

A Comparative Guide to Provenance Tracking in Scientific Workflow Management Systems

Author: BenchChem Technical Support Team. **Date:** December 2025

Compound of Interest

Compound Name: *Pegasus*

Cat. No.: *B039198*

[Get Quote](#)

In the realms of scientific research and drug development, the ability to meticulously track the origin and transformation of data—a practice known as provenance tracking—is not merely a feature but a cornerstone of reproducibility, validation, and regulatory compliance. Workflow Management Systems (WMS) are pivotal in automating and managing complex computational experiments, and their proficiency in capturing provenance is a critical factor for adoption. This guide provides an objective comparison of how **Pegasus** WMS and several popular alternatives—Nextflow, Snakemake, CWL (Common Workflow Language), and Galaxy—handle provenance tracking for experiments.

Comparison of Provenance Tracking Features

The following table summarizes the key provenance tracking capabilities of the discussed Workflow Management Systems.

Feature	Pegasus WMS	Nextflow	Snakemake	CWL (Common Workflow Language)	Galaxy
Provenance Capture	Automatic, via the "kickstart" process for every job. Captures runtime information, including executable, arguments, environment variables, and resource usage.[1][2]	Automatic. Captures task execution details, input/output files, parameters, and container information.[3]	Automatic. Tracks input/output files, parameters, software environments (Conda), and code changes.[4][5]	Not inherent to the language, but supported by runners like cwltool which can generate detailed provenance information.[6]	Automatic. Every analysis step and user action is recorded in a user's history, creating a comprehensive audit trail.[7]
Data Model	Stores provenance in a relational database (SQLite by default) with a well-defined schema.[8][9]	Has a native, experimental data lineage feature with a defined data model.[3][10] Also supports export to standard formats like RO-Crate and BioCompute Objects via the nf-prov plugin.[10]	Stores provenance information in a hidden .snakemake directory, tracking metadata for each output file.[4]	Promotes the use of the W3C PROV model through its CWLProv profile for a standardized representation of provenance.[6][11][12]	Maintains an internal data model that captures the relationships between datasets, tools, and parameters within a user's history. Can be exported to standard formats.[7]

Query & Exploration	Provides command-line tools (pegasus-statistics, pegasus-plots) and allows direct SQL queries on the provenance database for detailed analysis.[1][2]	The nextflow log command provides summaries of workflow executions. The experimental nextflow lineage command allows for more detailed querying of the provenance data.[3]	The --summary command-line option provides a concise overview of the provenance for each output file.[5] Generates interactive HTML reports for visual exploration of the workflow and results.	The CWLProv profile facilitates the use of standard RDF and PROV query languages (e.g., SPARQL) for complex provenance queries.	The web-based interface allows for easy exploration of the analysis history. Histories and workflows can be exported, and workflow invocation reports can be generated.[13]
Standardization & Export	Captures detailed provenance but does not natively export to a standardized format like PROV-O.	The nf-prov plugin enables the export of provenance information to standardized formats like BioCompute Objects and RO-Crate.[10]	Does not have a built-in feature for exporting to standardized provenance formats, though there are community efforts to enable PROV-JSON export.[14]	CWLProv is a profile for recording provenance as a Research Object, using standards like BagIt, RO, and W3C PROV.[6][11][12]	Can export histories and workflows in its own format, and increasingly supports standardized formats like RO-Crate for workflow invocations.[7]
User Interface	Provides a web-based dashboard for monitoring	Primarily command-line based. Visualization	Generates self-contained, interactive	As a specification, it does not have a user	A comprehensive web-based

workflows,	of the	HTML reports	interface.	graphical
which can	workflow	that visualize	Visualization	user interface
also be used	Directed	the workflow	depends on	is its core
to inspect	Acyclic Graph	and its	the	feature,
some	(DAG) can be	results. [5]	implementati	making
provenance	generated.		on of the	provenance
information.	[15]		CWL runner.	exploration
				highly
				accessible.

Experimental Protocols: A Representative Genomics Workflow

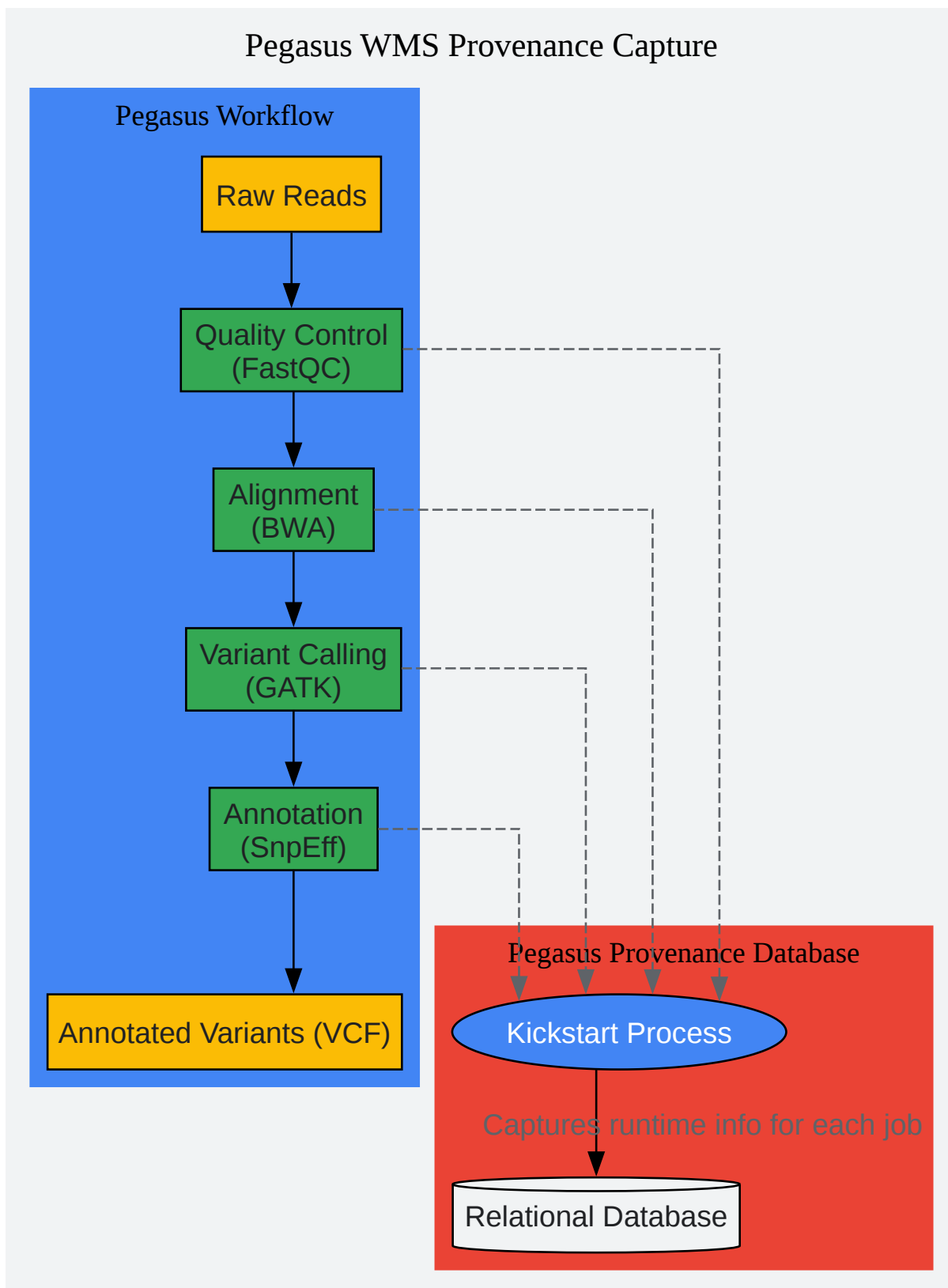
To illustrate and compare the provenance tracking mechanisms, we will consider a common genomics workflow for identifying genetic variants from sequencing data. This workflow typically involves the following key steps:

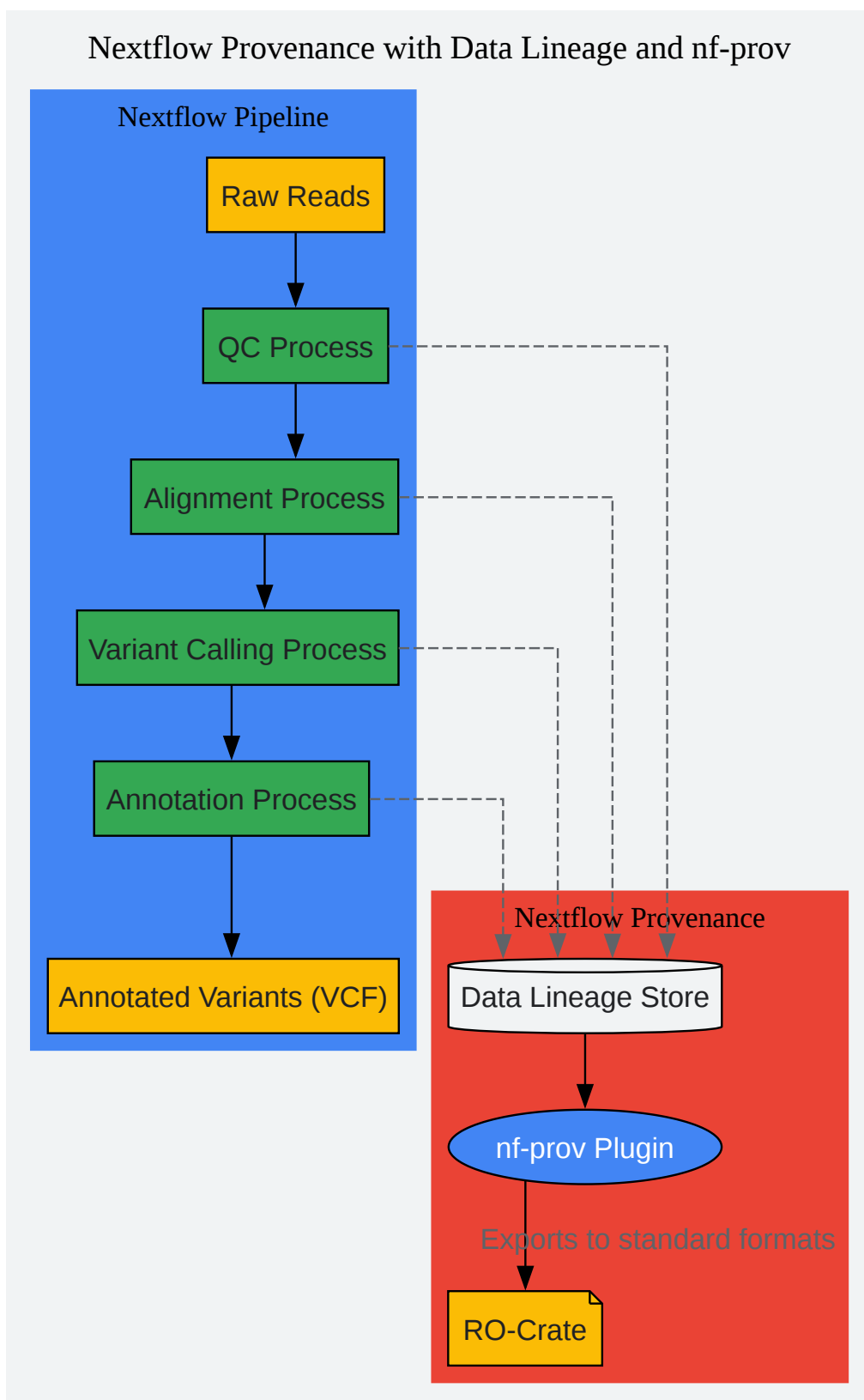
- Quality Control (QC): Assessing the quality of raw sequencing reads.
- Alignment: Aligning the sequencing reads to a reference genome.
- Variant Calling: Identifying differences between the aligned reads and the reference genome.
- Annotation: Annotating the identified variants with information about their potential functional impact.

Each of these steps involves specific software tools, parameters, and reference data files, all of which are critical pieces of provenance information that a WMS should capture.

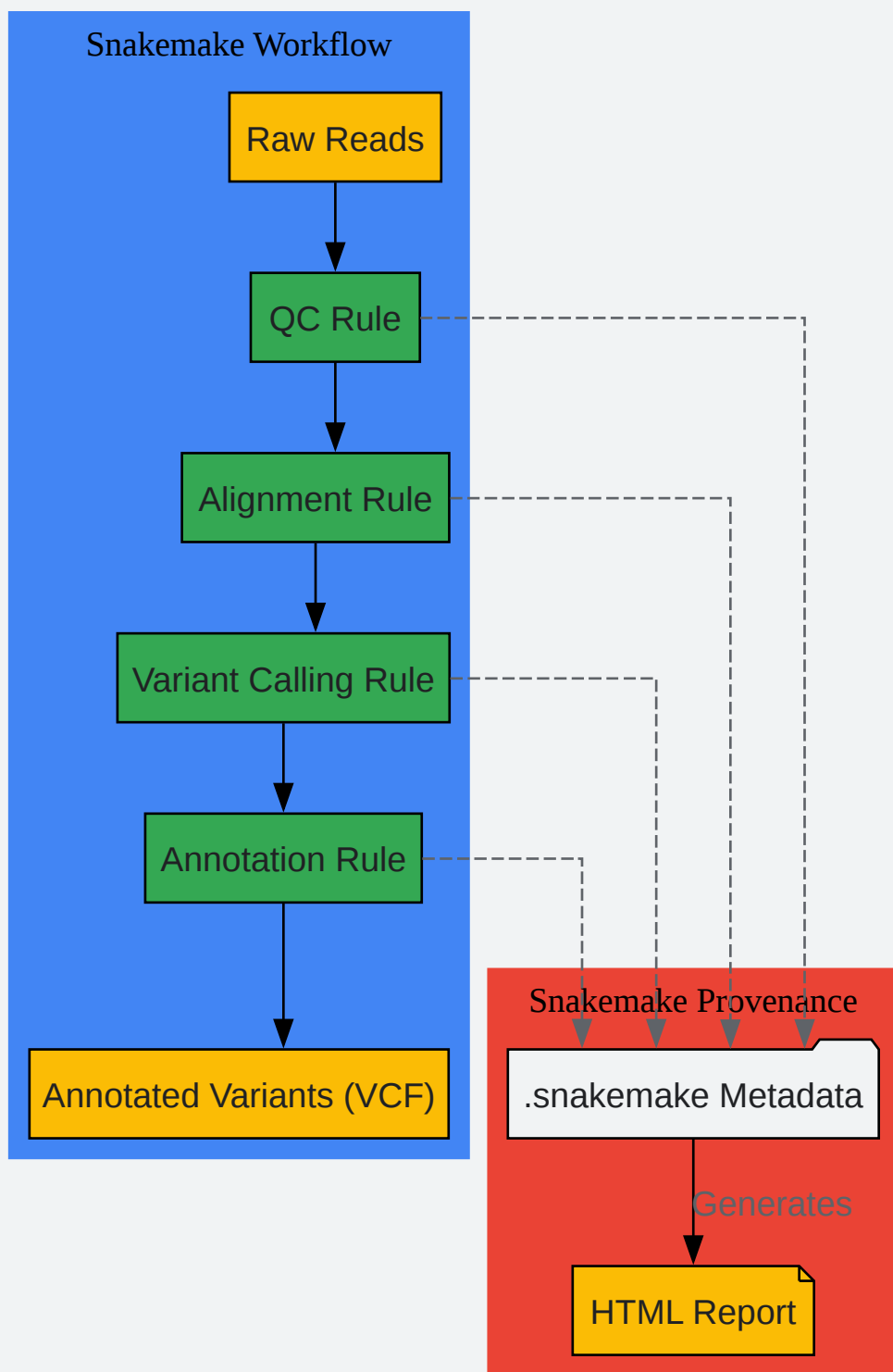
Visualizing Provenance Tracking in Action

The following diagrams, generated using the DOT language, illustrate a simplified conceptual model of how each WMS captures the provenance of this genomics workflow.

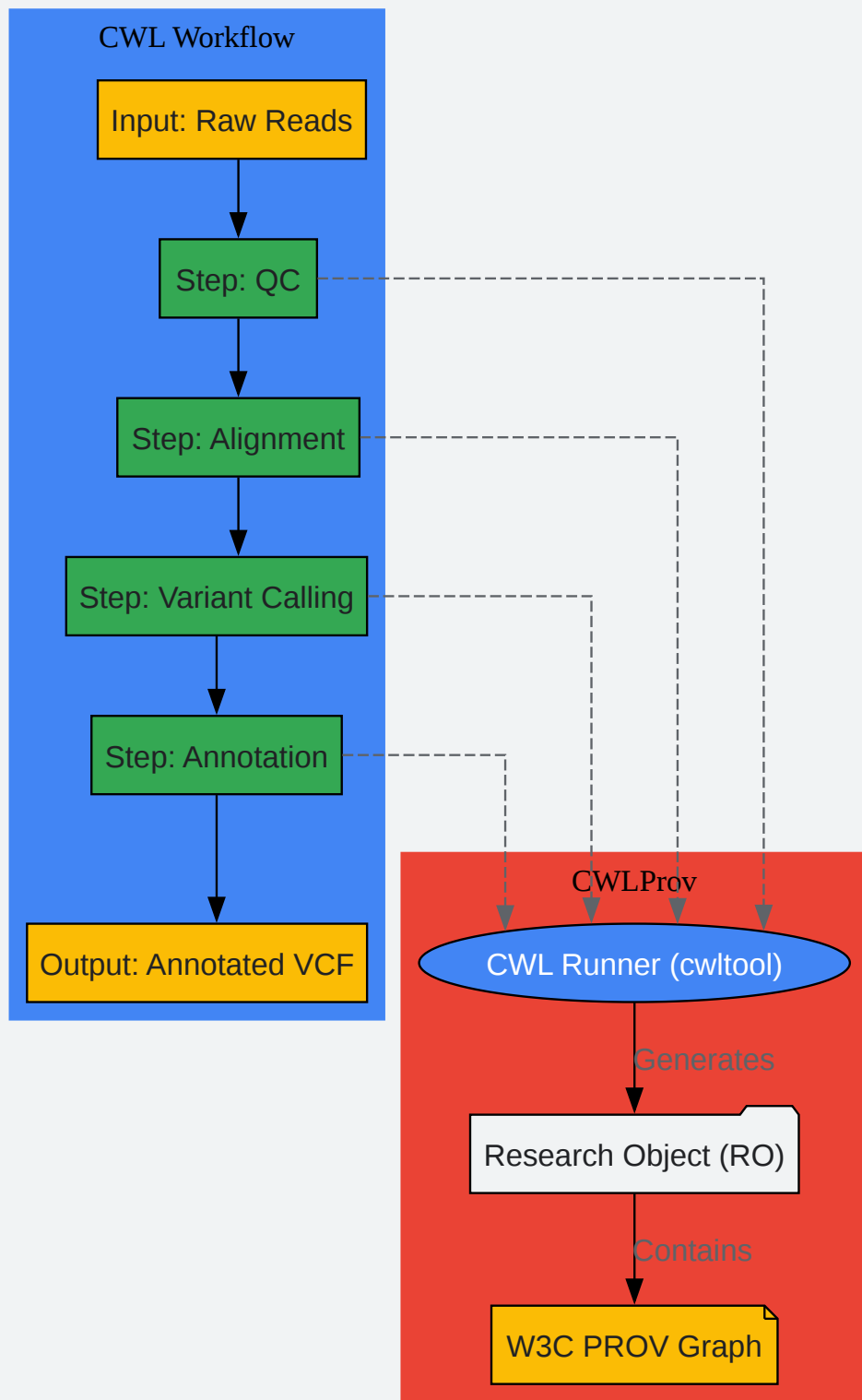




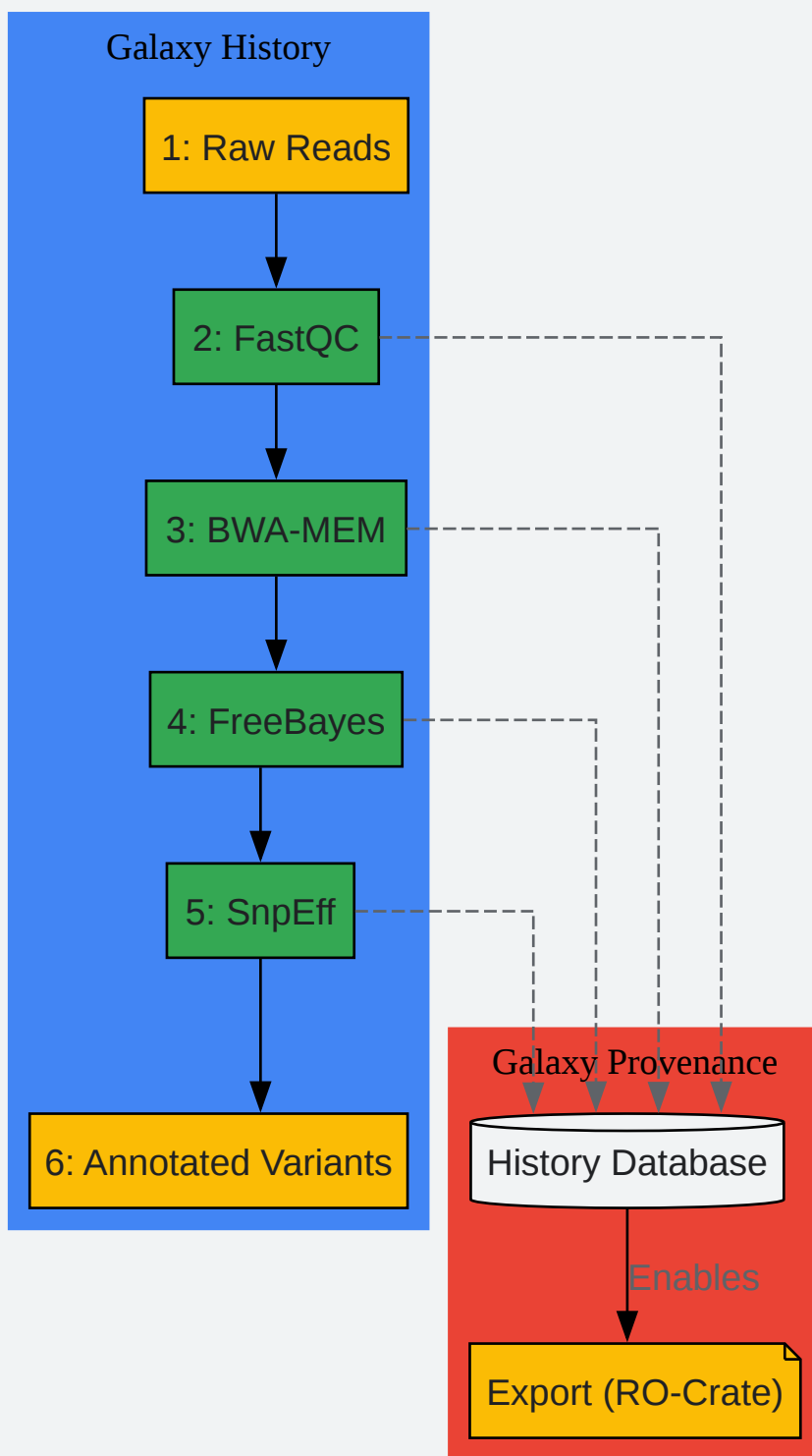
Snakemake Provenance Capture and Reporting



CWL Provenance via CWLProv and Research Objects



Galaxy's History-Based Provenance and Export

[Click to download full resolution via product page](#)

Need Custom Synthesis?

BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.

Email: info@benchchem.com or [Request Quote Online](#).

References

- 1. arokem.github.io [arokem.github.io]
- 2. 1. Introduction — Pegasus WMS 5.1.2-dev.0 documentation [pegasus.isi.edu]
- 3. Getting started with data lineage — Nextflow documentation [nextflow.io]
- 4. stackoverflow.com [stackoverflow.com]
- 5. Advanced: Decorating the example workflow | Snakemake 9.14.5 documentation [snakemake.readthedocs.io]
- 6. GitHub - common-workflow-language/cwlprov: Profile for provenance research object of a CWL workflow run [github.com]
- 7. direct.mit.edu [direct.mit.edu]
- 8. 14. Glossary — Pegasus WMS 5.1.2-dev.0 documentation [pegasus.isi.edu]
- 9. 13. Migration Notes — Pegasus WMS 5.1.2-dev.0 documentation [pegasus.isi.edu]
- 10. seqera.io [seqera.io]
- 11. academic.oup.com [academic.oup.com]
- 12. Sharing interoperable workflow provenance: A review of best practices and their practical application in CWLProv - PubMed [pubmed.ncbi.nlm.nih.gov]
- 13. Hands-on: Workflow Reports / Workflow Reports / Using Galaxy and Managing your Data [training.galaxyproject.org]
- 14. Allow generation of PROV-JSON for output files · Issue #2077 · snakemake/snakemake · GitHub [github.com]
- 15. Reports — Nextflow documentation [nextflow.io]
- To cite this document: BenchChem. [A Comparative Guide to Provenance Tracking in Scientific Workflow Management Systems]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b039198#how-does-pegasus-wms-handle-provenance-tracking-for-experiments]

Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

Technical Support: The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

Need Industrial/Bulk Grade? [Request Custom Synthesis Quote](#)

BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

Contact

Address: 3281 E Guasti Rd

Ontario, CA 91761, United States

Phone: (601) 213-4426

Email: info@benchchem.com