# Application Notes and Protocols for Pegasus Workflows

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

| Compound of Interest | | |
|---|---|---|
| Compound Name: | Pegasus | |
| Cat. No.: | B039198 | Get Quote |

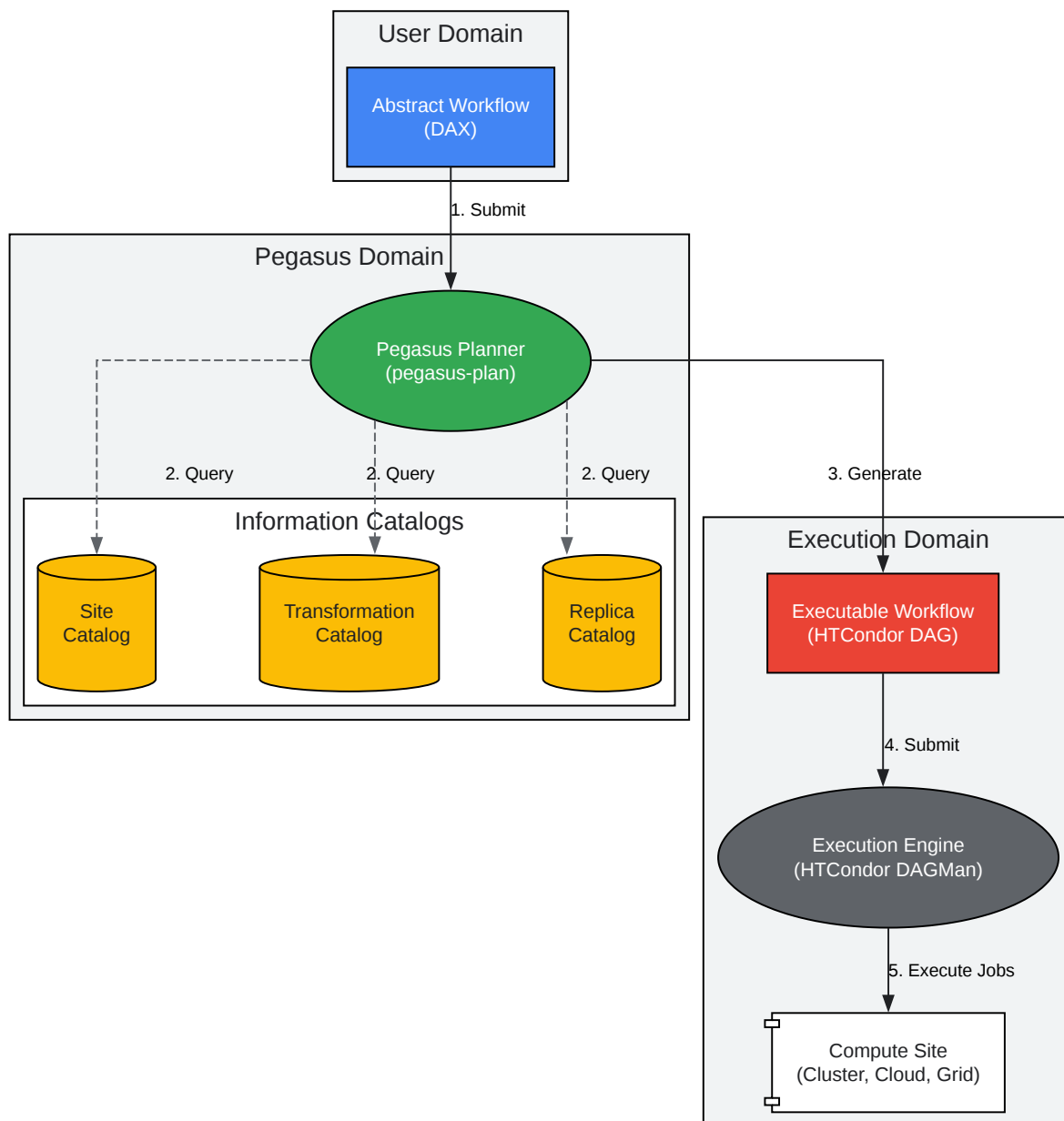Authored for: Researchers, Scientists, and Drug Development Professionals

Abstract: The automation of complex computational pipelines is critical in modern research, particularly in fields like bioinformatics and drug development, which involve large-scale data processing and analysis. The **Pegasus** Workflow Management System (WMS) provides a robust framework for defining, executing, and monitoring these complex scientific workflows across diverse computing environments.[1] By abstracting the logical workflow from the physical execution details, **Pegasus** enhances portability, reliability, and scalability.[2][3] This document provides a comprehensive guide to the core concepts of **Pegasus** and a step-by-step protocol for executing a sample bioinformatics workflow to identify mutational overlaps using data from the 1000 Genomes Project.

## Core Concepts of the **Pegasus** Workflow Management System

**Pegasus** is an open-source system that enables scientists to create abstract workflows that are automatically mapped and executed on a range of computational resources, including high-performance clusters, clouds, and grids.[1] The system is built on a key principle: the separation of the workflow description from the execution environment.[4] This allows the same workflow to be executed on a local machine, a campus cluster, or a national supercomputing facility without modification.[5]

The primary components of the **Pegasus** architecture are:

- Abstract Workflow (DAX): The scientist defines the workflow as a Directed Acyclic Graph (DAG), where nodes represent computational tasks and the edges represent dependencies. [4][6] This description, known as a DAX (Directed Acyclic Graph in XML), is abstract and does not specify where the code or data is located.[7] Users typically generate the DAX using a high-level API in Python, R, or Java.[1]

- **Pegasus** Planner (Mapper): This is the core engine that transforms the abstract DAX into a concrete, executable workflow.[8] It adds necessary auxiliary tasks for data management, such as staging input files, creating directories, and cleaning up intermediate data.[9]

- Information Catalogs: The planner consults three catalogs to resolve the physical details of the workflow:[9]

    - Site Catalog: Describes the computation and storage resources available (the "execution sites").[7][10]

    - Transformation Catalog: Maps the logical names of executables used in the workflow to their physical locations on the target sites.[10][11]

    - Replica Catalog: Maps the logical names of input files to their physical storage locations, which can include file paths or URLs.[3][7]

- Execution Engine: The resulting executable workflow is managed by an underlying execution engine, typically HTCondor's DAGMan, which handles job submission, dependency management, and error recovery.[8]

User Domain

Abstract Workflow
(DAX)

1. Submit

Pegasus Domain

Pegasus Planner
(pegasus-plan)

2. Query          2. Query          2. Query

Information Catalogs

Site
Catalog

Transformation
Catalog

Replica
Catalog

3. Generate

Execution Domain

Executable Workflow
(HTCondor DAG)

4. Submit

Execution Engine
(HTCondor DAGMan)

5. Execute Jobs

Compute Site
(Cluster, Cloud, Grid)

Pegasus Architecture: From Abstract to Executable Workflow

Click to download full resolution via product page
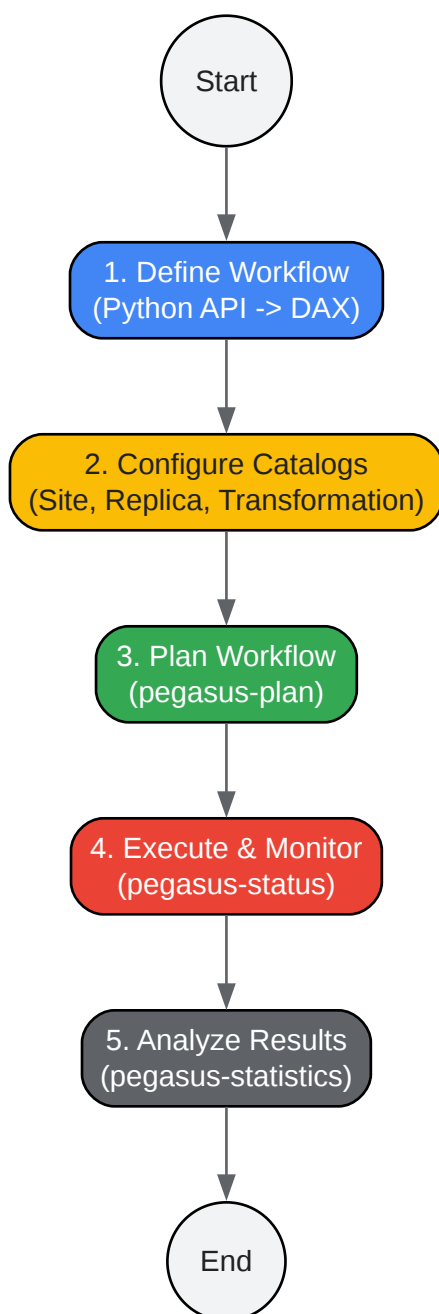
Caption: **Pegasus** architecture overview.

# General Protocol for Running a **Pegasus** Workflow

The following protocol outlines the high-level steps for executing any scientific workflow using **Pegasus** command-line tools.

Protocol Steps:

- Workflow Definition:

  - Write a script (e.g., dax-generator.py) using the **Pegasus** Python API to define the computational tasks, their dependencies, and their input/output files. This script generates the abstract workflow in a DAX file.[7]

- Catalog Configuration:

  - Site Catalog (sites.xml): Define the execution site(s), specifying the working directory, and the protocol for file transfers and job submission (e.g., local, HTCondor, SLURM).[7]

  - Replica Catalog (replicas.yml or .txt): For each logical input file name (LFN) required by the workflow, provide its physical file name (PFN), which is its actual location (e.g., file:///path/to/input.txt).[3][7]

  - Transformation Catalog (transformations.yml or .txt): For each logical executable name, define its physical path on the target site. Specify if the executable is pre-installed on the site or if it needs to be transferred.[11]

- Planning the Workflow:

  - Use the **pegasus**-plan command to map the abstract workflow to the execution site. This command takes the DAX file and catalogs as input and generates an executable workflow in a submit directory.

  - Command: **pegasus**-plan --dax my-workflow.dax --sites compute_site --output-site local --dir submit_dir --submit

- Execution and Monitoring:

- The --submit flag on **pegasus**-plan automatically sends the workflow to the execution engine.

- Monitor the workflow's progress using **pegasus**-status -v . This shows the status of jobs (e.g., QUEUED, RUNNING, SUCCEEDED, FAILED).[12]

- If the workflow fails, use **pegasus**-analyzer to diagnose the issue. The tool pinpoints the failed job and provides relevant error logs.[12]

- Analyzing Results and Provenance:

  - Once the workflow completes successfully, the final output files will be located in the directory specified during planning.

  - Use **pegasus**-statistics to generate a summary of the execution, including job runtimes, wait times, and data transfer volumes. This provenance data is crucial for performance analysis and reproducibility.[12]

Tech Support

General User Workflow Execution Flow

Caption: High-level steps for a **Pegasus** workflow.

# Application Protocol: 1000 Genomes Mutational Overlap Analysis

This protocol details a bioinformatics workflow that identifies mutational overlaps using data from the 1000 Genomes Project.[13] The workflow processes VCF (Variant Call Format) files to find common mutations across different individuals and chromosomes.
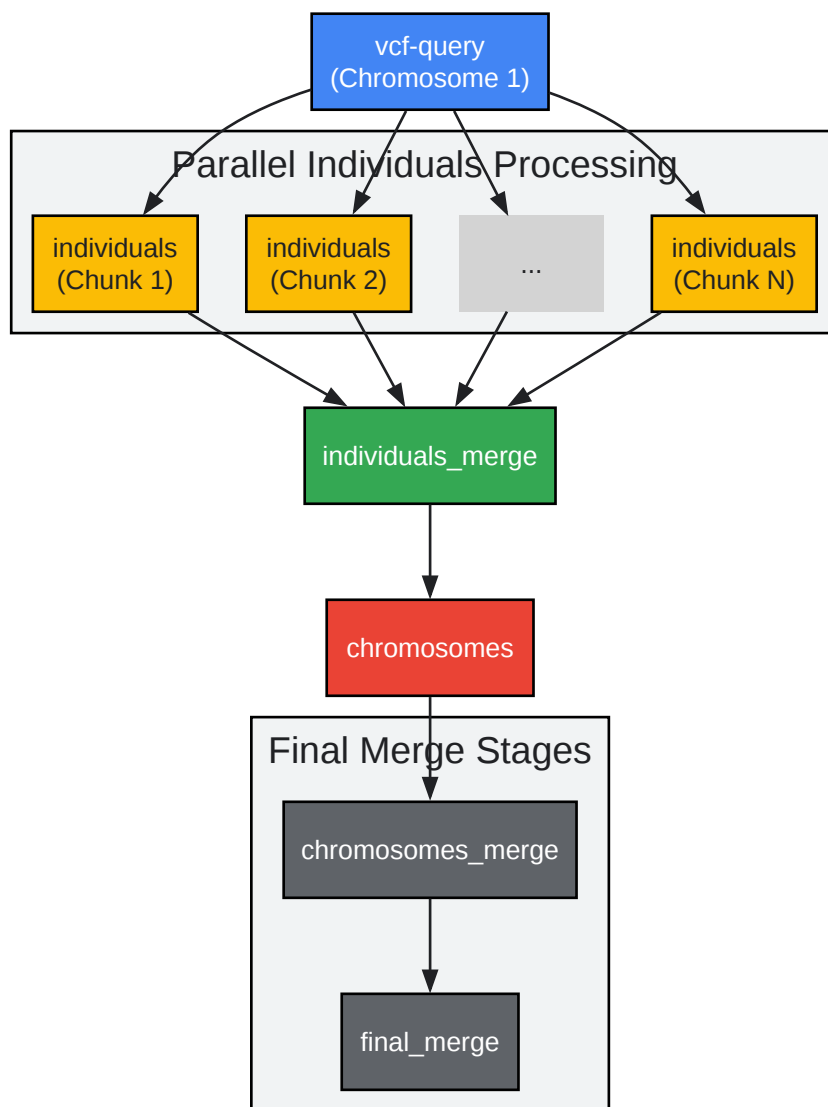
## Experimental Objective

To process a large genomic dataset in parallel to identify and merge mutational overlaps. The workflow is designed to be scalable, allowing for the processing of numerous chromosomes and individuals simultaneously.[13]

## Methodology and Workflow Structure

The workflow consists of several parallel and merge steps, creating a complex DAG structure.

Workflow Jobs:

- vcf-query: The initial step that queries a VCF file for a specific chromosome.

- individuals: This job processes chunks of the VCF file in parallel to identify mutations for a subset of individuals.[13]

- individuals_merge: Merges the parallel outputs from the individuals jobs for a single chromosome.

- chromosomes: Processes the merged data for each chromosome.

- chromosomes_merge: Merges the outputs from all chromosomes jobs.

- final_merge: A final step to combine all results into a single output.

1000 Genomes Mutational Overlap Workflow DAG

Click to download full resolution via product page

Caption: Job dependencies for the 1000 Genomes workflow.

# Execution Protocol

Prerequisites:

- **Pegasus** WMS version 5.0 or higher[13]

- Python version 3.6 or higher[13]

- HTCondor version 9.0 or higher[13]

- Access to an execution environment (e.g., local Condor pool, HPC cluster).

- Input data from the 1000 Genomes Project (VCF files).

Steps:

- Clone the Workflow Repository: git clone https://github.com/**pegasus**-isi/1000genome-workflow.git cd 1000genome-workflow

- Generate the Workflow (DAX):

  - A Python script (dax-generator.py) is provided to create the DAX file.

  - Execute the script, specifying the desired number of parallel individuals jobs and the target chromosome. For example, to create 10 parallel jobs for chromosome 22: ./dax-generator.py --individuals 10 --chromosome 22

- Configure Catalogs:

  - sites.xml: Modify this file to match your execution environment. The default is often a local HTCondor pool.

  - rc.txt: Update the replica catalog to point to the location of your input VCF files.

  - tc.txt: Ensure the transformation catalog correctly points to the paths of the workflow's executables (e.g., vcf-query).

- Plan and Submit:

  - Use the provided submit script or run **pegasus**-plan directly.

  - ./submit

  - This command plans the workflow, creating a submit directory (e.g., submit/user/**pegasus**/1000genome/run0001), and submits it to the local HTCondor scheduler.

Tech Support

- Monitor Execution:

  - Open a new terminal and monitor the workflow's progress: **pegasus**-status -v submit/user/**pegasus**/1000genome/run0001

  - Watch the jobs transition from READY to QUEUED, RUN, and finally SUCCESS.

# Quantitative Data Summary

The following table summarizes the execution time for a sample run of the 1000 Genomes workflow. The workflow was configured with 10 parallel individuals jobs for a single chromosome and executed on one Haswell node at the NERSC Cori supercomputer.[13]

| Job Class | Job Name | Wall Time (seconds) |
| --- | --- | --- |
| Compute | vcf-query | 13 |
| Compute | individuals | 10 |
| Compute | individuals_merge | 2 |
| Compute | chromosomes | 1 |
| Compute | chromosomes_merge | 1 |
| Compute | final_merge | 1 |
| Total Compute Time | 28 | |
| Auxiliary | Pegasus Internal Jobs | 10 |
| Total Workflow Time | 38 | |

Table Notes: For parallel jobs (e.g., individuals), the maximum duration among all parallel instances is reported. "Auxiliary" represents internal jobs managed by **Pegasus** for tasks like directory creation and cleanup. Data sourced from the **pegasus**-isi/1000genome-workflow GitHub repository.[13]

# References

- 1. Tutorials — Pegasus 1.10.2 documentation [pegasus.readthedocs.io]

- 2. arxiv.org [arxiv.org]

- 3. researchgate.net [researchgate.net]

- 4. Workflow gallery – Pegasus WMS [pegasus.isi.edu]

- 5. 5. Example Workflows — Pegasus WMS 5.1.2-dev.0 documentation [pegasus.isi.edu]

- 6. Pegasus Workflows with Application Containers — CyVerse Container Camp: Container Technology for Scientific Research 0.1.0 documentation [cyverse-container-camp-workshop-2018.readthedocs-hosted.com]

- 7. scielo.br [scielo.br]

- 8. Documentation – Pegasus WMS [pegasus.isi.edu]

- 9. youtube.com [youtube.com]

- 10. [2111.11624] Astronomical Image Processing at Scale With Pegasus and Montage [arxiv.org]

- 11. pegasus.isi.edu [pegasus.isi.edu]

- 12. GitHub - pegasus-isi/ACCESS-Pegasus-Examples: Pegasus Workflows examples including the Pegasus tutorial, to run on ACCESS resources. [github.com]

- 13. GitHub - pegasus-isi/1000genome-workflow: Bioinformatics workflow that identifies mutational overlaps using data from the 1000 genomes project [github.com]

- To cite this document: BenchChem. [Application Notes and Protocols for Pegasus Workflows]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b039198#step-by-step-guide-to-running-a-pegasus-workflow]

**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