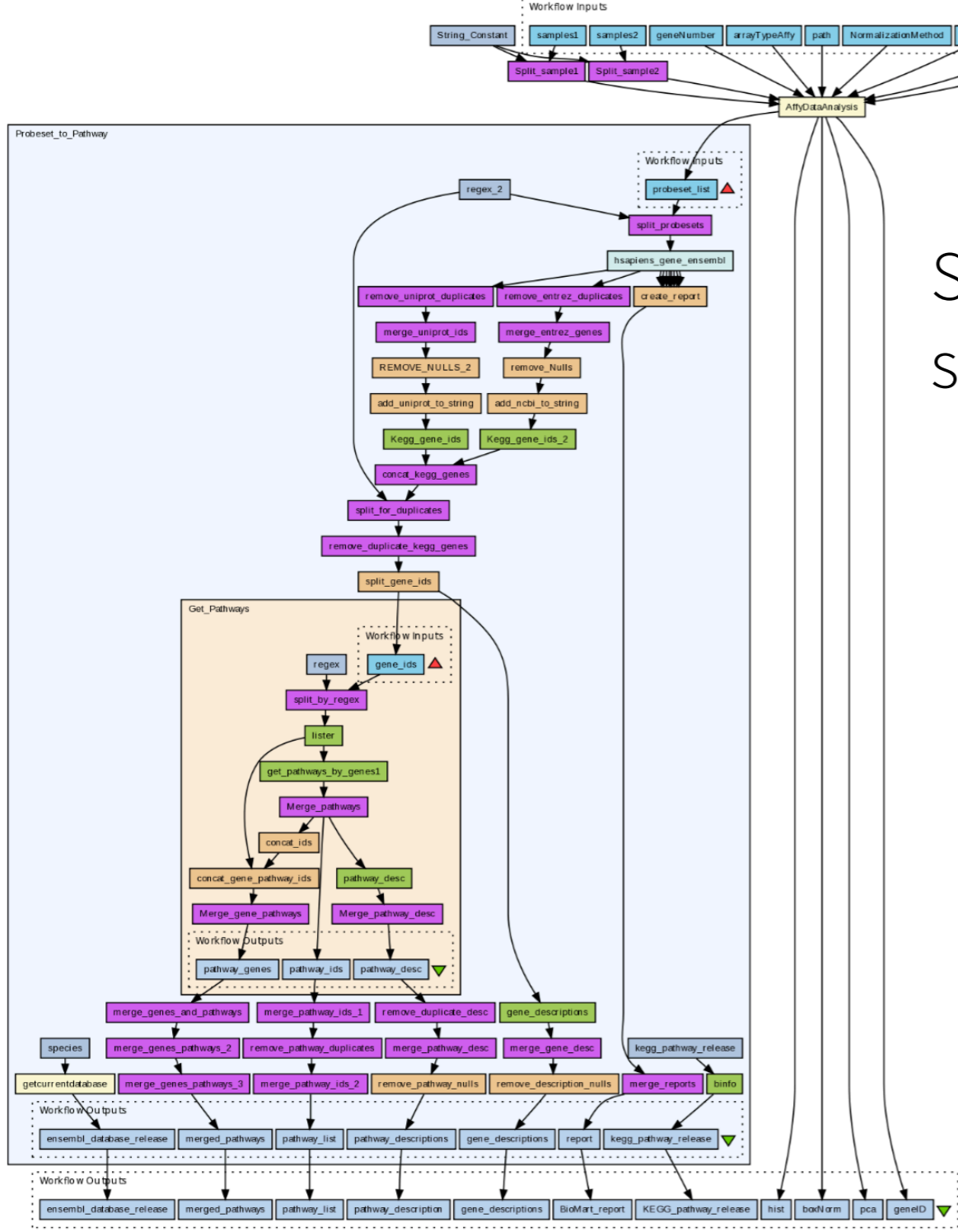
*some commonalities,*

*new experiment*

# Scientific **workflows** to the rescue ...

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

« Workflow management **systems** (WMS) **perform** the complex analyses on a variety of **distributed resources** »



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

- **abstraction** (describe/share methods)
- **automate** data analysis (at scale)
- **provenance** (~transparency)

nextflow



pdidommaso / awesome-pipeline

A curated list of awesome pipeline toolkits inspired by Awesome Sysadmin

#awesome-list #workflow

228 commits | 1 branch | 0 releases | 42 contributors

Branch: master | New pull request | Create new file | Upload files | Find file | Clone or download

pdidommaso Update README.md | Latest commit @hwk 25 days ago

CONTRIBUTING.md | Added contributing | 4 years ago

README.md | Update README.md | 25 days ago

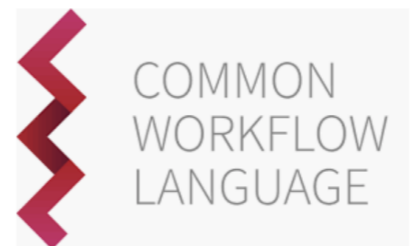
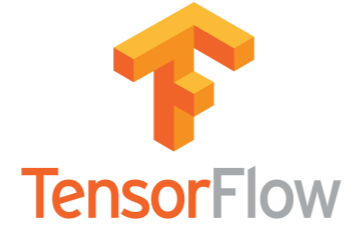
BI README.md

### Awesome Pipeline

A curated list of awesome pipeline toolkits inspired by [Awesome Sysadmin](#)

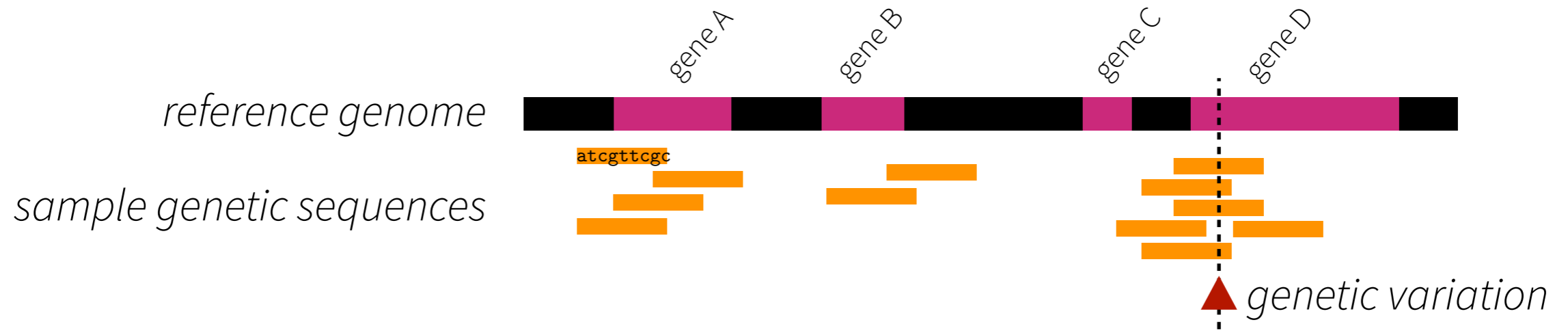
#### Pipeline frameworks & libraries

- **ActionChain** - A workflow system for simple linear success/failure workflows.
- **Adage** - Small package to describe workflows that are not completely known at definition time.
- **Airflow** - Python-based workflow system created by Airbnb.
- **Anduril** - Component-based workflow framework for scientific data analysis.
- **Anthra** - High-level language for biology.
- **Bids** - Scripting language for data pipelines.
- **BioMake** - GNU-Make-like utility for managing builds and complex workflows.
- **BioQueue** - Explicit framework with web monitoring and resource estimation.
- **Bistro** - Library to build and execute typed scientific workflows.
- **Bpipe** - Tool for running and managing bioinformatics pipelines.
- **Briefly** - Python Meta-programming Library for Job Flow Control.
- **Cluster Flow** - Command-line tool which uses common cluster managers to run bioinformatics pipelines.
- **Clusterjob** - Automated reproducibility, and hassle-free submission of computational jobs to clusters.
- **Compass** - Programming model for distributed infrastructures.
- **Conan2** - Light-weight workflow management application.
- **Consecution** - A Python pipeline abstraction inspired by Apache Storm topologies.
- **Cosmos** - Python library for massively parallel workflows.
- **Cromwell** - Workflow Management System geared towards scientific workflows from the Broad Institute.
- **Cuneiform** - Advanced functional workflow language and framework, implemented in Erlang.
- **Dagobah** - Simple DAG-based job scheduler in Python.
- **Dagr** - A scala based DSL and framework for writing and executing bioinformatics pipelines as Directed Acyclic Graphs.
- **Dask** - Dask is a flexible parallel computing library for analytics.
- **Dockerflow** - Workflow runner that uses Dataflow to run a series of tasks in Docker.
- **Dolt** - Task management & automation tool.
- **Drake** - Robust DSL akin to Make, implemented in Clojure.
- **Drake R package** - Reproducibility and high-performance computing with an easy R-focused interface. Unrelated to Factual's Drake.
- **Dray** - An engine for managing the execution of container-based workflows.
- **Fission Workflows** - A fast, lightweight workflow engine for serverless/FaaS functions.



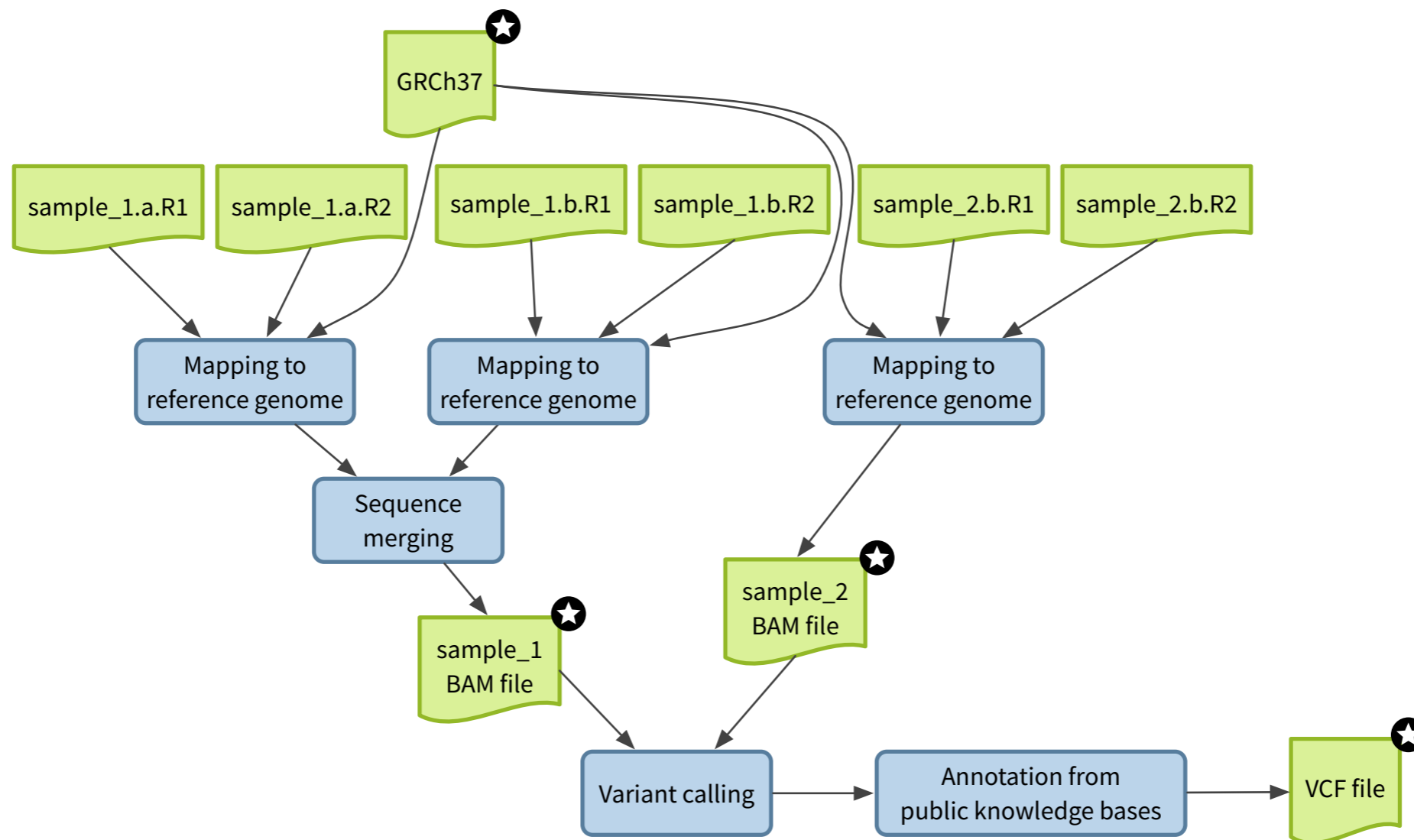
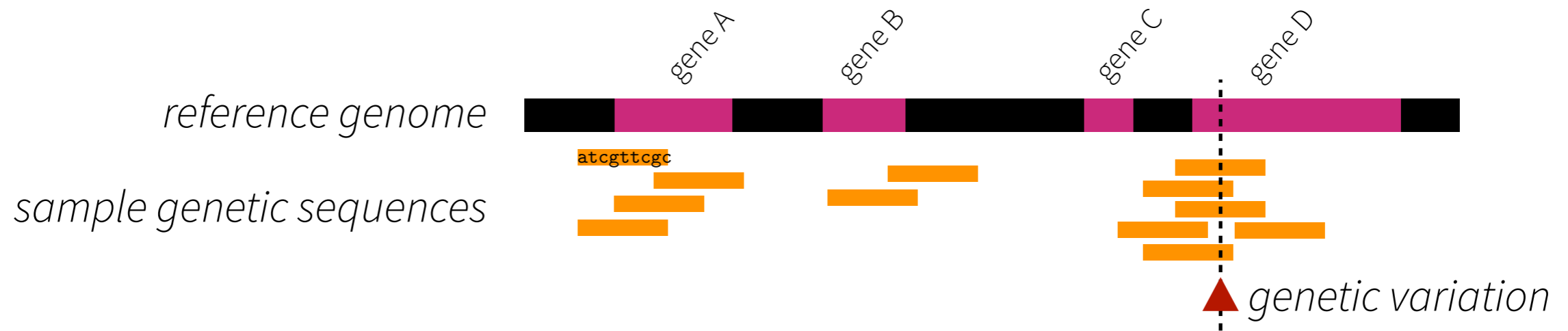
Problem statement

# Typical genomic variant detection workflow





# Typical genomic variant detection workflow



# Computational costs

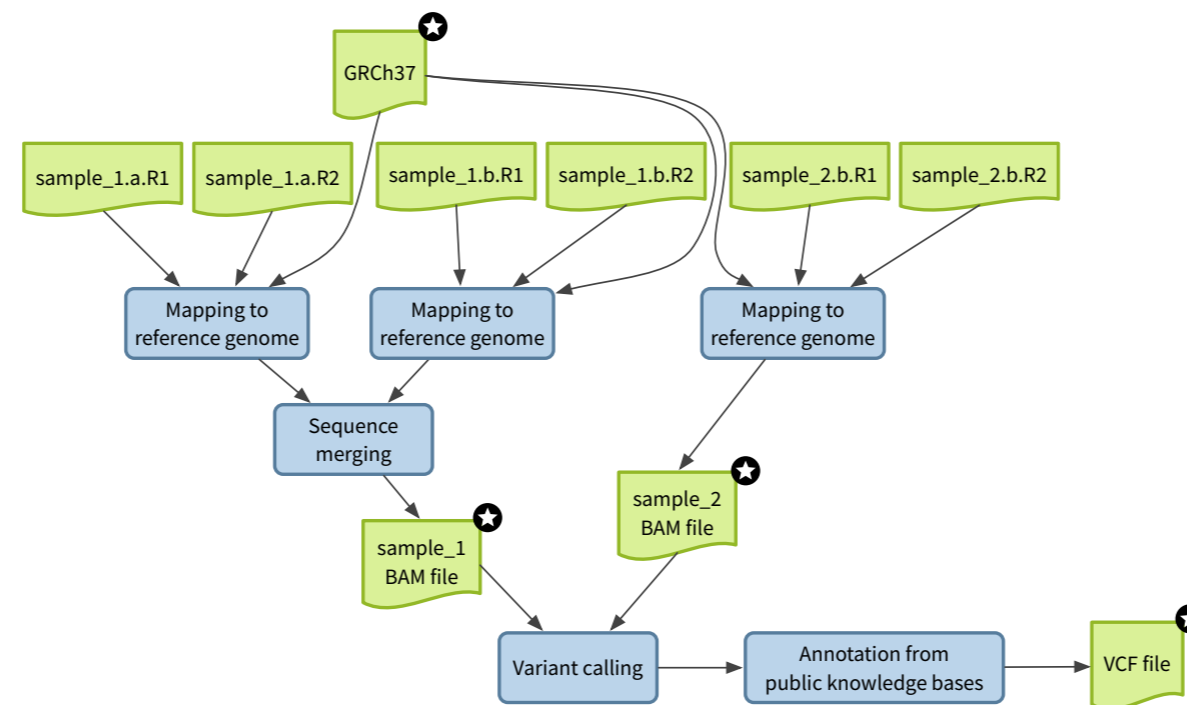
*as estimated by CNRGH, Evry*

	1 sample	100 samples	5000 samples
<b>Storing</b> processed exome data (compressed)	18.6 GB	1.8 TB	93 TB
<b>Computing</b> genomic variations (HPC, 196 cores)	2.5 hours	10+ days	1+ year
<b>Computing</b> (single CPU core)	20+ days / CPU	5+ years / CPU	250+ years / CPU

Re-computing: waste of computing time and storage

Can we better **reuse data** ?

# Issues when reusing bioinformatics data



« A new tool is available, which data subset should I reprocess ? »

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

→ need for an overall **tracking of provenance**

→ need for **domain-specific** contextual metadata

# Objectives

Limit the duplication of computing and storage efforts through better **data reuse**

1. Feed a **knowledge graph** with

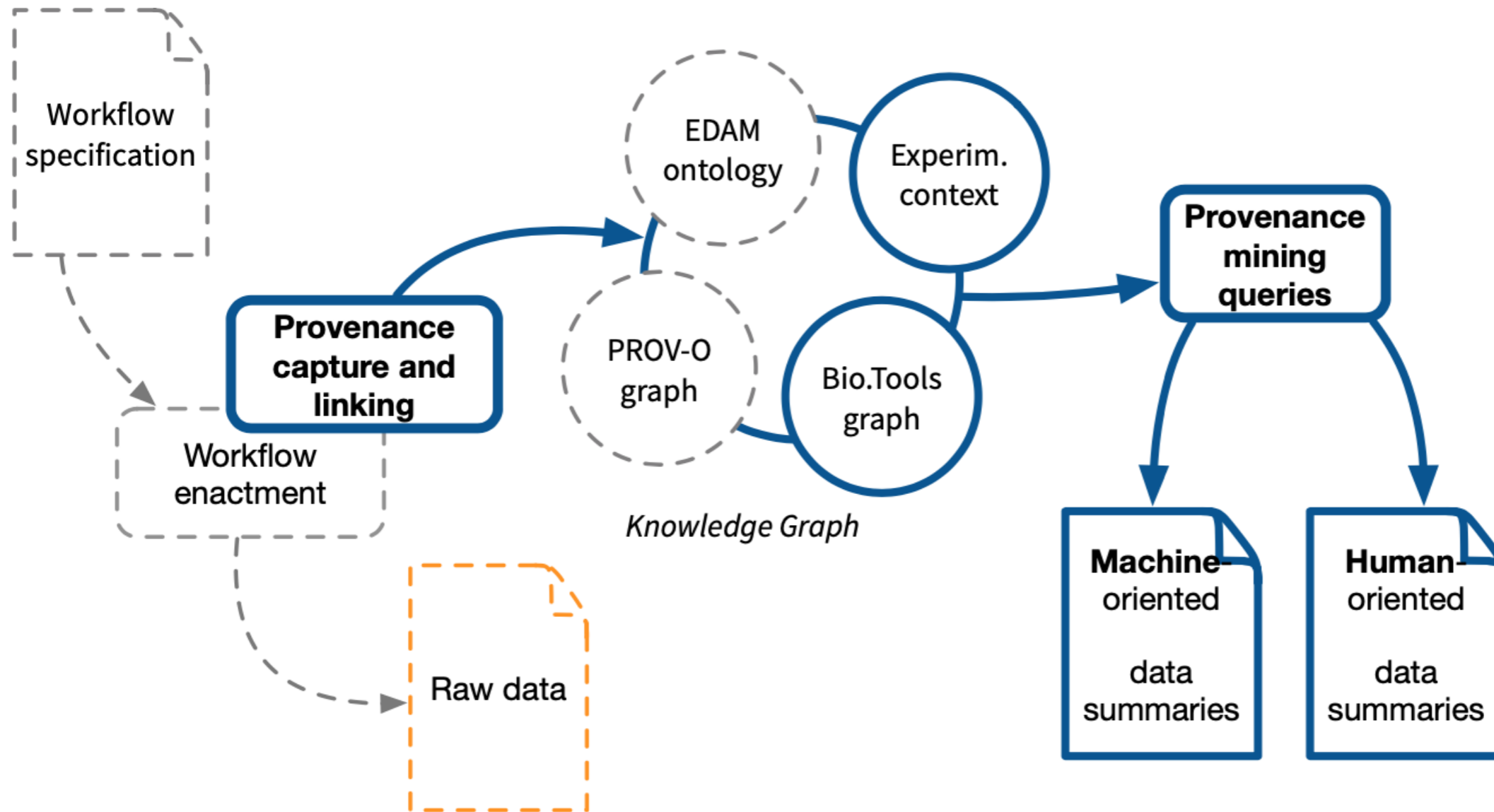
- generic **provenance** metadata
- **bioinformatics** community knowledge

2. Mine this knowledge graph to automatically produce

- **Machine**-oriented **data summaries**
- **Human**-oriented data **summaries**

Tooling a knowledge graph,  
to produce data summaries

# Approach



*“Which was the reference genome used to produce this VCF file ?”*

*“A new tool is available, which raw data should I reprocess ?”*

1. Capturing **provenance** at  
run-time

# What is provenance metadata

## Definitions in Computer Science

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

« Today, data is often made available on the Internet with no centralized control over its integrity: data is **constantly being created, copied, moved around, and combined** indiscriminately. Because information sources (or different parts of a single large source) may 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**. »



**Feature extraction**



**Learning task**



**Prediction task**

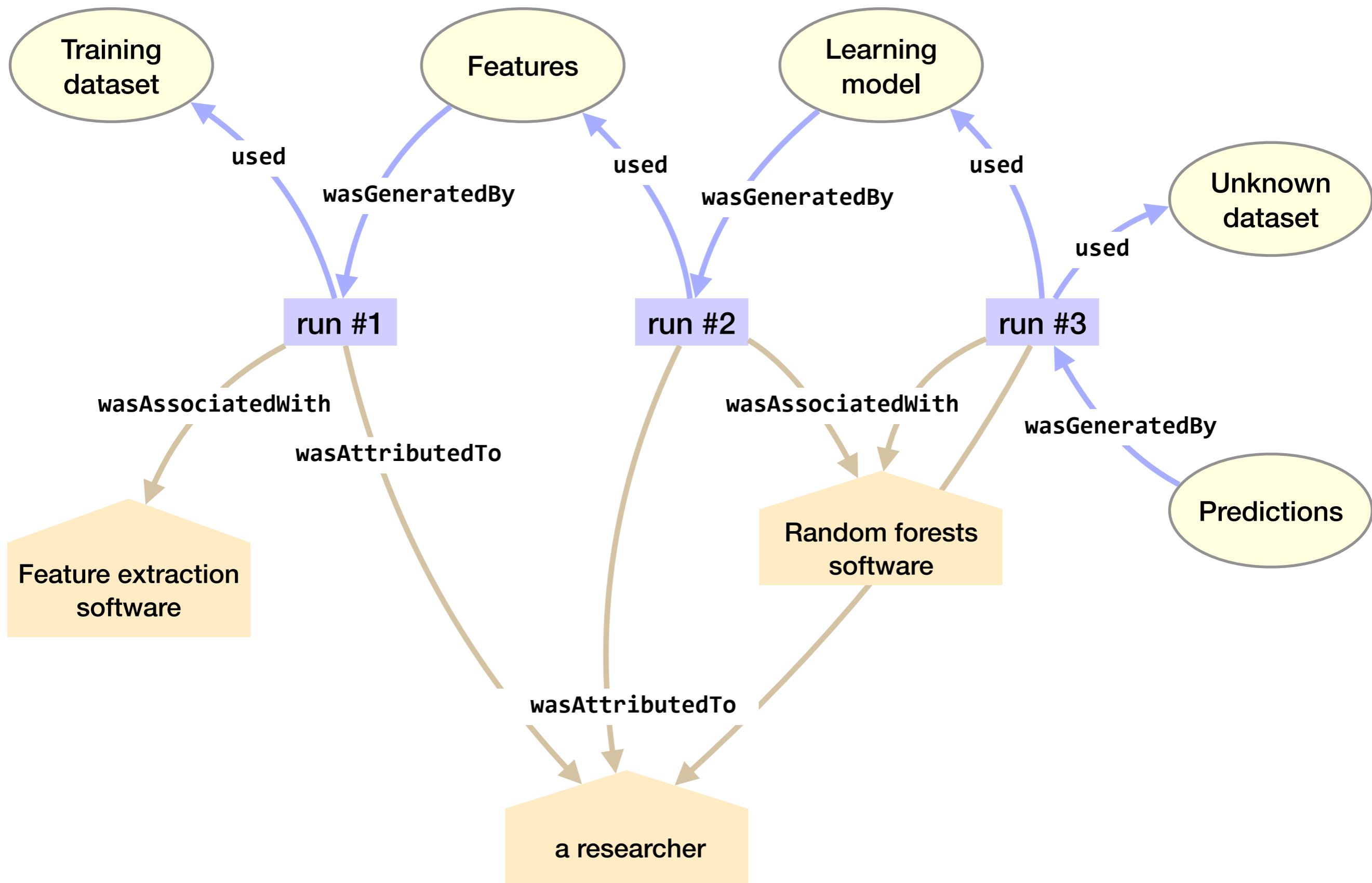
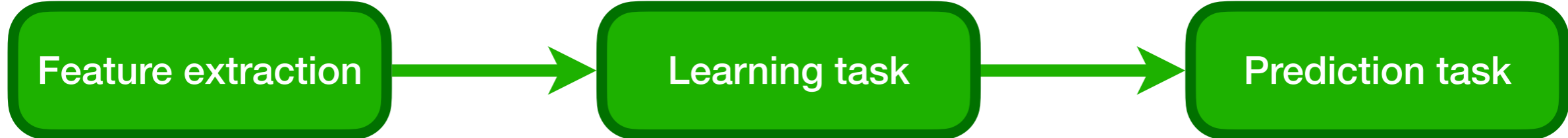
**Training  
dataset**

**Features**

**Learning  
model**

**Unknown  
dataset**

**Predictions**



# Representing provenance



## 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:**

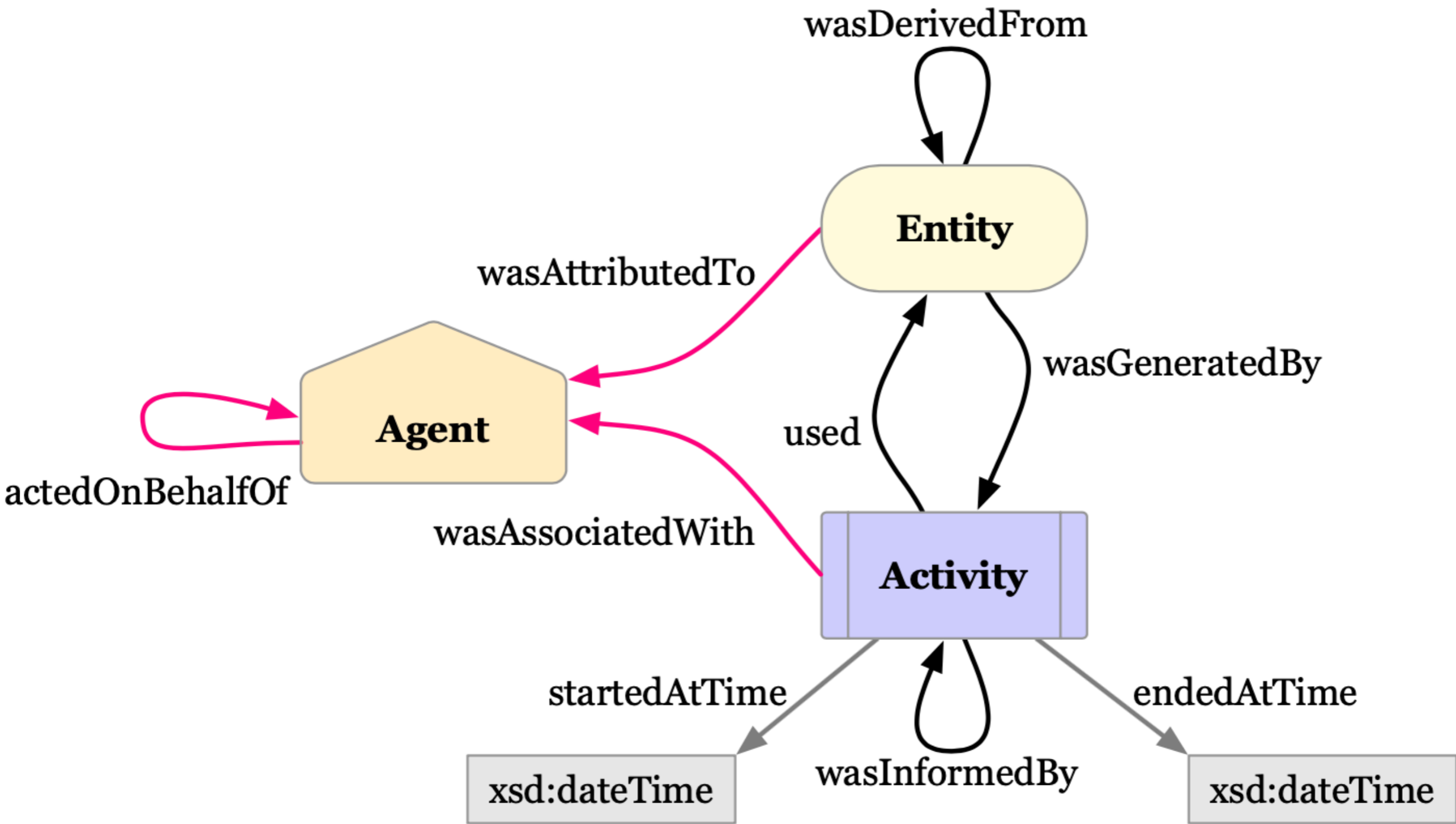
<http://www.w3.org/TR/2013/PR-prov-o-20130312/>

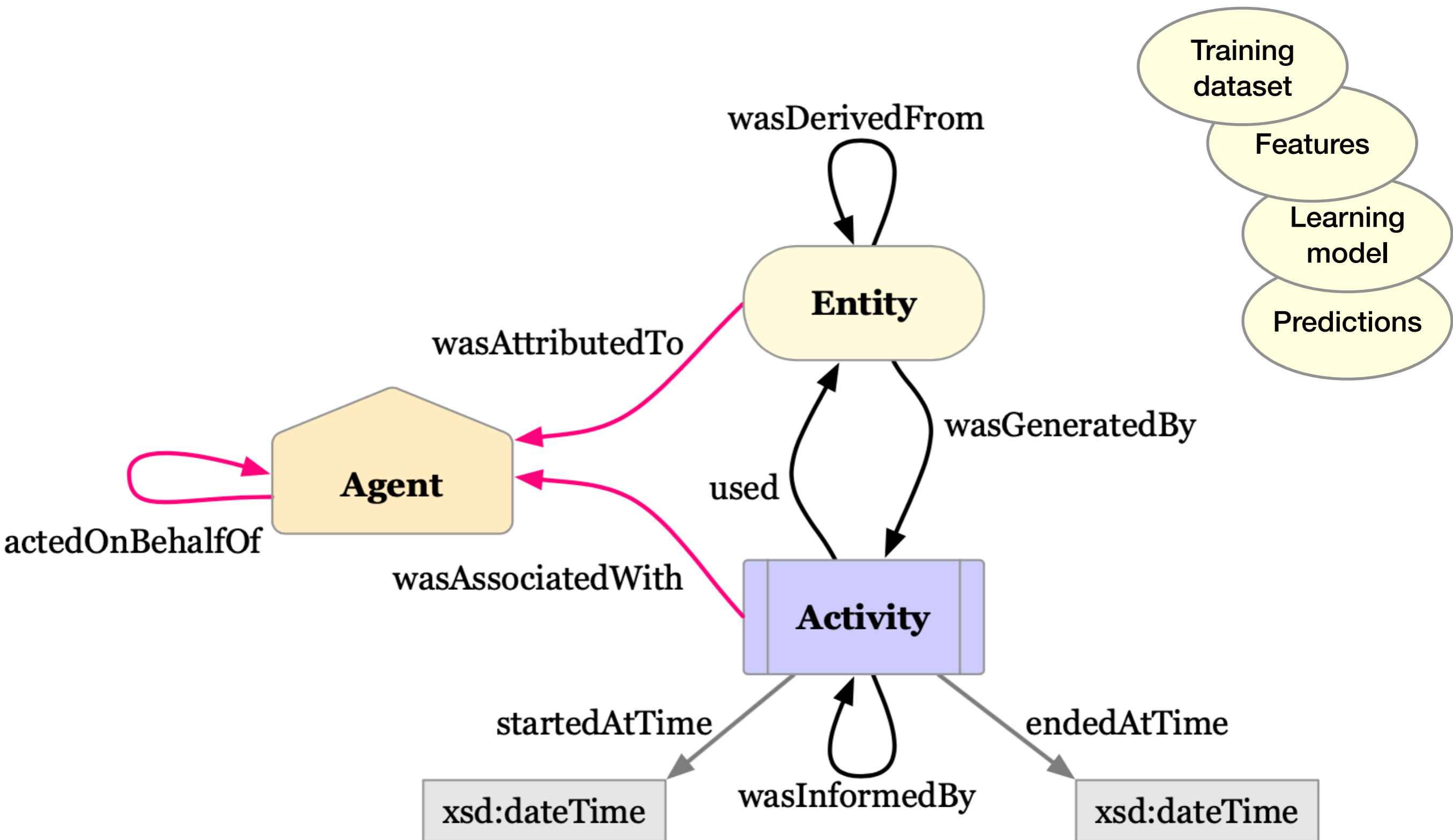
**Editors:**

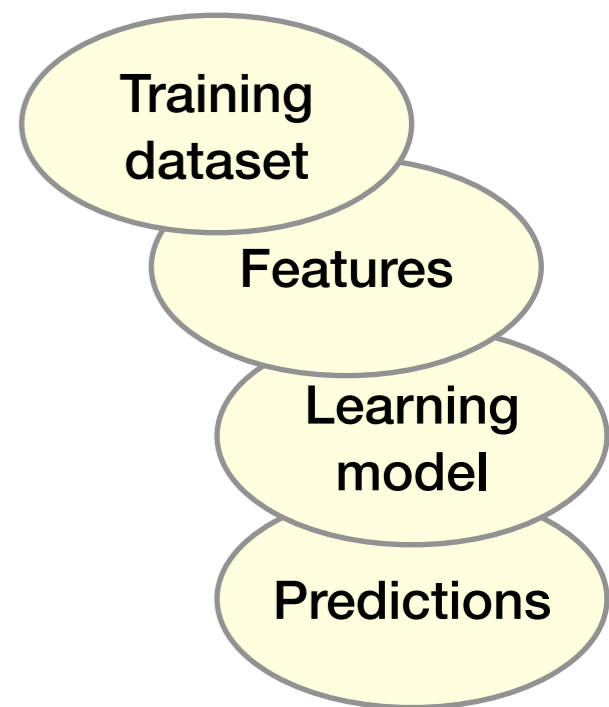
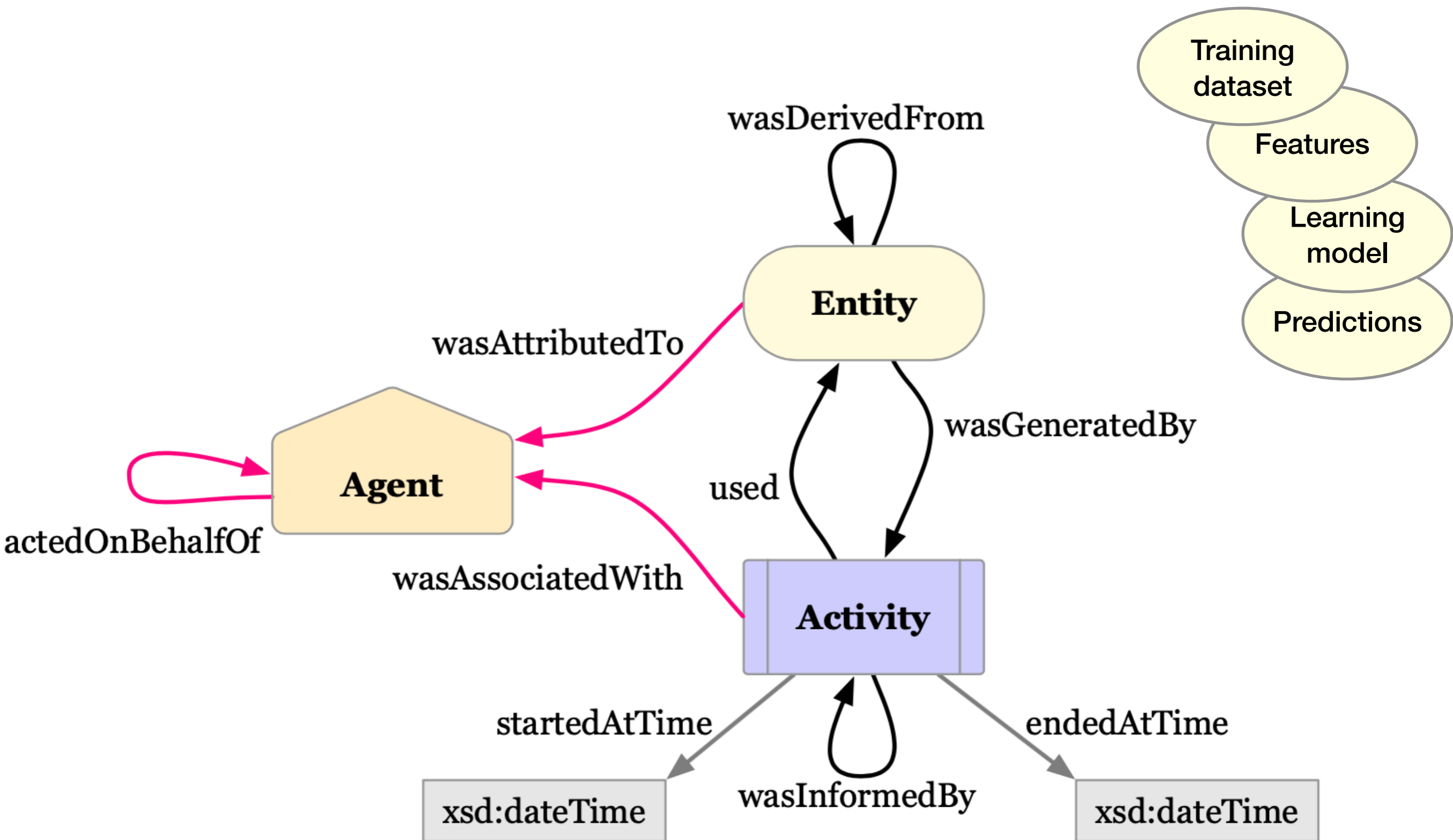
[Timothy Lebo](#), Rensselaer Polytechnic Institute, USA  
[Satya Sahoo](#), Case Western Reserve University, USA  
[Deborah McGuinness](#), Rensselaer Polytechnic Institute, USA

**Contributors:**

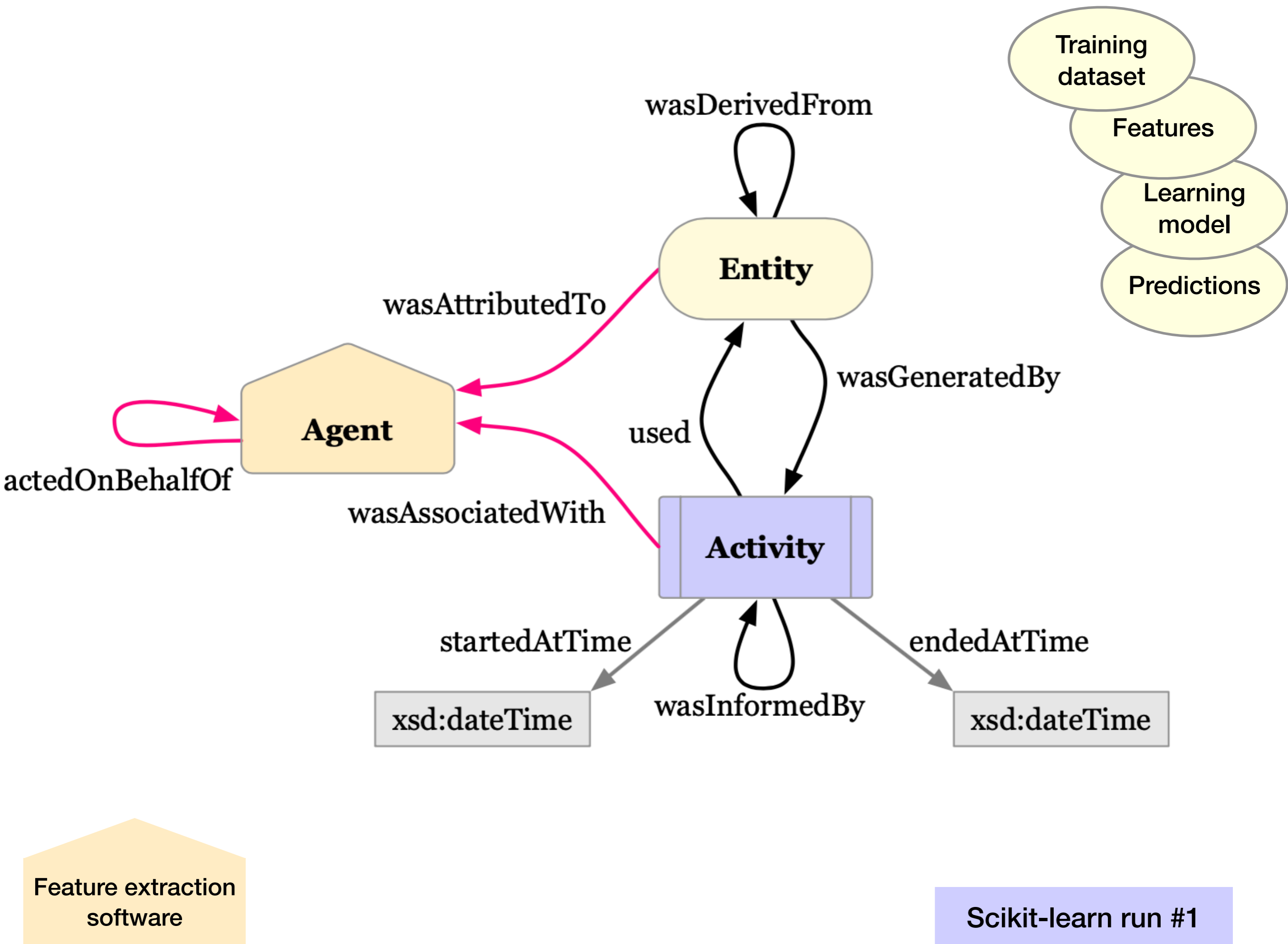
(In alphabetical order)  
[Khalid Belhajjame](#), University of Manchester, UK  
[James Cheney](#), University of Edinburgh, UK  
[David Corsar](#), University of Aberdeen, UK  
[Daniel Garijo](#), Ontology Engineering Group, Universidad Politécnica de Madrid, Spain  
[Stian Soiland-Reyes](#), University of Manchester, UK  
[Stephan Zednik](#), Rensselaer Polytechnic Institute, USA  
[Jun Zhao](#), University of Oxford, UK

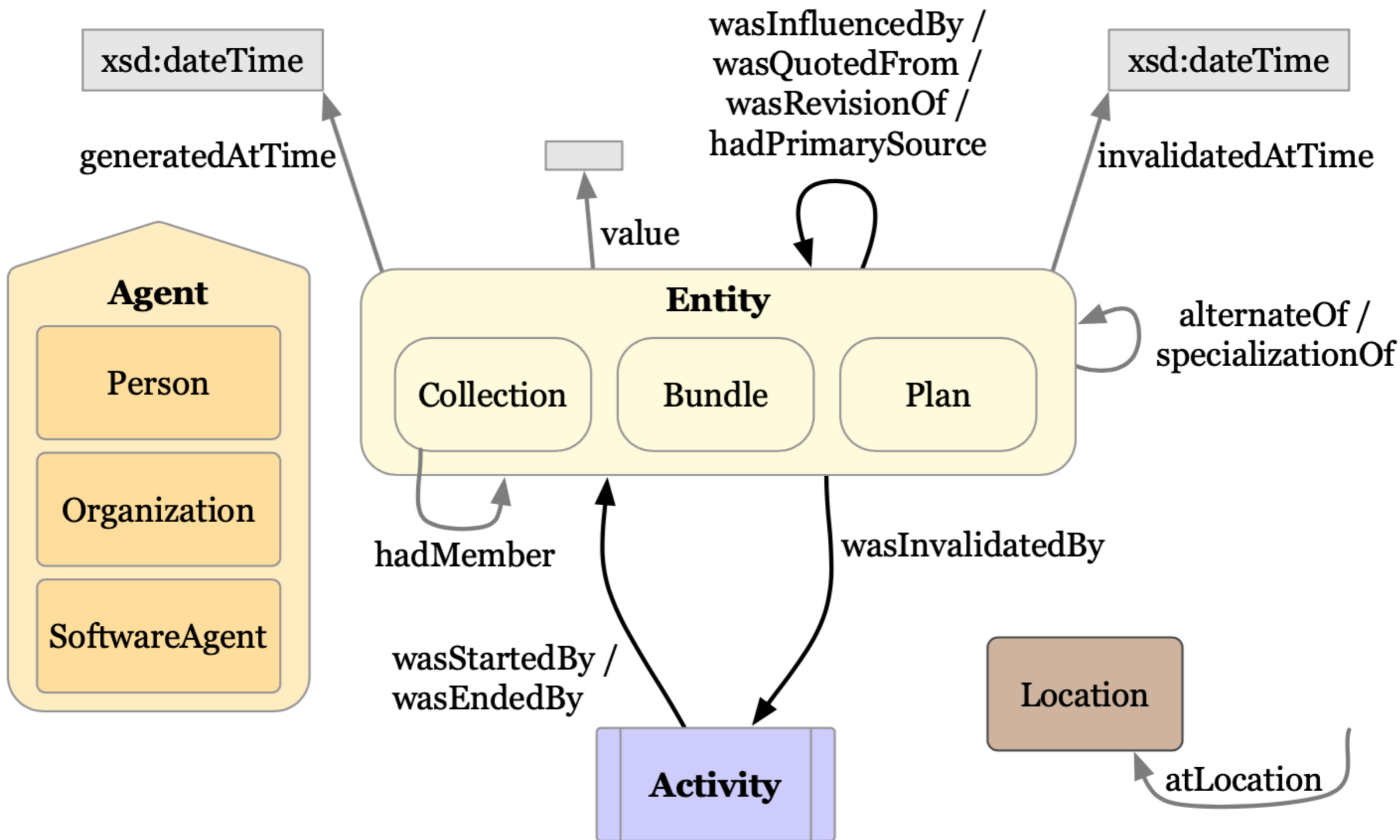






Scikit-learn run #1







# Reasoning with provenance



## 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.

# Reasoning with provenance

## 5.3 Derivations

Derivations with explicit activity, generation, and usage admit the following inference:

### Inference 11 (derivation-generation-use-inference)

In this inference, none of `a`, `gen2` or `use1` can be placeholders -.

**IF** `wasDerivedFrom(_id; e2,e1,a,gen2,use1,_attrs)`, **THEN** there exists `_t1` and `_t2` such that `used(use1; a,e1,_t1,[])` and `wasGeneratedBy(gen2; e2,a,_t2,[])`.

A revision admits the following inference, stating that the two entities linked by a revision are also alternates.

### Inference 12 (revision-is-alternate-inference)

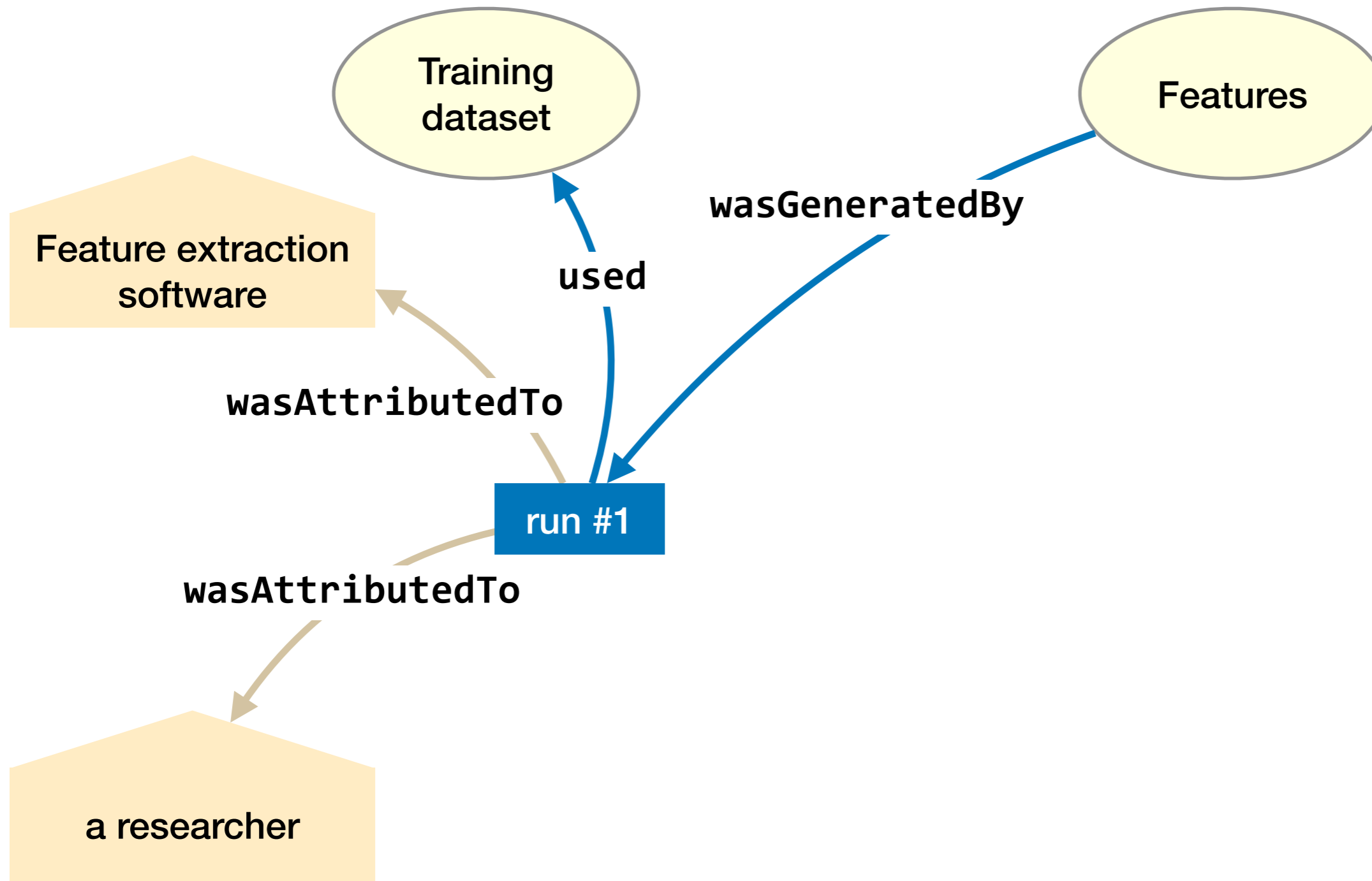
In this inference, any of `_a`, `_g` or `_u` **MAY** be placeholders.

**IF** `wasDerivedFrom(_id; e2,e1,_a,_g,_u,[prov:type='prov:Revision'])`, **THEN** `alternateOf(e2,e1)`.

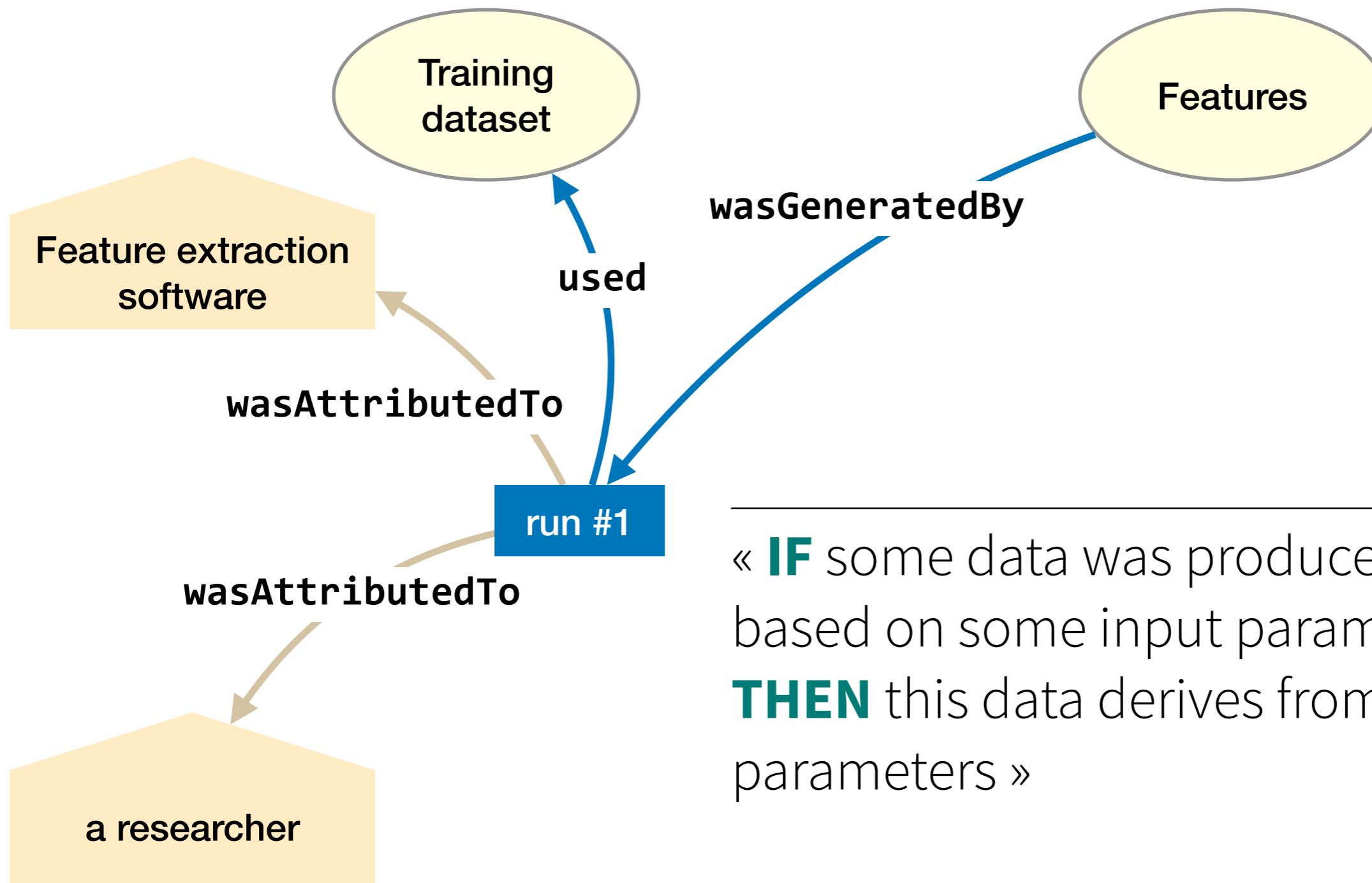
#### Remark

There is no inference stating that `wasDerivedFrom` is transitive.

# Reasoning with provenance



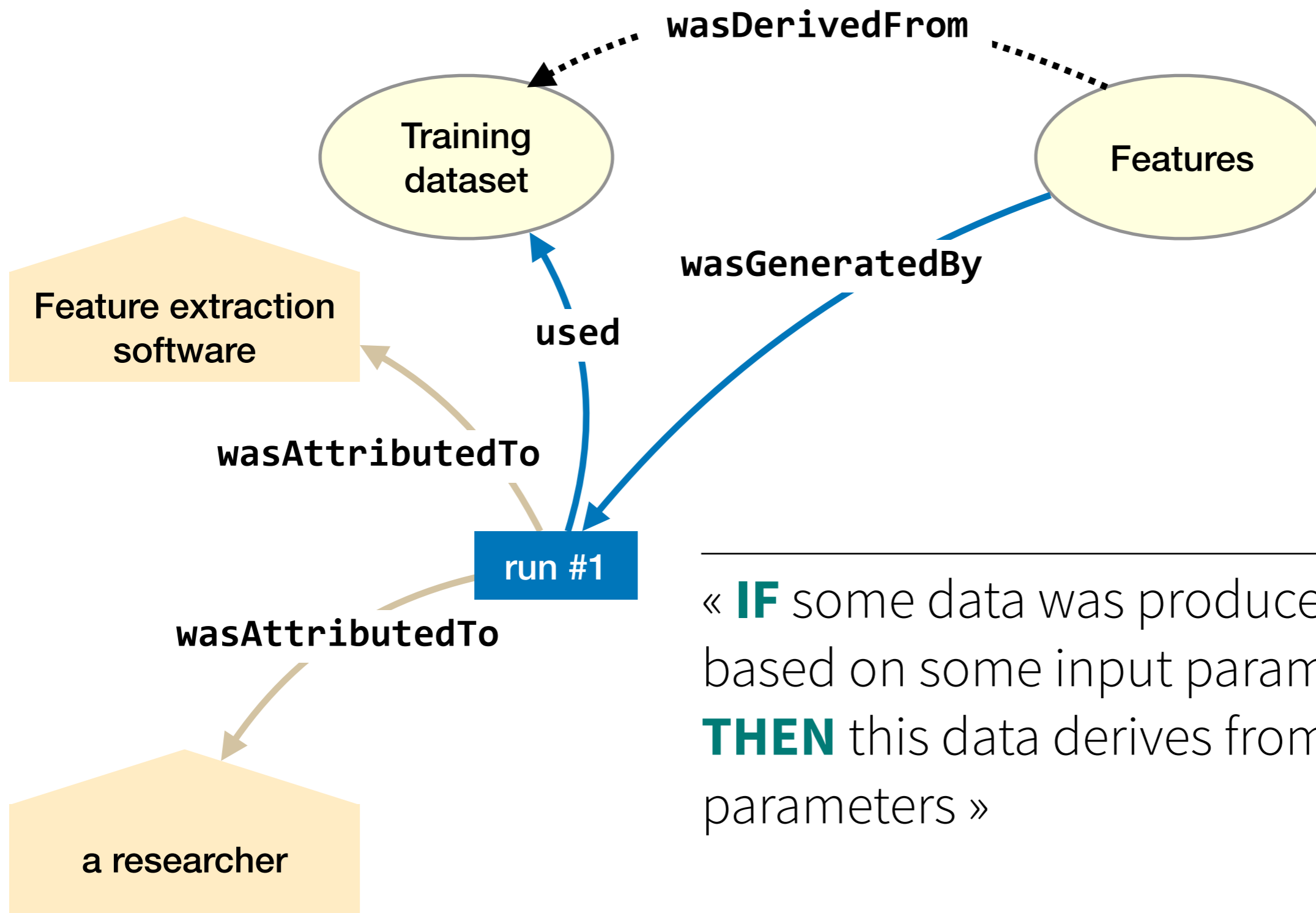
# Reasoning with provenance



---

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

# Reasoning with provenance



---

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



# Is provenance enough for reuse ?

Too fine-grained

No domain concepts

```
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;
```

Visualise

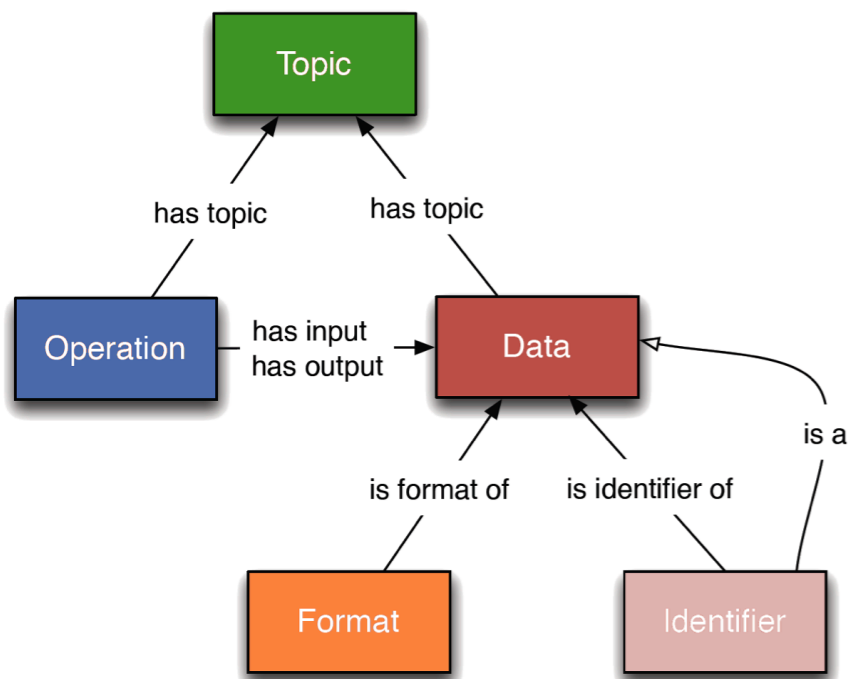


## 2. Bioinformatics

**domain-specific** concepts



# EDAM ontology



**BioPortal** Ontologies Search Annotator Recommender Mappings

## EDAM - Bioscientific data analysis ontology

Last uploaded: February 26, 2021

Summary **Classes** Properties Notes Mappings Widgets

Jump to:

- Data
- DeprecatedClass
- Format
- Operation
  - Alignment
  - Analysis
    - Enrichment analysis
    - Expression analysis
    - Genetic variation analysis
      - Collapsing methods
      - Genotyping
      - Reference identification
      - Structural variation detection
      - Variant calling**
        - Frameshift detection
        - Indel detection
        - SNP detection
        - Variant classification
        - Variant pattern analysis
        - Variant prioritisation
    - Image analysis
    - Network analysis

Details	Visualization	Notes ( 0 )	Class Mappings ( 3 )
Preferred Name	Variant calling		
Synonyms	Mutation detection Genome variant detection Somatic variant calling de novo mutation detection Germ line variant calling Allele calling Exome variant detection Variant mapping		
Definitions	Somatic variant calling is the detection of variations established by comparing the fluorescent traces produced by DNA sequencing to a reference genome. Variant detection Methods often utilise a database of aligned reference nucleotide polymorphisms, short indels and structural variants.		
ID	http://edamontology.org/operation_3227		

# bio.tools (semantic) catalog

**JASPAR** (biotools:jaspar) ID Verified  
<http://jaspar.genereg.net/>

Transcription factors and regulatory sites > Gene regulation > Genomics > Human biology > Plant biology > Model organisms >

Mature CC-BY-NC-4.0 Free of charge Open access

Web API Web application Database portal Python

The high-quality transcription factor binding profile database.

Transcription factor name > →  
Transcription factor identifier > →  
Taxon > →  
Family name > →  
Species name > →  
UniProt ID > →  
Text > →

DNA sequence > (FASTA >) →  
JASPAR profile ID > →

Database search > → JASPAR profile ID > (JASPAR format >, JSON >, Textual format >, CSV >, TSV >, YAML >)

Transcription factor binding site prediction > → JASPAR profile ID >

**Credits & Support**

- Albin Sandelin  
Primary contact | [albin@binf.ku.dk](mailto:albin@binf.ku.dk) | [Link](#)
- Boris Lenhard  
Primary contact | [b.lenhard@imperial.ac.uk](mailto:b.lenhard@imperial.ac.uk) | [Link](#)
- Wyeth Wasserman  
Primary contact | [wyeth@cmm.ubc.ca](mailto:wyeth@cmm.ubc.ca) | [Link](#)
- Anthony Mathelier  
Primary contact | [anthony.mathelier@ncmm.uio.no](mailto:anthony.mathelier@ncmm.uio.no) | [Link](#) | [ORCID](#)

**Documentation**

- <http://jaspar.genereg.net/docs/>  
General
- <http://jaspar.genereg.net/faq/>  
FAQ
- <http://jaspar.genereg.net/api/v1/docs/>  
API documentation

**Downloads**

- [Downloads page](#)

**Links**

- [https://twitter.com/jaspar\\_db](https://twitter.com/jaspar_db)  
Social media
- <https://bitbucket.org/CBGR/jaspar/src/master/>  
Repository

**Publication details**

38 586

Primary  
DOI: [10.1093/nar/gkx1126](https://doi.org/10.1093/nar/gkx1126)  
**JASPAR 2018: Update of the open-access database of transcription factor binding profiles and its web framework**  
Khan A. Fornes O. Stigliani A. Gheorghe M. Castro-Mondragon I.A. Van Der Lee R. Bessy A. Cheneby J. Kulkarni S.R. Tan G. Baranasic D. Arenillas D.I. Sandelin A. Vandepoele K.

Annotate workflow data with EDAM concepts ?

3. From knowledge graph to  
**data summaries**

# SPARQL queries to produce data summaries

jupyter FRESH-notebook (autosaved)



Visit repo

Copy Binder link

File Edit View Insert Cell Kernel Widgets Help

Not Trusted

Python 3

Code Download GitHub Binder

Memory: 351.3 MB / 2 GB

## 2. Human-oriented data summaries

### Sentence-based data explanations

Here, the goal is to describe a piece of data with the consensual definition of **what** does the tool that generates this piece of data.

Technically, this is done with a SPARQL query that combine the provenance information (*prov:wasGeneratedBy*), the description of the tool (*biotools:has\_function*), and the domain knowledge on the nature of the processing (*oboInOwl:hasDefinition*).

```
In [6]: %%time
query = """
SELECT ?d_label ?title ?f_def ?st WHERE {
  ?d rdf:type prov:Entity ;
  prov:wasGeneratedBy ?x ;
  prov:wasAssociatedWith ?tool ;
  rdfs:label ?d_label .

  ?tool dc:title ?title ;
  biotools:has_function ?f .

  ?f rdfs:label ?f_label ;
  oboInOwl:hasDefinition ?f_def .

  ?c rdf:type mp:Claim ;
  mp:statement ?st .
}
"""

results = g.query(query)
for r in results :
  display(Markdown('The file `` + str(r['d_label']) + ``
                  + ' **results from tool** ' + str(r['title'])
                  + ' **which** ' + str(r['f_def'])))
  display(Markdown(' **It was produced in the context of** ' + str(r['st']) ))
```

launch binder

<https://github.com/albangaingard/fresh-toolbox>

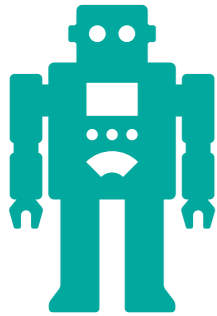
The file `VCF/hapcaller.snv.recal.filter.vcf.gz.tbi` **results from tool** `gatk2_variant_filtration-IP` **which** Analyse a genetic variation, for example to annotate its location, alleles, classification, and effects on individual transcripts predicted for a gene model.

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

The file `VCF/hapcaller.snv.recal.filter.vcf.gz.tbi` **results from tool** `gatk2_variant_filtration-IP` **which** Filter a set of files or data items according to some property.

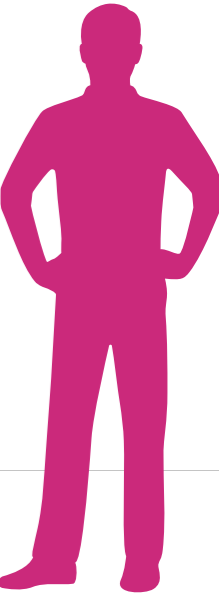
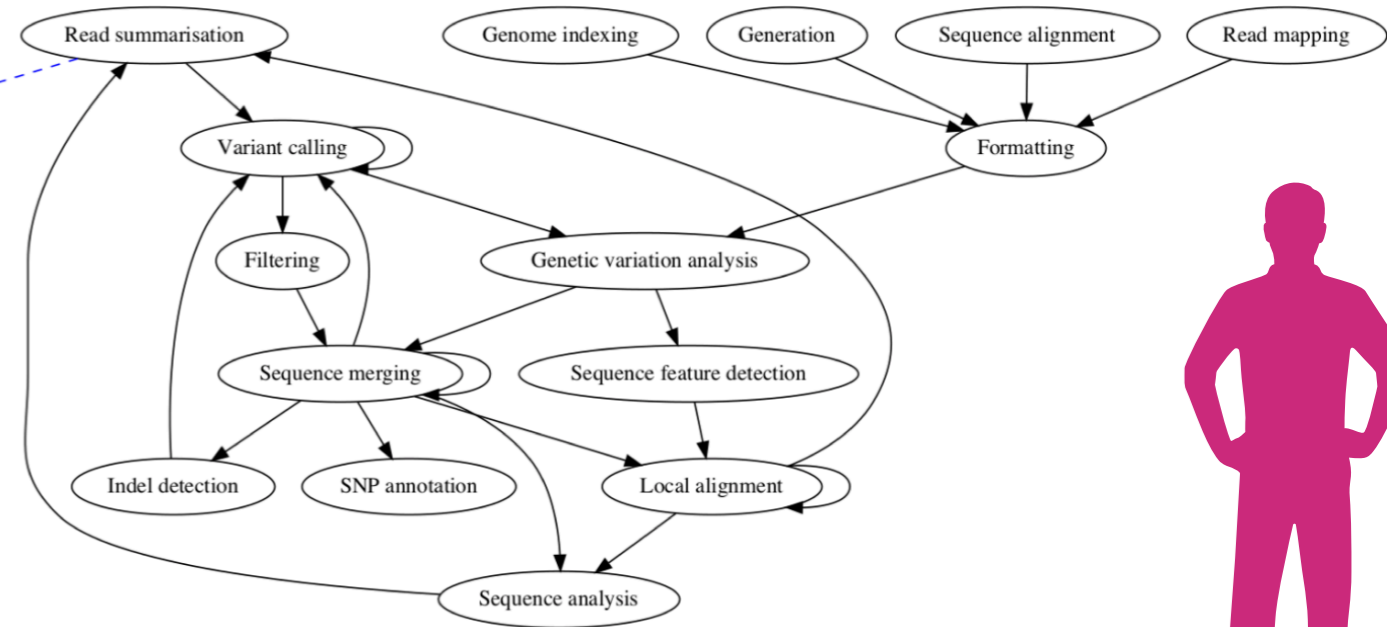
**It was produced in the context of** Rare Coding Variants in ANGPTL6 Are Associated with Familial Forms of Intracranial Aneurysm

# Data summaries



```
[...]  
:head {  
  _:np1 a np:Nanopublication .  
  _:np1 np:hasAssertion :assertion .  
  _:np1 np:hasProvenance :provenance .  
  _:np1 np:hasPublicationInfo :pubInfo .  
}  
  
:assertion {  
  <http://snakemake-provenance/Samples/Sample1/  
  BAM/Sample1.merged.bai> rdfs:seeAlso  
  <http://edamontology.org/operation_3197> .  
  
  <http://snakemake-provenance/VCF/hapcaller.  
  indel.recal.filter.vcf.gz> rdfs:seeAlso  
  <http://edamontology.org/operation_3695> .  
}  
[...]
```

Samples/Sample1/BAM/Sample1.final.bam



```
...  
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  
...
```

1. It's possible to automatically produce **machine-oriented nanoPublications**

2. It's possible to automatically display the **typical bioinformatics tasks** data originate from
3. It's possible to document data with **text** leveraging ontology definitions (EDAM)

# Perspectives

# Take home message & perspectives

**Scientific Workflows** → automation, abstraction, provenance

Standards for **provenance representation** and **reasoning** (PROV-O)

Contributions

- feeding a **knowledge graph** with generic provenance metadata captured at runtime and domain knowledge (EDAM)
- generating **domain-specific** machine and human-oriented **data summaries**
- in line with **F - - R principles** : machine-readable data (F2), community standards (R1.3), provenance (R1.2)

Future works

- distributed data sources → **distributed provenance, reasoning**
- application to ML workflows → better interpret/explain predictions ?
- **evaluation** → large bioinformatics communities through the MuDiS4LS Equipex+ project federating HPC clusters, workflow developers and biologists

# Acknowledgments



Hala Skaf-Molli, LS2N,  
University of Nantes



Khalid Belhajjame,  
LAMSADE, University of  
Paris-Dauphine, PSL

contact: [alban.gaignard@univ-nantes.fr](mailto:alban.gaignard@univ-nantes.fr)

