

# Biomedical provenance metadata to support FAIR genomic research data: *an intracranial aneurysms use case*

BC2 workshop T7  
8th of September 2025

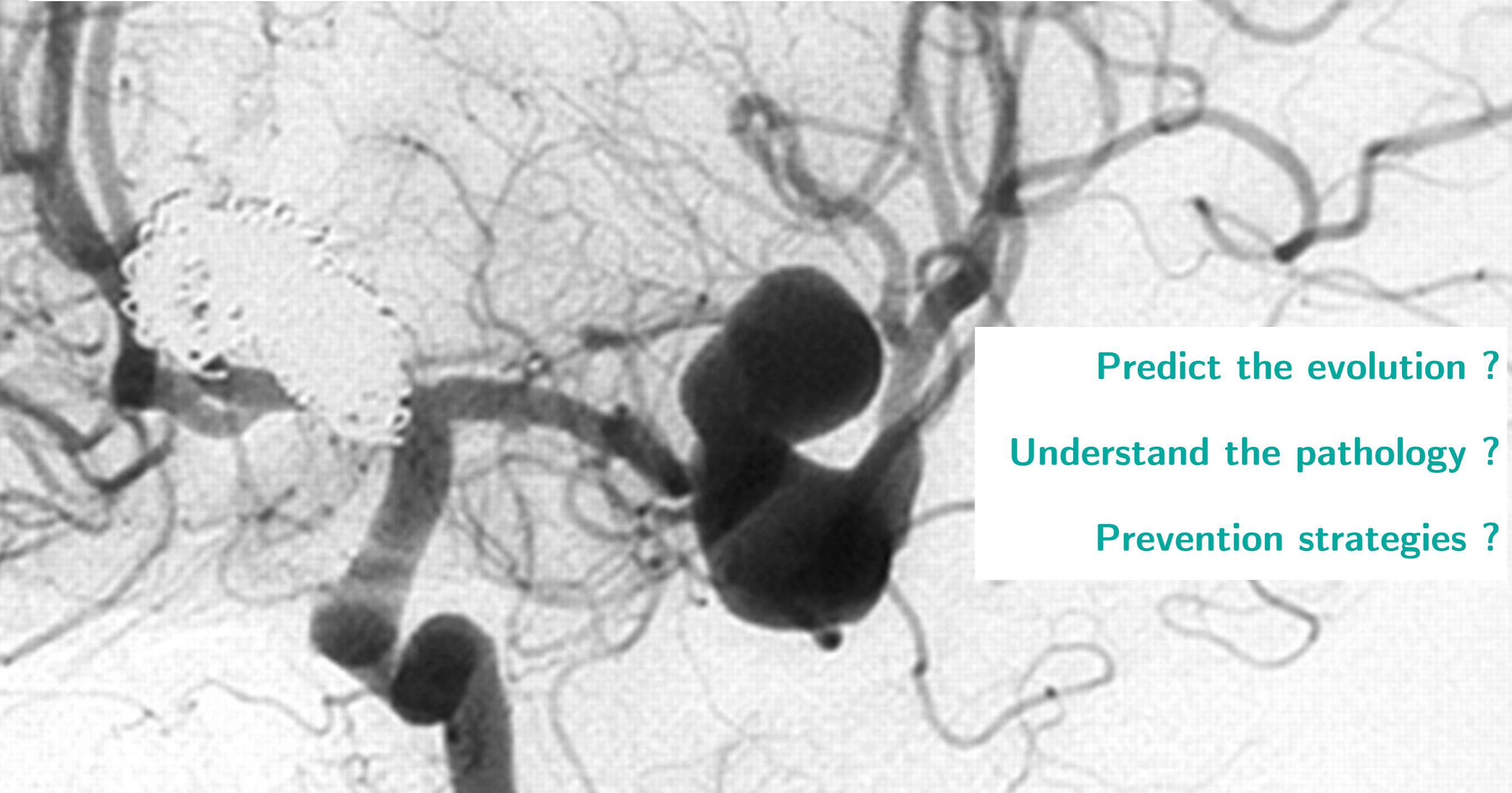
Alban Gaignard, Institut du Thorax, CNRS, Nantes, France

# Biomedical context

ICAN cohort: 34 univ. hospitals / 3000 subjects

3.000 ToF MRIs

600 whole genomes (analysis in progress)



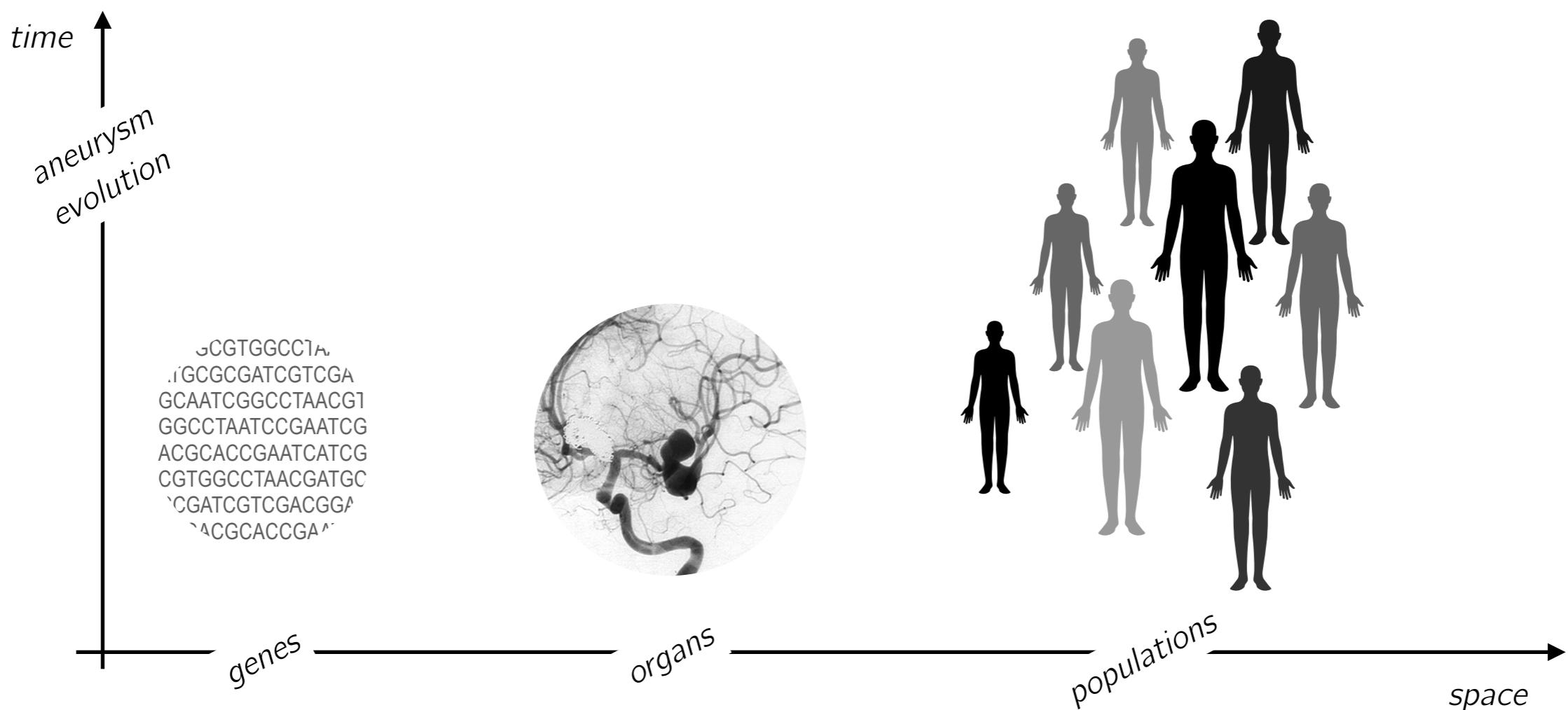
Predict the evolution ?

Understand the pathology ?

Prevention strategies ?

# Multi-factorial disease → multi-scale data

- ▶ Intracranial aneurysms: a complex & multifactorial disease
- ▶ Inter-disciplinary efforts needed for a better understanding of the pathology
- ▶ Specific data produced at very specific scales



# Multi-modal data analysis

Using deep learning for an automatic detection and classification of the vascular bifurcations along the Circle of Willis

Rafic Nader<sup>a</sup>, Romain Bourcier<sup>a</sup>, Florent Autrusseau<sup>a,b,\*</sup>

<sup>a</sup> Nantes Université, CHU Nantes, CNRS, INSERM, l'institut du thorax, F-44000 Nantes, France

<sup>b</sup> Nantes Université, Polytech'Nantes, LTeN, U-6607, Rue Cl. Pauc, 44306, Nantes, France

## ARTICLE INFO

### MSC:

41A05

41A10

65D05

65D17

### Keywords:

Vascular bifurcations  
Circle of Willis  
Deep learning

## ABSTRACT

Most of the intracranial aneurysms (ICAs) occur on a specific portion of the cerebral vascular tree named the Circle of Willis (CoW). More particularly, they mainly arise onto fifteen of the major arterial bifurcations constituting this circular structure. Hence, for an efficient and timely diagnosis it is critical to develop some methods being able to accurately recognize each Bifurcation of Interest (BoI). Indeed, an automatic extraction of the bifurcations presenting the higher risk of developing an ICA would offer the neuroradiologists a quick glance at the most alarming areas. Due to the recent efforts on Artificial Intelligence, Deep Learning turned out to be the best performing technology for many pattern recognition tasks. Moreover, various methods have been particularly designed for medical image analysis purposes. This study intends to assist the neuroradiologists to promptly locate any bifurcation presenting a high risk of ICA occurrence. It can be seen as a Computer Aided Diagnosis scheme, where the Artificial Intelligence facilitates the access to the regions of interest within the MRI. In this work, we propose a method for a fully automatic detection and recognition of the bifurcations of interest forming the Circle of Willis. Several neural networks architectures have been tested, and we thoroughly evaluate the bifurcation recognition rate.

### MRA-TOF

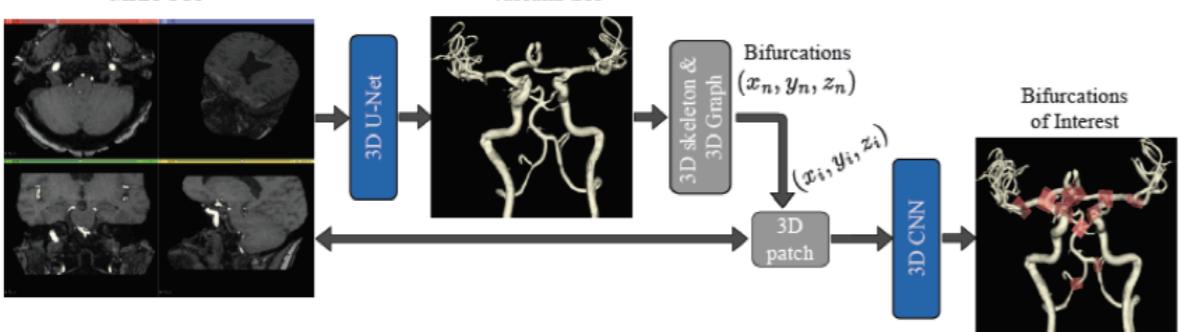


Fig. 1: General flowchart of the bifurcation recognition process.

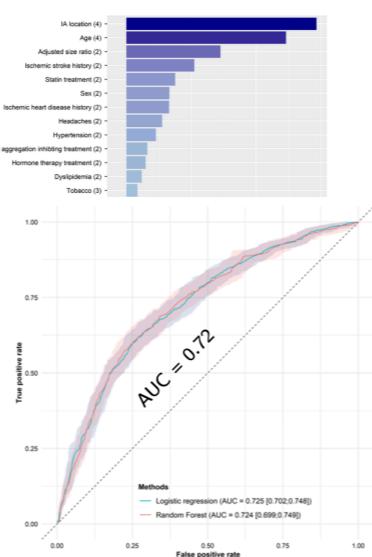
Observational Study > J Neurol Neurosurg Psychiatry. 2021 Feb;92(2):122-128.  
doi: 10.1136/jnnp-2020-324371. Epub 2020 Oct 23.

Location of intracranial aneurysms is the main factor associated with rupture in the ICAN population

Olivia Rousseau<sup>1</sup>, Matilde Karakachoff<sup>1</sup>, Alban Gaignard<sup>2</sup>, Lise Bellanger<sup>3</sup>, Philippe Bijlenga<sup>4</sup>, Pacôme Constant Dit Beaufils<sup>1</sup>, Vincent L'Allinec<sup>5</sup>, Olivier Levrier<sup>6</sup>, Pierre Aguettaz<sup>7</sup>, Jean-Philippe Desilles<sup>8</sup>, Caterina Michelozzi<sup>9</sup>, Gaultier Marnat<sup>10</sup>, Anne-Clémence Vion<sup>2</sup>, Gervaise Loirand<sup>2</sup>, Hubert Desal<sup>11</sup>, Richard Redon<sup>2</sup>, Pierre-Antoine Gourraud<sup>1</sup>, Romain Bourcier<sup>12</sup>; ICAN Investigators

Collaborators, Affiliations + expand

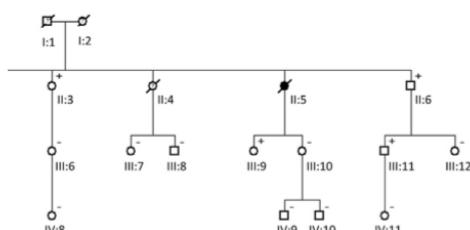
PMID: 33097563 DOI: 10.1136/jnnp-2020-324371



## Rare Coding Variants in ANGPTL6 Are Associated with Familial Forms of Intracranial Aneurysm

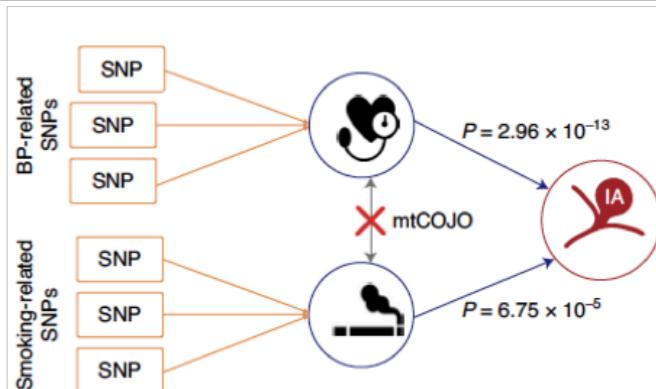
Romain Bourcier,<sup>1,2</sup> Solena Le Scouarnec,<sup>1</sup> Stéphanie Bonnaud,<sup>1,3</sup> Matilde Karakachoff,<sup>1,3</sup> Emmanuelle Bourcereau,<sup>3</sup> Sandrine Heurtelise-Chrétien,<sup>1</sup> Céline Menguy,<sup>1</sup> Christian Dina,<sup>1,3</sup> Floriane Simonet,<sup>1,3</sup> Alexis Moles,<sup>4</sup> Cédric Lenoble,<sup>2</sup> Pierre Lindenbaum,<sup>1</sup> Stéphanie Chatel,<sup>1,3</sup> Bertrand Isidor,<sup>5</sup> Emmanuelle Génin,<sup>6</sup> Jean-François Deleuze,<sup>7</sup> Jean-Jacques Schott,<sup>1,3</sup> Hervé Le Marec,<sup>1,3</sup> ICAN Study Group, Gervaise Loirand,<sup>1,3,8</sup> Hubert Desal,<sup>1,2,8,\*</sup> and Richard Redon<sup>1,3,8,\*</sup>

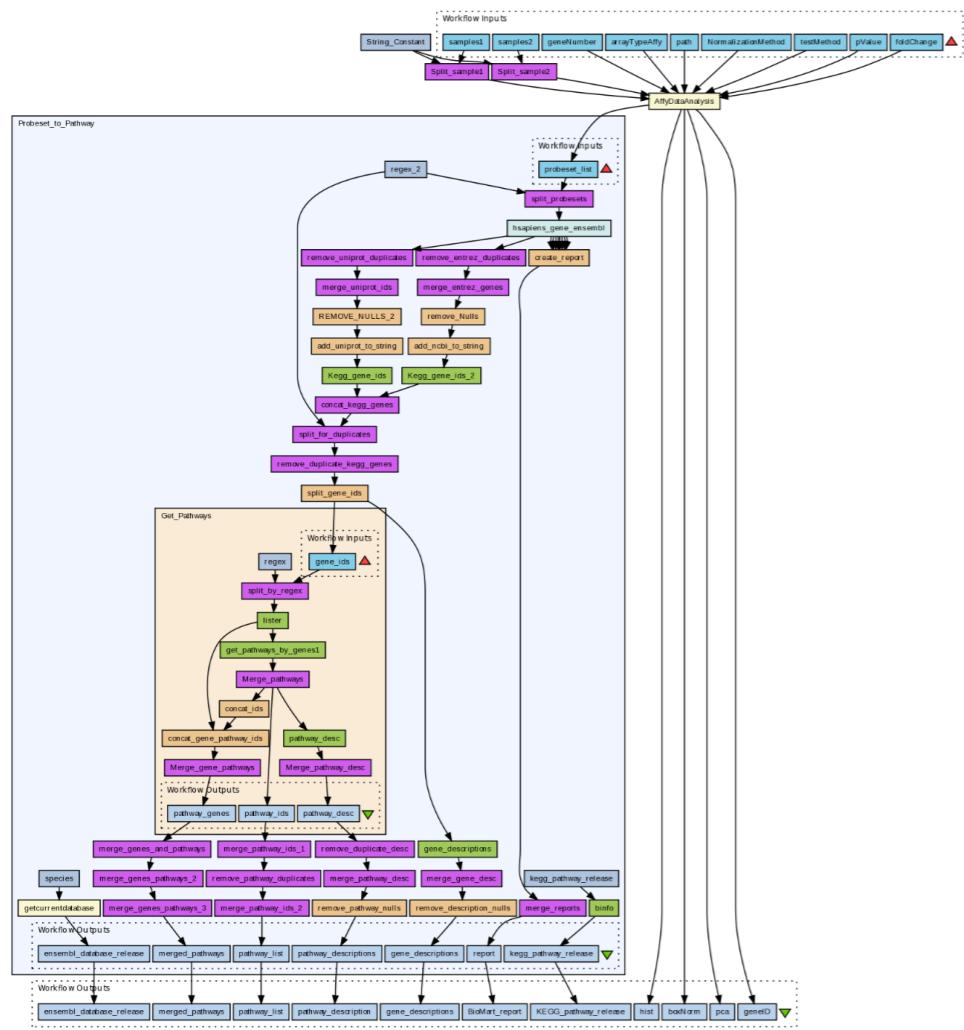
Intracranial aneurysms (ICAs) are acquired cerebrovascular abnormalities characterized by localized dilation and wall thinning in intracranial arteries, possibly leading to subarachnoid hemorrhage and severe outcome in case of rupture. Here, we identified one rare nonsense variant (c.1378A>T) in the last exon of ANGPTL6 (Angiopoietin-Like 6)—which encodes a circulating pro-angiogenic factor mainly secreted from the liver—shared by the four tested affected members of a large pedigree with multiple IA-affected case subjects. We showed a 50% reduction of ANGPTL6 serum concentration in individuals heterozygous for the c.1378A>T allele (p.Lys460Ter) compared to relatives homozygous for the normal allele, probably due to the non-secretion of the truncated protein produced by the c.1378A>T transcript. *Concomitant ANGPTL6 in a series of 94 additional index case subjects with familial IA identified three other*



| Filtering steps          | Remaining variants |             |
|--------------------------|--------------------|-------------|
|                          | Individuals :      | III-1 III-5 |
| Functional variants      | 25,674             | 23,456      |
| MAF < 0,1% in ExAC (NFE) | 456                | 436         |
| Shared by III.1 & III.5  |                    | 29          |
| In IBD haplotypes        |                    | 10          |
| Shared by all affected   |                    | 8           |
| Predicted LOF            | 1 (ANGPTL6)        |             |

## Genome-wide association study of intracranial aneurysms identifies 17 risk loci and genetic overlap with clinical risk factors

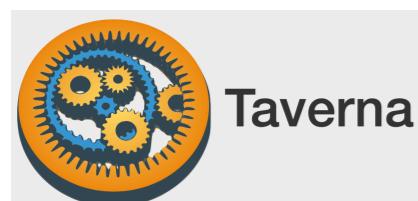




# Workflows to enhance trust in scientific results :

- automation (scalability)
- abstraction (methods sharing)
- traceability (provenance)

**nextflow**



**TensorFlow**



**snakemake**

**COMMON  
WORKFLOW  
LANGUAGE**

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

[Branch: master | New pull request](#) [Create new file](#) [Upload files](#) [Find file](#) [Clone or download](#)

| Branch: master                  | New pull request   | Create new file | Upload files | Find file | Clone or download |
|---------------------------------|--------------------|-----------------|--------------|-----------|-------------------|
| <a href="#">CONTRIBUTING.md</a> | Added contributing | 4 years ago     |              |           |                   |
| <a href="#">README.md</a>       | Update README.md   | 25 days ago     |              |           |                   |
| <a href="#">README.md</a>       | Update README.md   | 25 days ago     |              |           |                   |

**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.
- Ageage - 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.
- Antha - High-level language for biology.
- Bds - 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.
- Compss - 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.
- Dagobian - 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.
- Doit - 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.

# Galaxy platform

Galaxy

Workflow Preview

Download Import Run

Preprocessing

1: Reverse raw reads

2: Forward raw reads

3: FastQC

4: FastQC

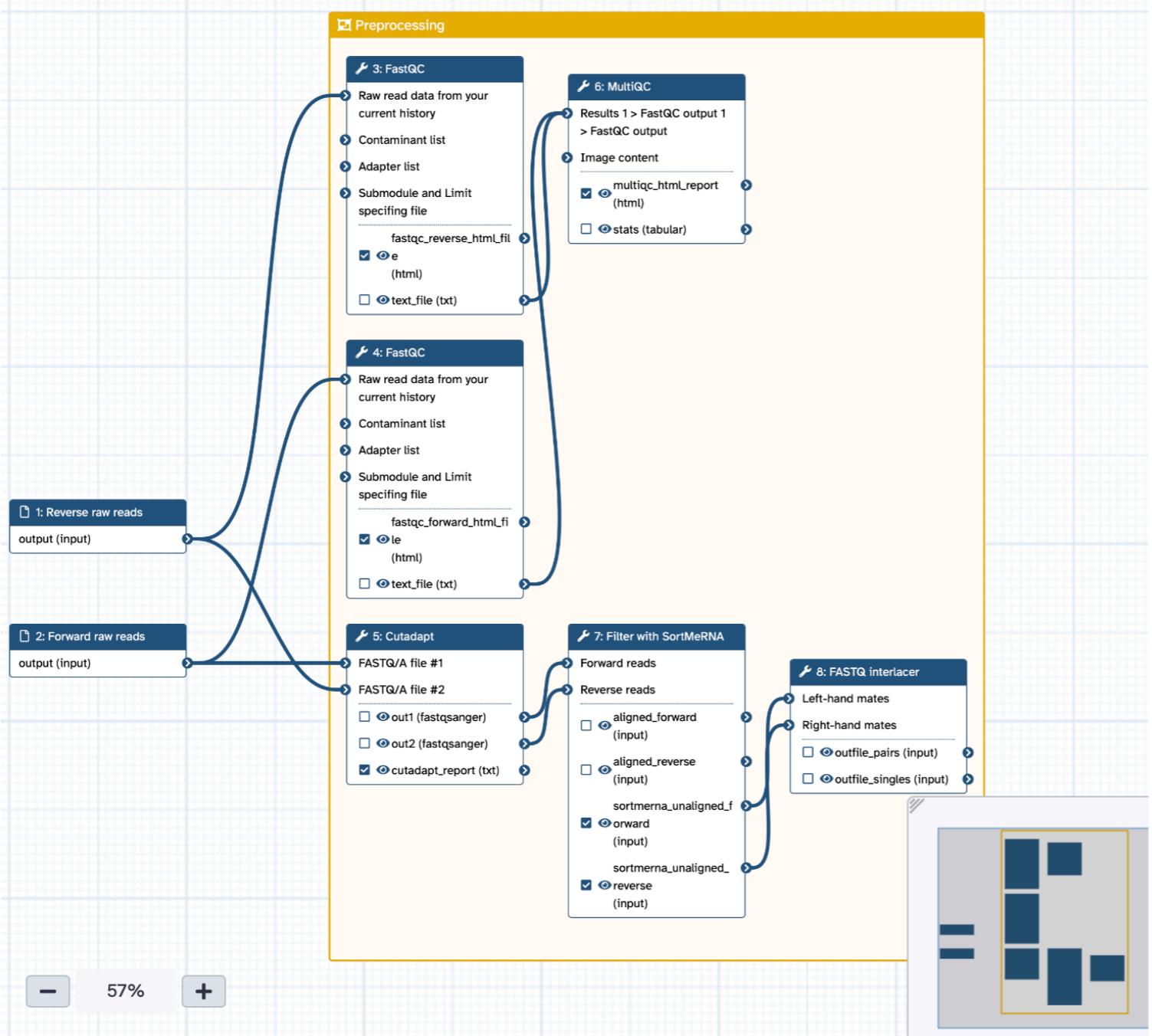
5: Cutadapt

6: MultiQC

7: Filter with SortMeRNA

8: FASTQ interlacer

57%



About This Workflow

Metatranscriptomics analysis using microbiome RNA-seq data - Workflow 1: Preprocessing - Version 2

Author

bebatus



All published Workflows by bebatus

Creators

- Bérénice Batut
- Pratik Jagtap
- Subina Mehta
- Ray Sajulga
- Emma Leith
- Praveen Kumar
- Saskia Hiltemann
- Paul Zierep
- Christine Oger

Description

Metatranscriptomics analysis using microbiome RNA-seq data (short)

Tags

microbiome

License

MIT License

Last Updated

Wednesday Jun 18th 8:56:01 2025 GMT+2

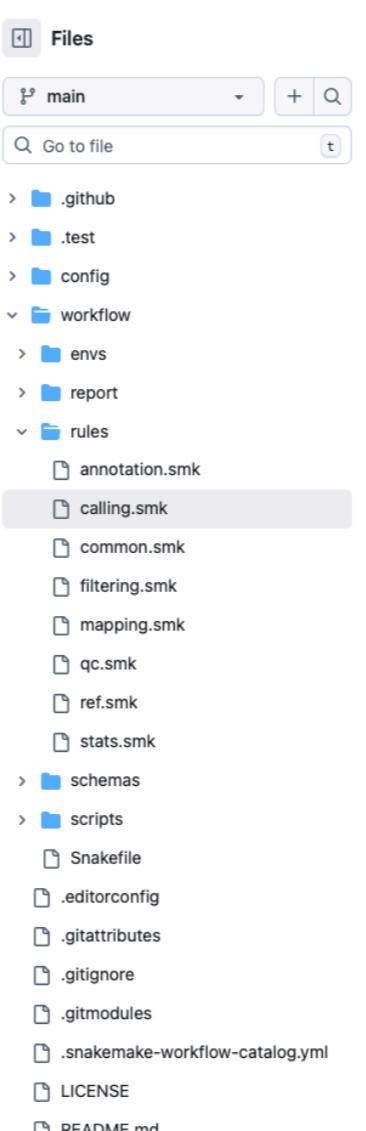
Sharing

<http://usegalaxy.fr>

<http://usegalaxy.eu>

# Other workflow languages / engines

- ▶ When more control is needed on
  - data organization
  - storage
  - HPC clusters / GPUs ..
- ▶ Rule-based data transformation
- ▶ Programming skills are needed



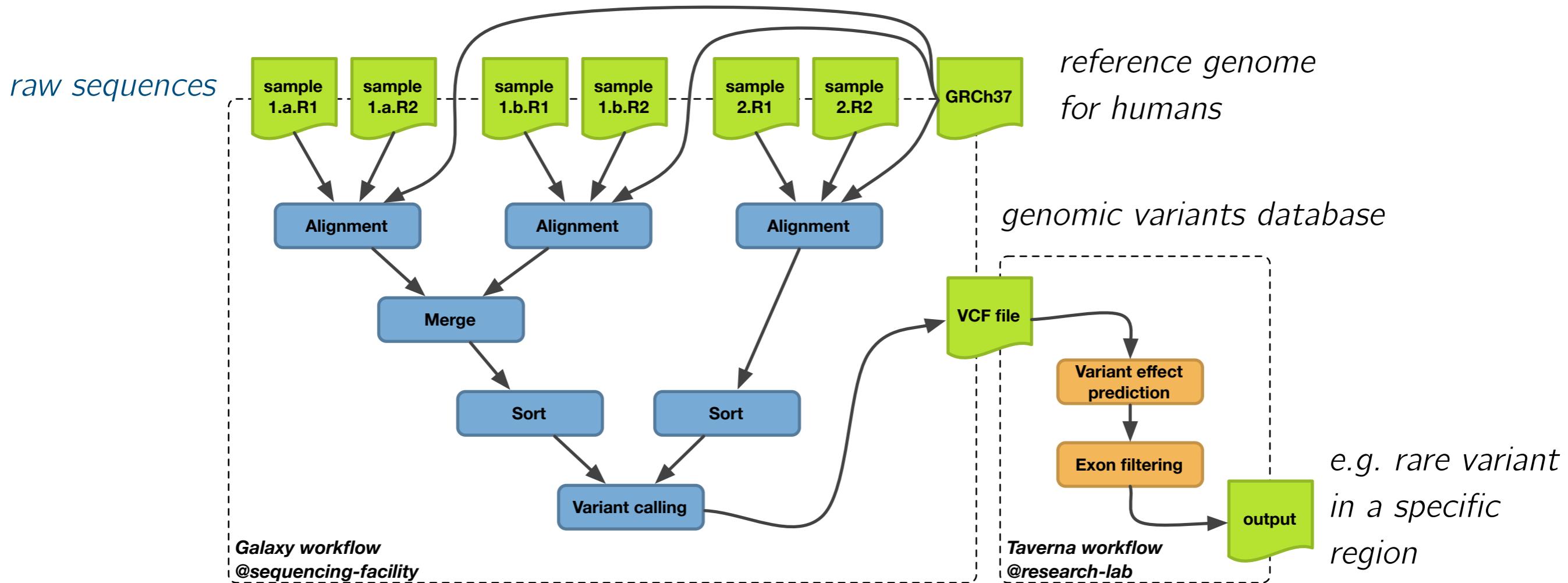
The screenshot shows a file browser interface with a sidebar and a main content area. The sidebar contains a 'Files' section with a search bar and a 'main' folder selected. Below the search bar is a 'Go to file' input field. The main content area shows a directory structure under 'main': '.github', '.test', 'config', 'workflow' (which contains 'envs', 'report', and 'rules'), and several smk files: 'annotation.smk', 'calling.smk' (which is highlighted), 'common.smk', 'filtering.smk', 'mapping.smk', 'qc.smk', 'ref.smk', 'stats.smk', 'schemas', 'scripts' (containing 'Snakefile', '.editorconfig', '.gitattributes', '.gitignore', '.gitmodules', and '.snakemake-workflow-catalog.yml'), and 'LICENSE' and 'README.md'. To the right of the file browser is a code editor window titled 'dna-seq-gatk-variant-calling / workflow / rules / calling.smk'. The code is a Snakemake rule definition for calling variants using GATK. It includes imports for resources like genome.fasta, genome.dict, variation.noipac.vcf.gz, and variation.noipac.vcf.tbi. The rule defines inputs (bam, ref, gvcf, log, params, wrapper), outputs (gvcf), and log files (logs/gatk/haplotypecaller/{sample}.{contig}.log). It also includes rules for combining calls, genotype variants, and merging variants using GATK tools.

```
12
13
14 rule call_variants:
15     input:
16         bam=get_sample_bams,
17         ref="resources/genome.fasta",
18         idx="resources/genome.dict",
19         known="resources/variation.noipac.vcf.gz",
20         tbi="resources/variation.noipac.vcf.gz.tbi",
21         regions=(
22             "results/called/{contig}.regions.bed"
23             if config["processing"].get("restrict-regions")
24             else []
25         ),
26     output:
27         gvcf=protected("results/called/{sample}.{contig}.g.vcf.gz"),
28     log:
29         "logs/gatk/haplotypecaller/{sample}.{contig}.log",
30     params:
31         extra=get_call_variants_params,
32     wrapper:
33         "0.59.0/bio/gatk/haplotypecaller"
34
35
36 rule combine_calls:
37     input:
38         ref="resources/genome.fasta",
39         gvcfs=expand(
40             "results/called/{sample}.{{contig}}.g.vcf.gz", sample=samples.index
41         ),
42     output:
43         gvcf="results/called/all.{contig}.g.vcf.gz",
44     log:
45         "logs/gatk/combinegvcfs.{contig}.log",
46     wrapper:
47         "0.74.0/bio/gatk/combinegvcfs"
48
49
50 rule genotype_variants:
51     input:
52         ref="resources/genome.fasta",
53         gvcf="results/called/all.{contig}.g.vcf.gz",
54     output:
55         vcf=temp("results/genotyped/all.{contig}.vcf.gz"),
56     params:
57         extra=config["params"]["gatk"]["GenotypeGVCFs"],
58     log:
59         "logs/gatk/genotypegvcfs.{contig}.log",
60     wrapper:
61         "0.74.0/bio/gatk/genotypegvcfs"
62
63
64 rule merge_variants:
65     input:
66         vcfs=lambda w: expand(
67             "results/genotyped/all.{contig}.vcf.gz", contig=get_contigs()
68         ),
69     output:
70         vcf="results/genotyped/all.vcf.gz",
```

# Massively produced data

Can we **reuse** data  
rather than **re-compute** ?

# Reusing processed data ?

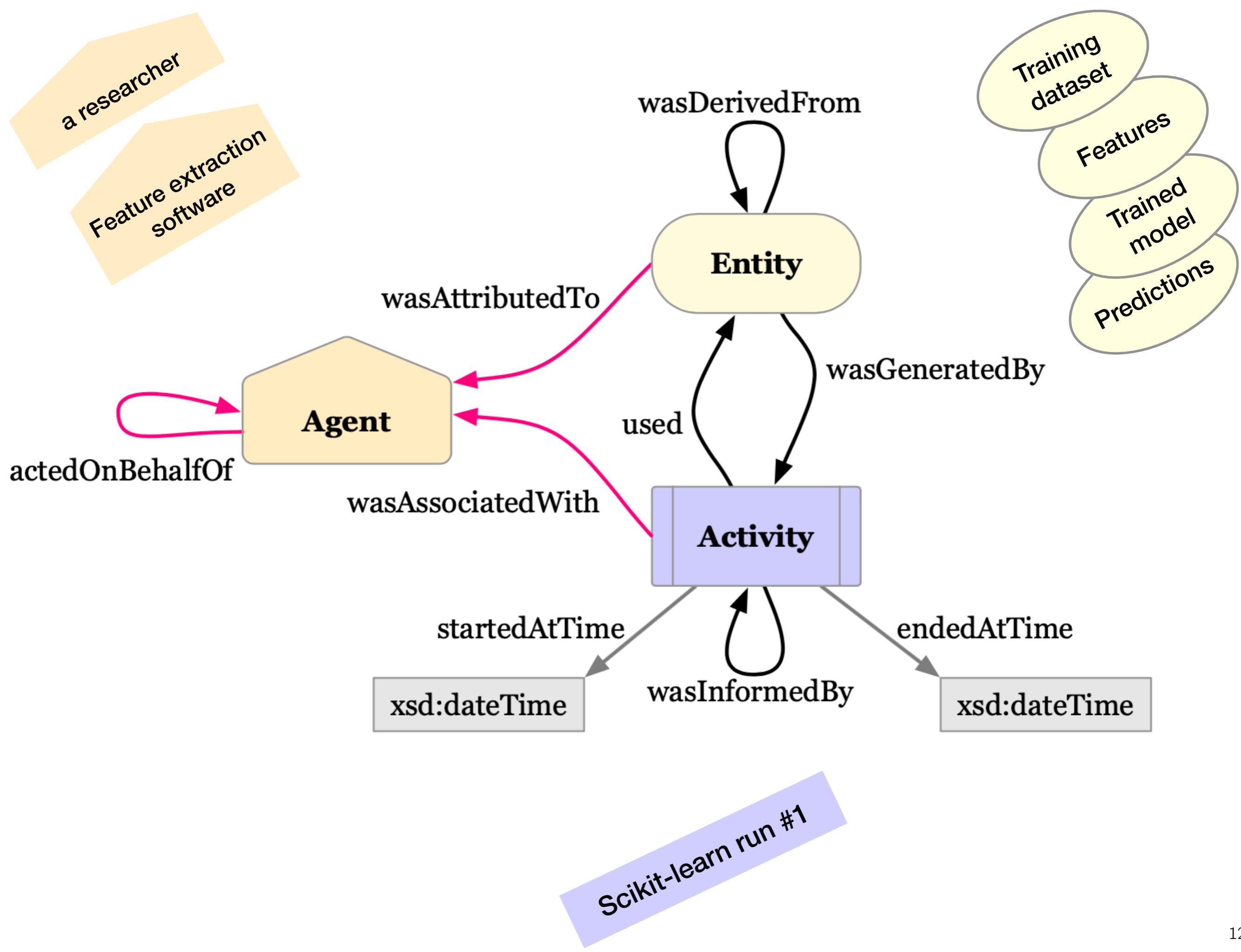


« Which alignment algorithm was used when predicting this pathogenic score ? »

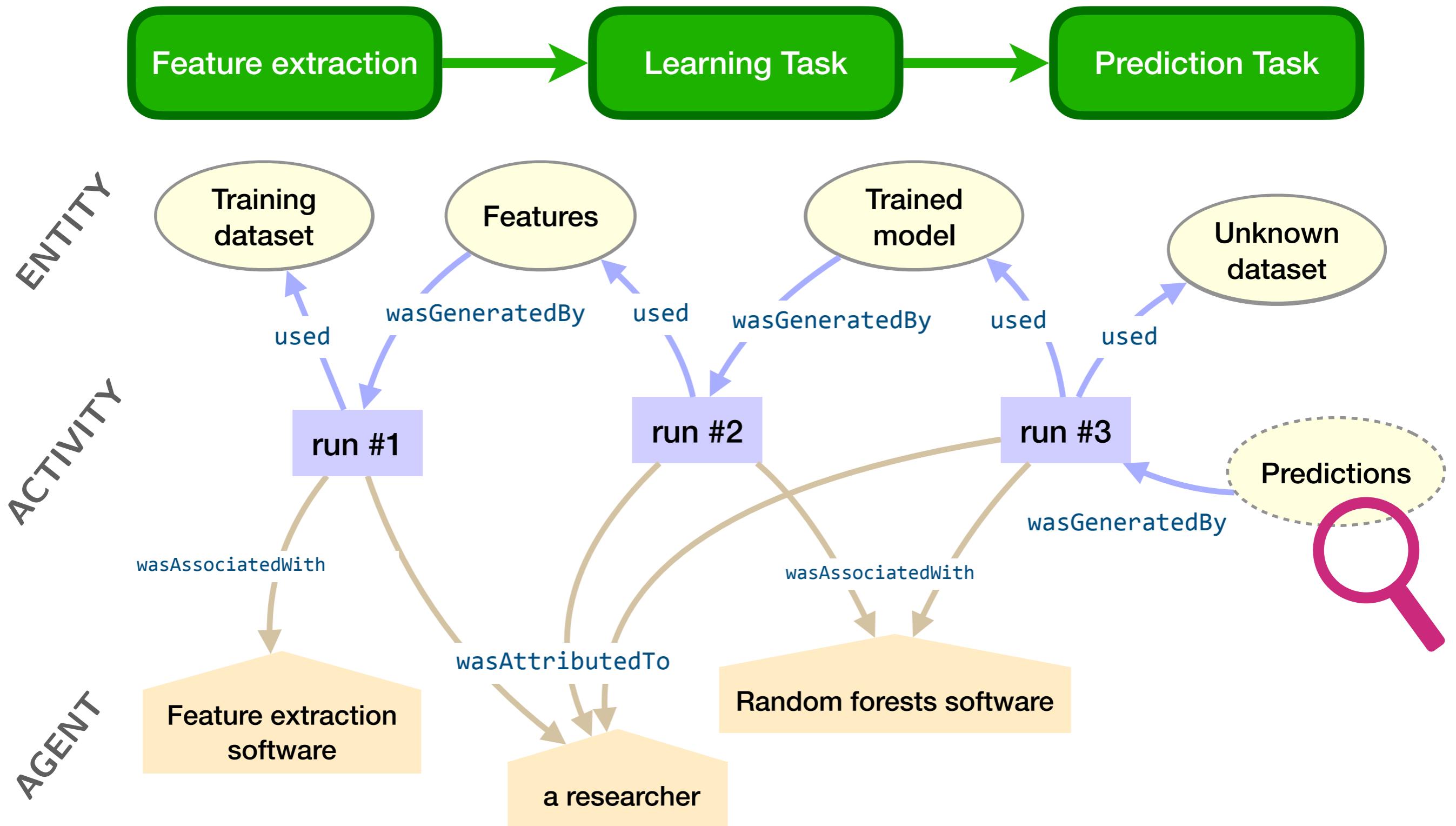
« A new version of a reference genome is available, which genome was used when predicting these phenotypes ? »

**Need for an overall tracking of provenance over multiple workflows !**

Hopefully, we have  
W3C PROV



# Linked data-algorithms-software-researchers



# PROV is domain-agnostic 😊🤔

[nature](#) > [scientific data](#) > [articles](#) > [article](#)

[Open Access](#) | Published: 06 December 2016

## Sharing brain mapping statistical results with the neuroimaging data model

[Camille Maumet](#)✉, [Tibor Auer](#), [Alexander Bowring](#), [Gang Chen](#), [Samir Das](#), [Guillaume Flandin](#), [Satrajit Ghosh](#), [Tristan Glatard](#), [Krzysztof J. Gorgolewski](#), [Karl G. Helmer](#), [Mark Jenkinson](#), [David B. Keator](#), [B. Nolan Nichols](#), [Jean-Baptiste Poline](#), [Richard Reynolds](#), [Vanessa Sochat](#), [Jessica Turner](#) & [Thomas E. Nichols](#)

[Scientific Data](#) 3, Article number: 160102 (2016) | [Cite this article](#)

5277 Accesses | 30 Citations | 42 Altmetric | [Metrics](#)

<https://doi.org/10.1038/sdata.2016.102>



Research | [Open Access](#) | Published: 31 January 2022

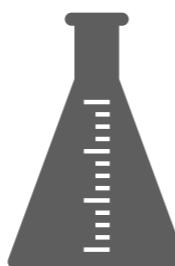
## Structure-based knowledge acquisition from electronic lab notebooks for research data provenance documentation

[Max Schröder](#)✉, [Susanne Staehlke](#), [Paul Groth](#), [J. Barbara Nebe](#), [Sascha Spors](#) & [Frank Krüger](#)

[Journal of Biomedical Semantics](#) 13, Article number: 4 (2022) | [Cite this article](#)

4265 Accesses | 6 Citations | 9 Altmetric | [Metrics](#)

<https://doi.org/10.1186/s13326-021-00257-x>



TaPP 2021  
PAPERS  
A USENIX Publication  
ACCEPTED PAPERS

## Astronomical Pipeline Provenance: A Use Case Evaluation

Authors:

Michael A. C. Johnson, *Institute of Data Science (DLR)* and *Max Planck Institute for Radio Astronomy*; Marcus Paradies and Marta Dembska, *Institute of Data Science (DLR)*; Kristen Lackeos, Hans-Rainer Klöckner, and David J. Champion, *Max Planck Institute for Radio Astronomy*; Sirk Schindler, *Institute of Data Science (DLR)*

<https://doi.org/10.48550/arXiv.2109.10759>

[Home](#) > [Provenance and Annotation of Data and Processes](#) > Conference paper

## Towards a Provenance Management System for Astronomical Observatories

[Mathieu Servillat](#)✉, [François Bonnarel](#), [Catherine Boisson](#), [Mireille Louys](#), [Jose Enrique Ruiz](#) & [Michèle Sanguillon](#)

Conference paper | [First Online: 09 July 2021](#)

515 Accesses | 1 Citations | 8 Altmetric

Part of the [Lecture Notes in Computer Science](#) book series (LNISA, volume 12839)

[https://doi.org/10.1007/978-3-030-80960-7\\_20](https://doi.org/10.1007/978-3-030-80960-7_20)

# Many expectations ...

Comparability, transparency, confidence +

- ▶ **Citing** researchers and organisations
- ▶ Identifying **critical** data / software **resources** associated to scientific results
- ▶ Identifying possible **bias** when reusing / sharing datasets / pre-trained models

We need **guidance** for  
documenting  
Life-Science provenance

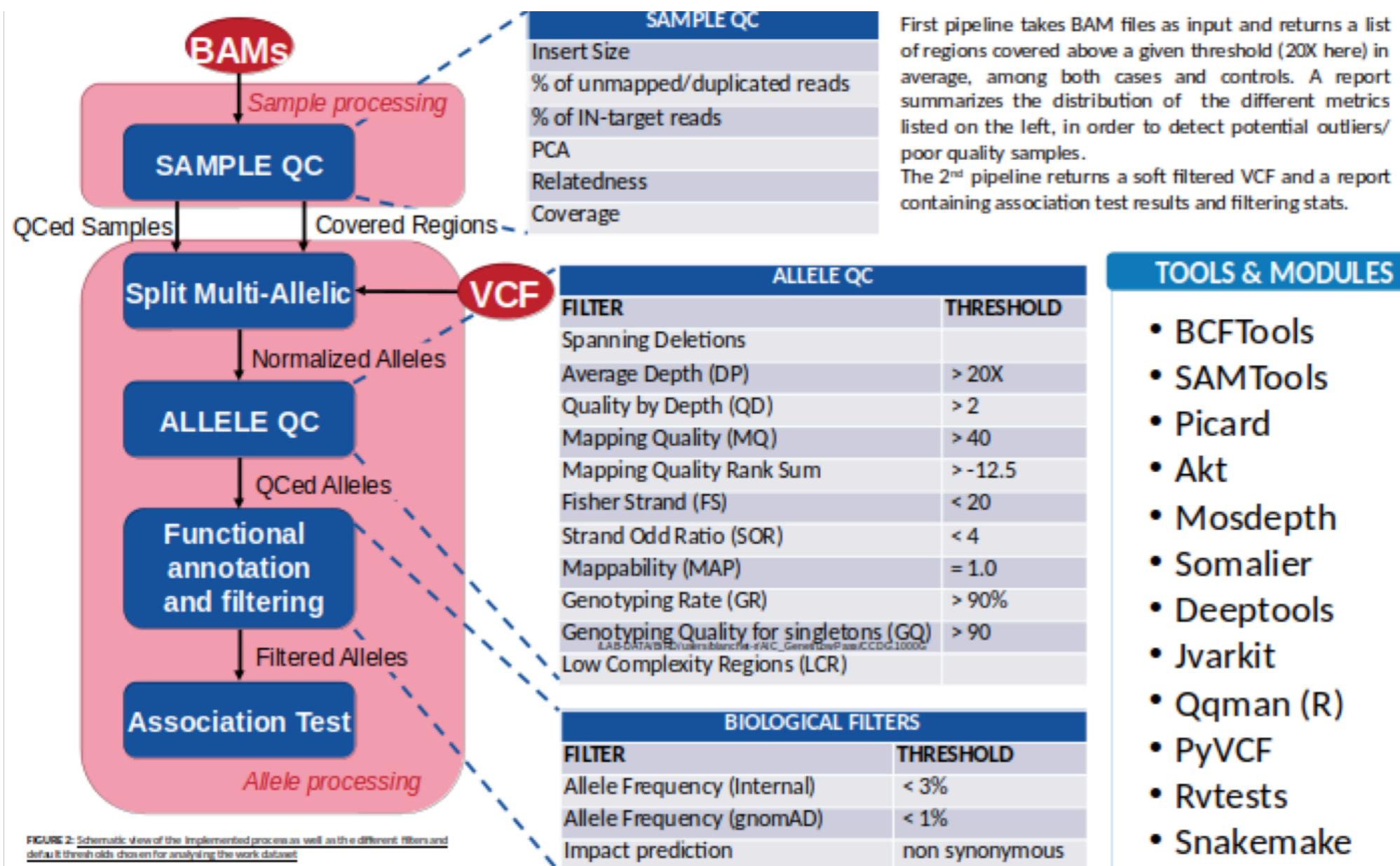
# MIRAPIE

for genomic workflow ?



# WF#1 Rare variants analysis

Raphaël Blanchet



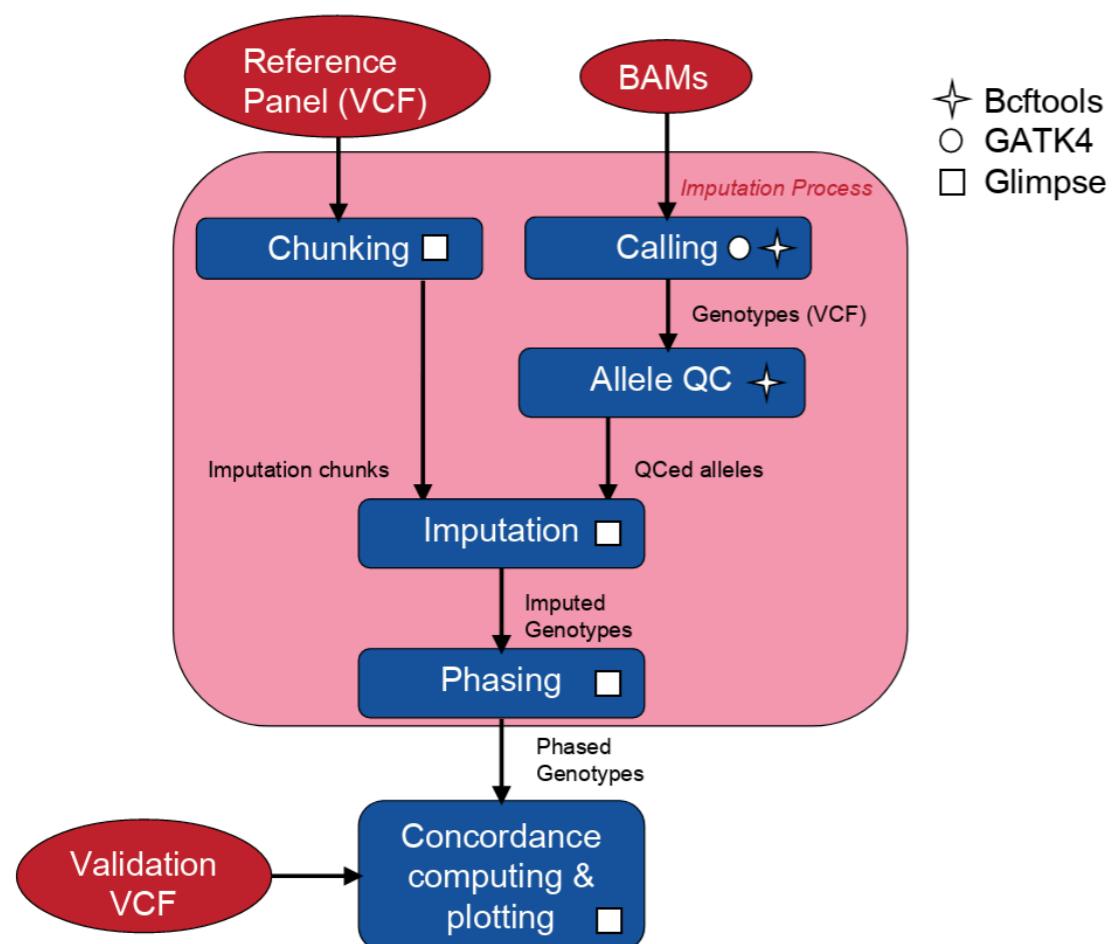
Provenance issue : tracking all these  
(hidden) domain-specific parameters ?



# WF#2 Common variants analysis (LPS)



Raphaël Blanchet



- ▶ Low Pass Sequencing
- ▶ WF designed on 2600 whole genomes with very low coverage (1x, 2x, 4x)
- ▶ 1000 whole genomes for reference panel
- ▶ Statistical imputation

Provenance issues: which reference panel ? which imputation method ?

# W<sub>7</sub> = Who - When - Where - Why - What - hoW - Which ?

## Who ?

- ▶ “Raphael Blanchet”, as a PhD Student, working at ITX, funded by Inserm
- ▶ ITX, as a research lab
- ▶ Inserm, as a research organization, funding institute

## When ?

- ▶ Acquisition of input data
- ▶ Sequencing of biological samples
- ▶ Design of the data analysis workflow
- ▶ Run of the data analysis workflow

## Where ?

- ▶ Location of the researcher
- ▶ Location of the biobank
- ▶ Sequenced biological samples
- ▶ Designed of the data analysis workflow
- ▶ Run of data analysis workflow

# W<sub>7</sub> = Who - When - Where - Why - What - hoW - Which ?

## What ?

*WHAT* is a sequence of events <e1, e2, ..., en> that affect a data object during its life time.

- ▶ workflow execution, fine-grained

## Why ?

- ▶ Are rare variant associated with the formation of intracranial aneurysms ?
- ▶ Are rare variant associated with the rupture of intracranial aneurysms ?

## Which (device) ?

- ▶ bcftools (<http://bio.tools/bcftools>)
- ▶ picard (<http://bio.tools/picard>)
- ▶ WF#1 rare variant analysis workflow
- ▶ **GLICID HPC infrastructure**

## hoW ?

"How" documents actions that lead to the occurrence of an event."

- ▶ WF#1 rare variant analysis workflow
- ▶ **reference guidelines, protocols, papers ?**

# (partial) Implementation with PROV-O

PROV-O metadata captured at runtime with a slightly modified implementation of the Snakemake workflow engine.

```
--  
108    <http://snakemake-provenance#activity-e5a4336a-3495-4821-9bd2-24d444b61637>  
109        a prov:Activity ;  
110        rdfs:comment """  
111            bwa mem -t 1 -M           -H '@CO\tProject:test Sample:Sample1 Date:2019-01-23 CWD:/Users/gaignard-a/D  
112            """ ;  
113            prov:wasAssociatedWith <https://bio.tools/bwa_cloudfb> ;  
114            prov:startedAtTime "2019-01-23T15:27:59.661594"^^xsd:dateTime;  
115            prov:endedAtTime "2019-01-23T15:27:59.661606"^^xsd:dateTime;  
116            prov:used <../testdata/human_g1k_v37.chr22.fasta> ;  
117            prov:used <../testdata/Sample1_ATGCCTAA_L001_R1_001.fastq.gz> ;  
118            prov:used <../testdata/Sample1_ATGCCTAA_L001_R2_001.fastq.gz> ;  
119            prov:used <../testdata/human_g1k_v37.chr22.fasta.bwt> ;  
120            .  
121  
122    <Samples/Sample1/BAM/Sample1.p1.aligned.sam>  
123        a prov:Entity;  
124        prov:wasGeneratedBy <http://snakemake-provenance#activity-e5a4336a-3495-4821-9bd2-24d444b61637>;  
125        prov:wasAssociatedWith <https://bio.tools/bwa_cloudfb> ;  
126        rdfs:label "Samples/Sample1/BAM/Sample1.p1.aligned.sam";  
127        prov:wasDerivedFrom <../testdata/human_g1k_v37.chr22.fasta> ;  
128        prov:wasDerivedFrom <../testdata/Sample1_ATGCCTAA_L001_R1_001.fastq.gz> ;  
129        prov:wasDerivedFrom <../testdata/Sample1_ATGCCTAA_L001_R2_001.fastq.gz> ;  
130        prov:wasDerivedFrom <../testdata/human_g1k_v37.chr22.fasta.bwt> ;  
131        .
```

W3C Working Group Note



PROV-Overview

An Overview of the PROV Family of Documents

W3C Working Group Note 30 April 2013

<https://github.com/albangaillard/fresh-toolbox/blob/master/provenance.ttl>



[https://github.com/albangaillard/galaxy-PROV/blob/master/notebooks/galaxy\\_prov.ttl](https://github.com/albangaillard/galaxy-PROV/blob/master/notebooks/galaxy_prov.ttl)



# What is missing in PROV for biomedical use cases ?

... to increase findability, trust and reuse

- ▶ Link with the software containers / computing platform / **hardware**
- ▶ Fine-grained **biological parameters** (thresholds) ?
- ▶ Data & software **access conditions**
  - out of the scope of PROV (DUO ontology ?)
- ▶ **Reference methods / papers** associated with data analysis
  - PAV ? SIO ? MicroPublication ? CITO ?
- ▶ **Research context** ? questions ? outcomes ? data management plans ?

PEPR Santé numérique

- ▶ National project funded on Digital Health : ShareFAIR
- ▶ Focus on workflows for reproducibility / provenance / FAIR protocols & data
- ▶ 6 PhD projects in progress + 1 software developer

# Any questions ?