

Des annotations de provenance, aux résultats de workflows “Linked Data”

Alban Gaignard

Nantes Academic Hospital
(CHU de Nantes), CNRS,
France

Hala Skaf-Molli

LINA,
Nantes University, CNRS,
France

Audrey Bihouée

Institut du Thorax, Nantes
University, INSERM, CNRS,
France

Travaux publiés à “Theory and Practice of Provenance” (TaPP’16 workshop), Washington DC



Needs for linked experiment reports

Motivations: reusing (massive) RNA-seq data

TopHat: algorithm to align multiple sequence reads to a reference genome (known genes).

Motivations: reusing (massive) RNA-seq data

TopHat: algorithm to align multiple sequence reads to a reference genome (known genes).

	1 sample
Input data	2 x 17 Gb
1-core CPU	170 hours
32-cores CPU	32 hours
Output data	12 Gb

Motivations: reusing (massive) RNA-seq data

TopHat: algorithm to align multiple sequence reads to a reference genome (known genes).

	1 sample	300 samples
Input data	2 x 17 Gb	10.2 Tb
1-core CPU	170 hours	5.9 years
32-cores CPU	32 hours	14 months
Output data	12 Gb	3.6 Tb

Motivations: reusing (massive) RNA-seq data

TopHat: algorithm to align multiple sequence reads to a reference genome (known genes).

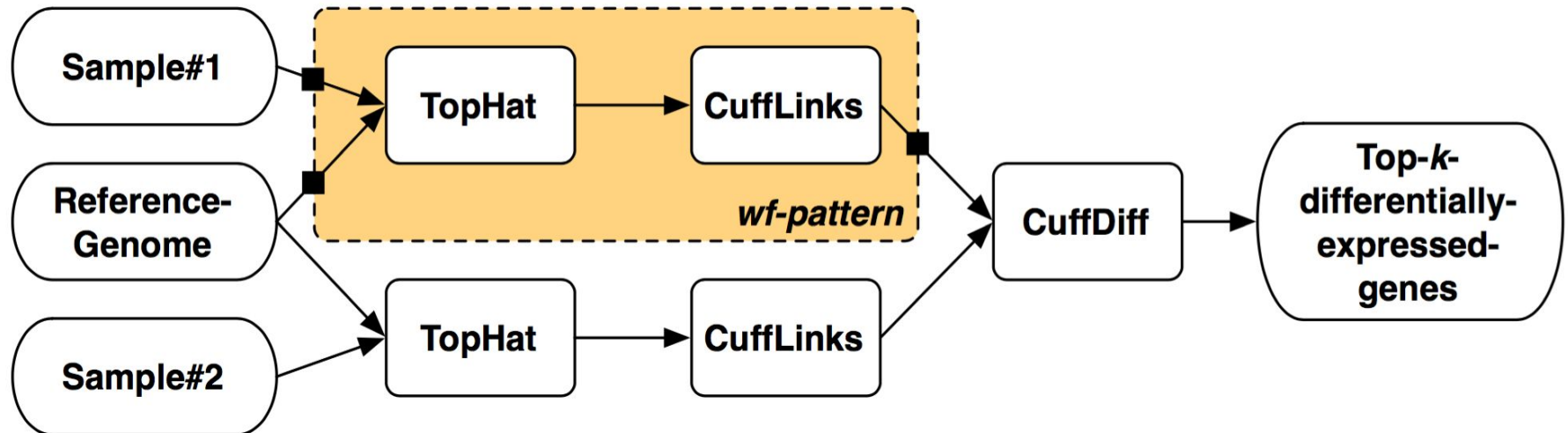
	1 sample	300 samples
Input data	2 x 17 Gb	10.2 Tb
1-core CPU	170 hours	5.9 years
32-cores CPU	32 hours	14 months
Output data	12 Gb	3.6 Tb

Challenges

Algorithmic performance, storage, preservation,
reuse (limit recompute) & share.

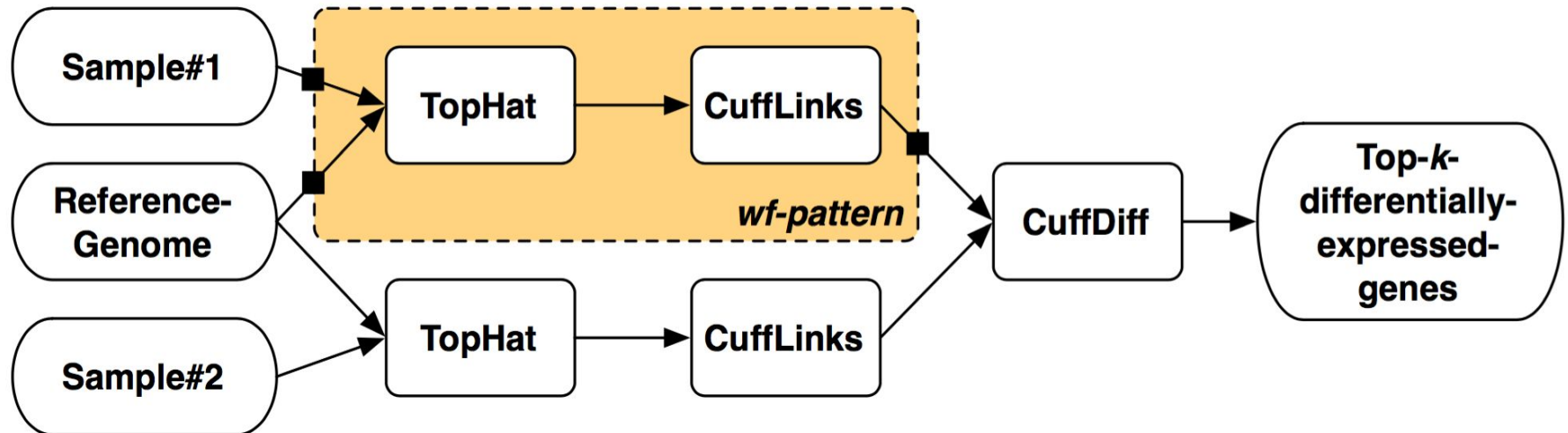
Motivations: reusing experiment results

Scientific experiment: RNA sequencing to quantify gene expression levels under multiple biological conditions.



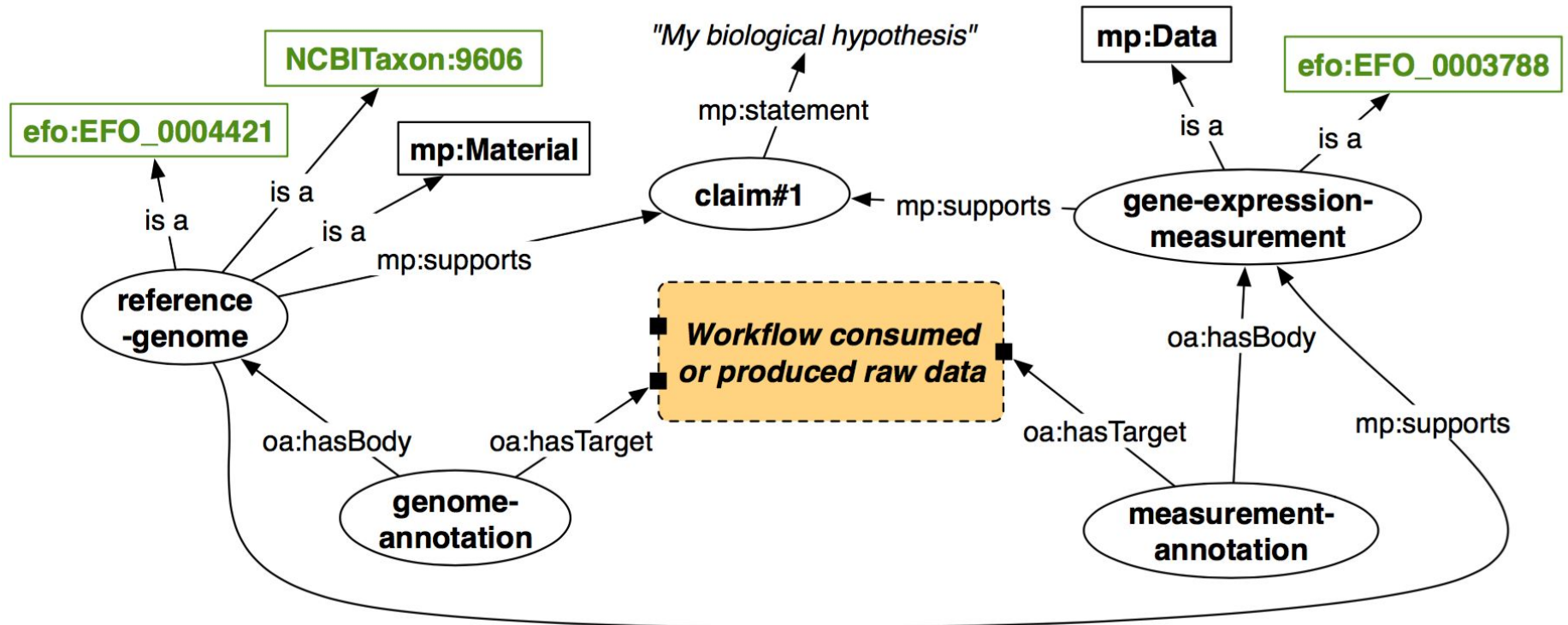
Motivations: reusing experiment results

Scientific experiment: RNA sequencing to quantify gene expression levels under multiple biological conditions.



**Interpreting and sharing results :
need for scientific context (metadata)**

Expected result: human+machine tractable reports



Annotated "Material & Methods"

Links to **some** workflow artifacts (algorithms, data)

5-star Linked Open Data

W3C standards for **machine** and **human** readable data on the web.

★★★★★ : time and expertise !



5-star Linked Open Data

W3C standards for **machine** and **human** readable data on the web.

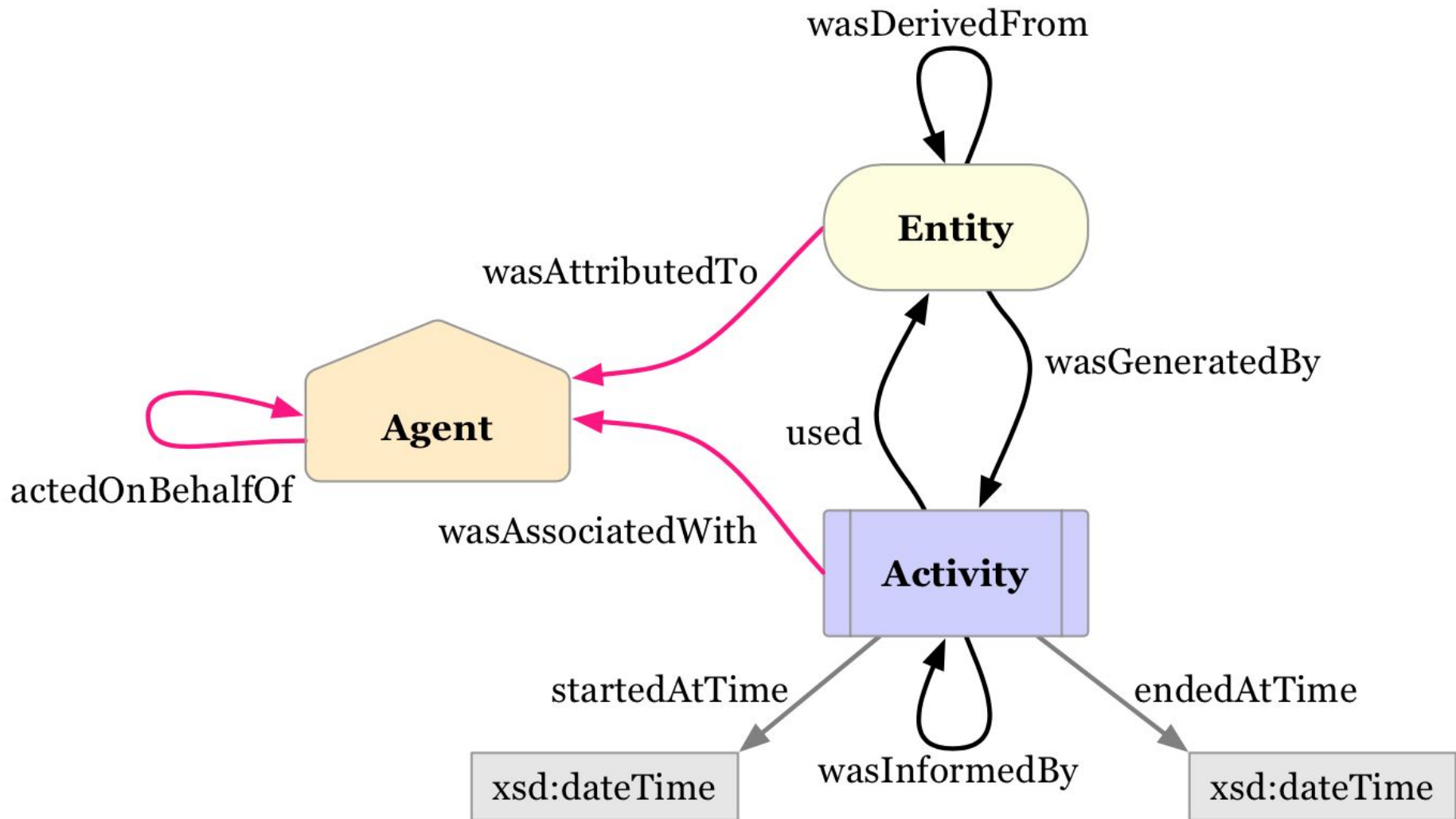
★★★★★ : time and expertise !

How to ease this process ?

- Workflow engines → **automation**
- PROV → **workflow runs as linked data**



W3C Recommendation : the PROV Ontology



PROV traces for a WF run

```
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



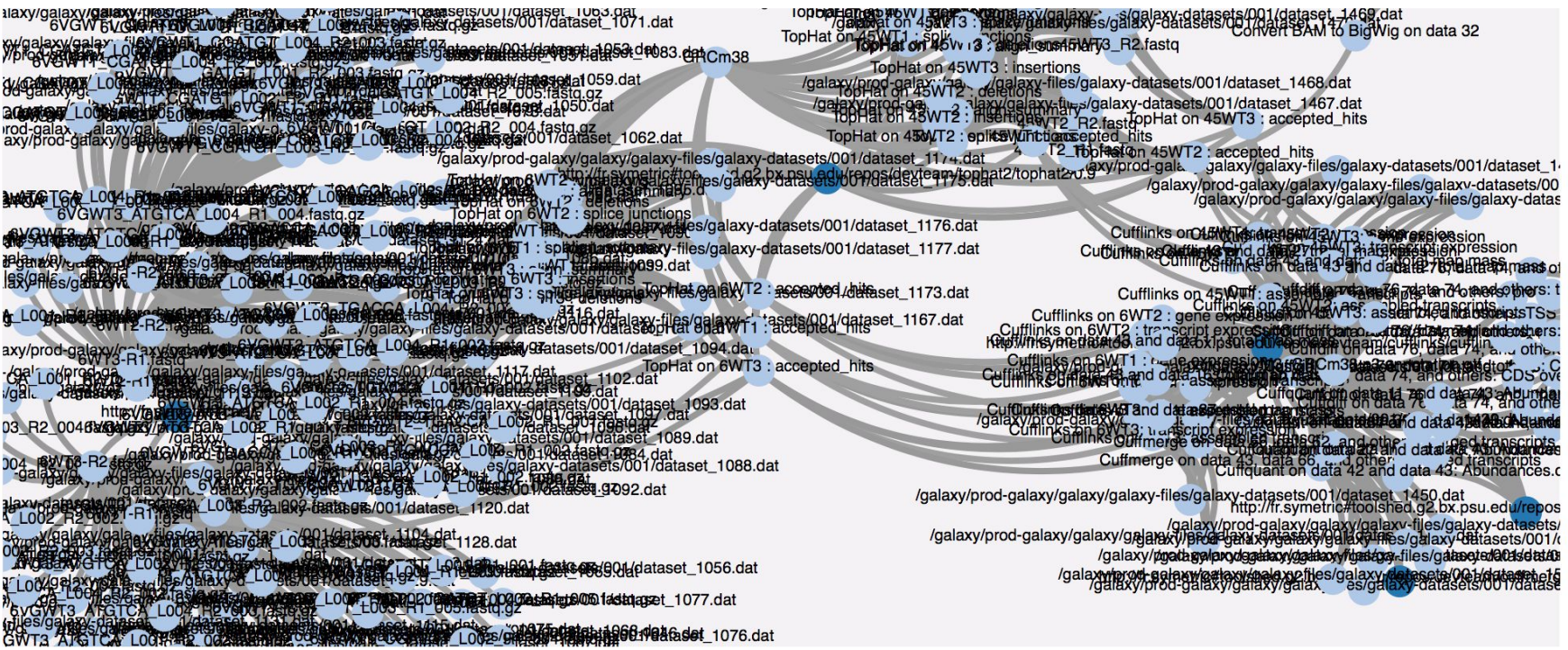
PROV traces for a WF run

```
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;
```

granularity ?

domain concepts ?

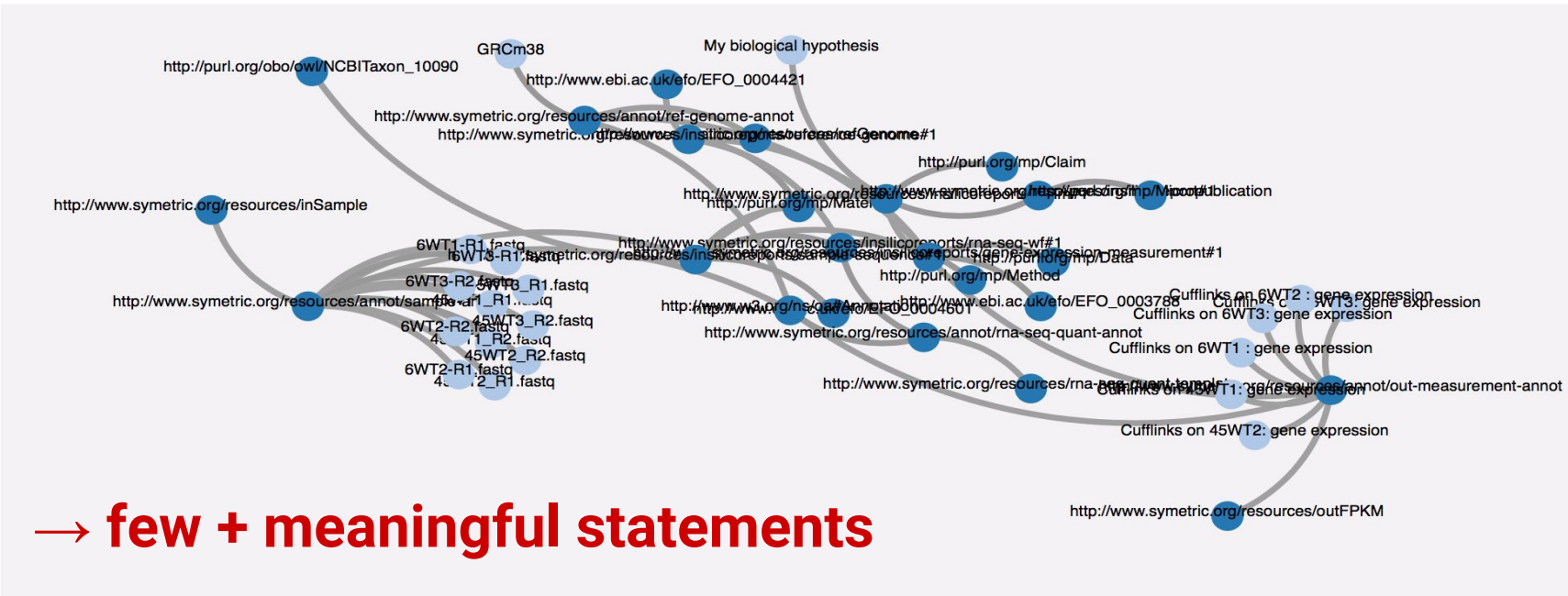
Visualise



Provenance as a Linked Experiment Report

```
40 <http://www.symetric.org/resources/insilicoreports/claim#1>
41   a          mp:Claim ;
42   mp:statement "My biological hypothesis" .
43
44 <http://www.symetric.org/resources/insilicoreports/rna-seq-wf#1>
45   a          mp:Method ;
46   mp:supports <http://www.symetric.org/resources/insilicoreports/gene-expression-measurement#1> .
47
48 <http://www.symetric.org/resources/annot/ref-genome-annot>
49   a          oa:Annotation ;
50   oa:hasBody <http://www.symetric.org/resources/insilicoreports/gene-expression-measurement#1> ;
51   oa:hasTarget "GRCm38" ;
52   oa:hasTarget <http://www.symetric.org/resources/refGenome> .
53
54 <http://www.symetric.org/resources/annot/rna-seq-quant-annot>
```

Visualise



→ few + meaningful statements

Problem statement

Scientific workflows produce massive raw results. Their publication into curated query-able linked data repositories requires lot of time and expertise.

Can we exploit provenance traces to ease the publication of scientific results as Linked Data ?

Problem statement & objectives

Problem statement

Scientific workflows produce massive raw results. Their publication into curated query-able linked data repositories requires lot of time and expertise.

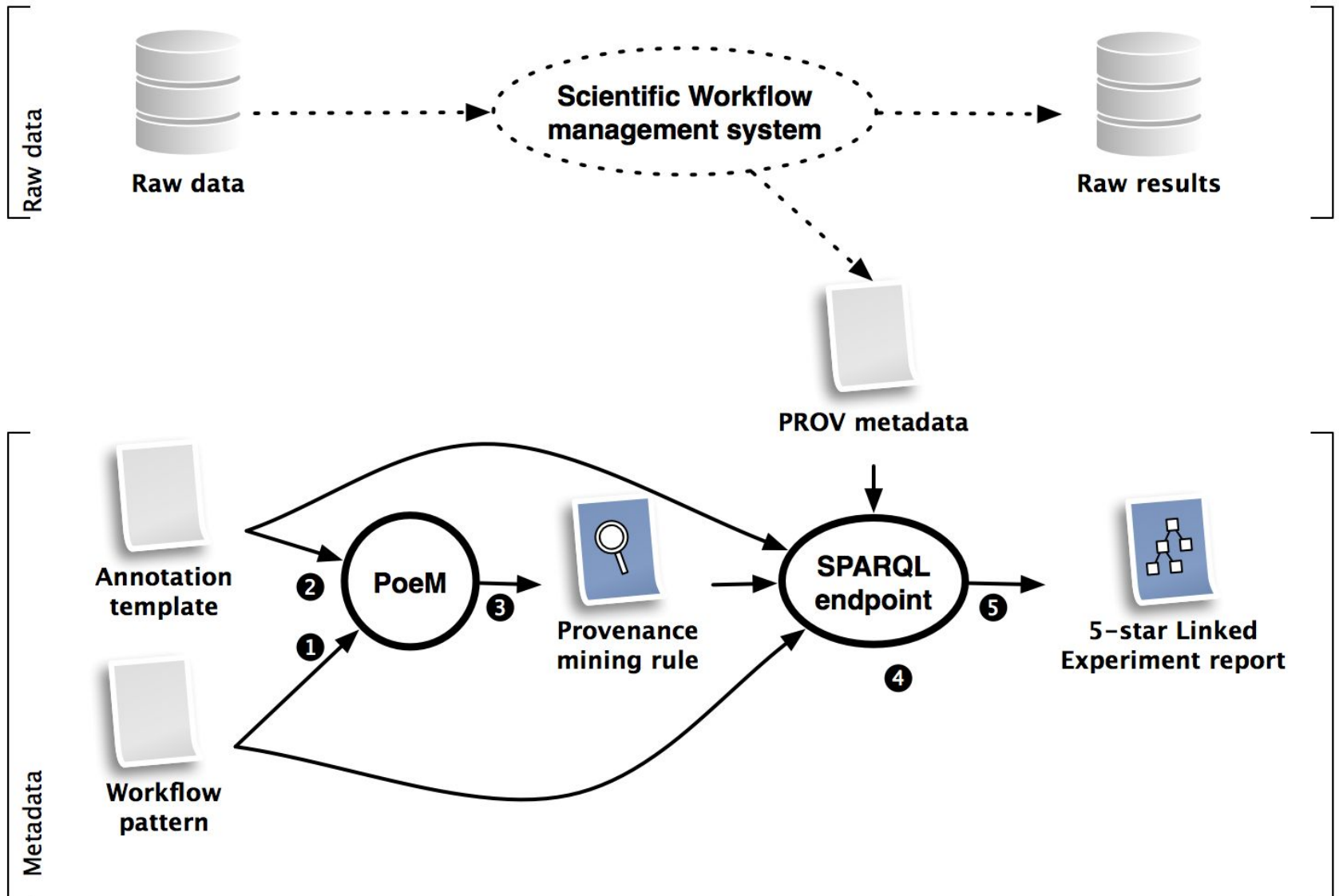
Can we exploit provenance traces to ease the publication of scientific results as Linked Data ?

Objectives

- (1) Leverage **annotated workflow patterns** to generate **provenance mining rules**.
- (2) Refine **provenance** traces into **linked experiment reports**.

Rules generation

Approach



Input domain-specific annotations (1,2)

Workflow patterns ①

Sequence patterns, with possibly intermediate steps

- P-PLAN ontology: *Step, Variable, hasInputVar, hasOutputVar*
- EDAM ontology: *hasFunction, RNA sequence, Genome map*

Input domain-specific annotations (1,2)

Workflow patterns ①

Sequence patterns, with possibly intermediate steps

- P-PLAN ontology: *Step, Variable, hasInputVar, hasOutputVar*
- EDAM ontology: *hasFunction, RNA sequence, Genome map*

Experiment report template ②

Link scientific claims, statements, material and methods

- MicroPublication ontology: *Material, Method, Claims*
- Experimental factor ontology: *Transcriptome, Gene expression*
- NCBI taxonomy: *Homo Sapiens*
- Open Annotation model: *hasBody, hasTarget*

PoeM: generating Provenance Mining rules ③

Input : W : Workflow annotated pattern ①,
 S_1 : First step of W ,
 S_2 : Last step of W ,
 A : Annotation template ②.

Output: $Rule$: Provenance mining rule.

1 **begin**

2 $IN_{S_1} \leftarrow getInputs(S_1)$

3 $OUT_{S_2} \leftarrow getOutputs(S_2)$

4

5 $provGraph \leftarrow genDataLineage(OUT_{S_2}, IN_{S_1})$

(SPARQL Property Path)

6 $reportGraph \leftarrow bindReportTargets(provGraph, A)$

(SPARQL Basic Graph Pattern)

7

8 $Rule \leftarrow \frac{provGraph.edge_1 \wedge \dots \wedge provGraph.edge_N}{reportGraph}$ (SPARQL Construct query)

PoeM: sample generated rule 3

```
1 PREFIX ...
2
3 CONSTRUCT {
4
5     <http://.../insilicoreport#1> rdf:type <http://purl.org/mp/Micropublication> .
6 # ...
7     <http://.../out-measurement-annot> rdf:type <http://www.w3.org/ns/oa#Annotation> .
8     <http://.../out-measurement-annot> oa:hasTarget ?lout0 .
9     <http://.../out-measurement-annot> oa:hasBody <http://.../gene-expression-measurement#1> .
10    <http://.../gene-expression-measurement#1> rdf:type <http://purl.org/mp/Data> .
11    <http://.../gene-expression-measurement#1> rdf:type <http://www.ebi.ac.uk/efo/EFO_0003788> .
12    ...
13
14 } WHERE {
15
16 # ...
17     ?outS1 prov:wasGeneratedBy ?step1 .
18     ?in0 rdfs:label | ^rdfs:label | prov:wasDerivedFrom)* ?outS0 .
19
20     ?in2 prov:used ?in2 .
21     ?soft2 prov:wasAssociatedWith ?soft2 .
22     ?planStep2 a p-plan:Step .
23     ?planStep2 rdfs:label ?planStep2Label .
24     FILTER (contains(lcase(str(?soft2)),lcase(str(?planStep2Label)))) .
25
26     ?out0 prov:wasGeneratedBy ?step2 .
27     ?out0 rdfs:label ?lout0 .
28     FILTER (contains(lcase(str(?lout0)), lcase("gene expression")) .
29 }
30
```

<Then> part

<If> part

First experiments & results

Material & methods

- Real-life RNA-seq workflow to study 3 mice populations
- WF implemented in Galaxy, run on 2 biological samples
- PROV traces exported from Galaxy Histories (API)

Material & methods

- Real-life RNA-seq workflow to study 3 mice populations
- WF implemented in Galaxy, run on 2 biological samples
- PROV traces exported from Galaxy Histories (API)

Results (for 1 biological sample)

- 60h CPU (12 cores for genome alignment), 21Gb storage
- 3s to export 81 PROV triples from the Galaxy history
- 2s to apply the rule and produce 35 Micropublication triples

Results (for 4 biological samples)

- 1232 PROV triples from the Galaxy history
- 49 Micropublication triples (4%)

Semi-automated approach

- (1) PoeM generates **semantic web rules**
- (2) PoeM rules applied on **PROV** traces to assemble **linked experiment reports** (MicroPublication)

Limitations:

- Sequence workflow patterns only
- *SPARQL property paths* with complex WF patterns ?
- Syntactic matching between WF patterns and PROV labels

Usage scenarios:

- **Query** workflow datasets with domain concepts
- **Populate** RDF repositories with WF results

Future works

- (1) **WF patterns**: split-merge, “common motifs”
- (2) **Genericity**: other domains / other reports (RO, Nanopub.)
- (3) **PROV heterogeneity**: multi-systems PROV reconciliation
- (4) **Evaluation**: involving biologists, at larger scale

Questions ?

Demo: <http://poem.univ-nantes.fr>

Contact: alban.gaignard@univ-nantes.fr

Acknowledgments



BiRD bioinformatics facility



Connect Talent Call

Experiment context

SyMeTRIC: **systems medicine** project (2015-2017, call “Connect Talent”), funded by the french region Pays de la Loire.

