

Automating Data Processing Pipelines with Pegasus WMS: Application Notes and Protocols

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

Compound of Interest

Compound Name: Pegasus

Cat. No.: B039198

[Get Quote](#)

For Researchers, Scientists, and Drug Development Professionals

This document provides detailed application notes and protocols for leveraging the **Pegasus** Workflow Management System (WMS) to automate complex data processing pipelines in scientific research and drug development. **Pegasus** WMS is a powerful open-source platform that enables the definition, execution, and monitoring of complex, multi-stage computational workflows across a variety of computing environments, from local clusters to national supercomputing centers and commercial clouds.^{[1][2][3]}

By abstracting the workflow from the underlying execution infrastructure, **Pegasus** allows researchers to focus on the scientific aspects of their data analysis, while the system handles the complexities of job scheduling, data management, fault tolerance, and provenance tracking.^{[4][5][6]} This leads to increased efficiency, reproducibility, and scalability of scientific computations.

Core Concepts of Pegasus WMS

Pegasus workflows are described as Directed Acyclic Graphs (DAGs), where nodes represent computational tasks and edges represent the dependencies between them.^[5] This model allows for the clear definition of complex data processing pipelines. Key features of **Pegasus** WMS include:

- **Portability and Reuse:** Workflows are defined in a resource-independent manner, allowing them to be executed on different computational infrastructures without modification.^{[1][3]}

- Scalability: **Pegasus** is designed to handle workflows of varying scales, from a few tasks to millions, processing terabytes of data.[\[1\]](#)[\[3\]](#)
- Data Management: The system automates the transfer of input and output data required by the different workflow tasks.[\[7\]](#)
- Performance Optimization: **Pegasus** can optimize workflow execution by clustering small, short-running jobs into larger ones to reduce overhead.[\[1\]](#)[\[8\]](#)
- Reliability and Fault Tolerance: It automatically retries failed tasks and can provide a "rescue" workflow for the remaining tasks in case of unrecoverable failures.[\[2\]](#)
- Provenance Tracking: Detailed information about the workflow execution, including the software, parameters, and data used, is captured to ensure reproducibility.[\[1\]](#)[\[3\]](#)[\[9\]](#)

Application Note 1: High-Throughput DNA Sequencing Analysis

This application note details a protocol for a typical high-throughput DNA sequencing (HTS) data analysis pipeline, automated using **Pegasus** WMS. This workflow is based on the practices of the USC Epigenome Center and is applicable to various research areas, including genomics, epigenomics, and personalized medicine.[\[10\]](#)

Experimental Protocol: DNA Sequencing Data Pre-processing

This protocol outlines the steps for pre-processing raw DNA sequencing data, starting from unmapped BAM files to produce an analysis-ready BAM file. The workflow leverages common bioinformatics tools like BWA for alignment and GATK4 for base quality score recalibration.[\[11\]](#)
[\[12\]](#)

1. Data Staging:

- Input: Unmapped BAM (.ubam) files.
- Action: Transfer the input files to the processing cluster's storage system. This is handled automatically by **Pegasus**.

- Tool: **Pegasus** data management tools.

2. Parallelization:

- Action: The input data is split into smaller chunks to be processed in parallel. This is a key feature of **Pegasus** for handling large datasets.[\[10\]](#)
- Tool: **Pegasus** job planner.

3. Sequence Alignment:

- Action: Each chunk of the unmapped data is aligned to a reference genome.
- Tool: BWA (mem)
- Exemplar Command:

4. Mark Duplicates:

- Action: Duplicate reads, which can arise from PCR artifacts, are identified and marked.
- Tool: GATK4 (MarkDuplicates)
- Exemplar Command:

5. Base Quality Score Recalibration (BQSR):

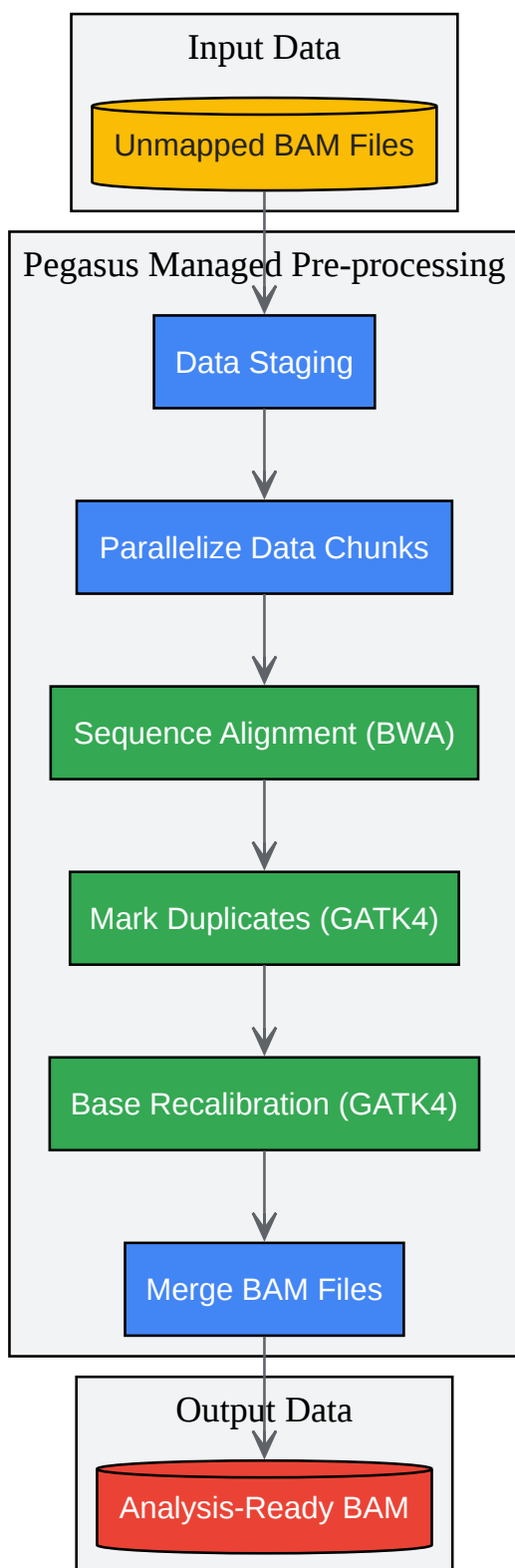
- Action: The base quality scores are recalibrated to provide more accurate quality estimations. This involves two steps: building a recalibration model and applying it.
- Tool: GATK4 (BaseRecalibrator, ApplyBQSR)
- Exemplar Commands:

6. Merge and Finalize:

- Action: The processed BAM files from all the parallel chunks are merged into a single, analysis-ready BAM file.[\[10\]](#)

- Tool: Samtools (merge)
- Exemplar Command:

Workflow Visualization



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

Caption: High-throughput DNA sequencing pre-processing workflow.

Application Note 2: Large-Scale Astronomical Image Mosaicking

This application note describes the use of **Pegasus** WMS to automate the creation of large-scale astronomical image mosaics using the Montage toolkit. This is a common task in astronomy for combining multiple smaller images into a single, scientifically valuable larger image.^[2]

Experimental Protocol: Astronomical Image Mosaicking with Montage

This protocol details the steps involved in creating a mosaic from a collection of astronomical images in the FITS format.

1. Define Region of Interest:

- Action: Specify the central coordinates and the size of the desired mosaic.
- Tool:montage-workflow.py script.^[13]
- Exemplar Command:

2. Data Discovery and Staging:

- Action: **Pegasus**, through the mArchiveList tool, queries astronomical archives to find the images that cover the specified region of the sky. These images are then staged for processing.
- Tool:mArchiveList

3. Re-projection:

- Action: Each input image is re-projected to a common coordinate system and pixel scale. This step is highly parallelizable and **Pegasus** distributes these tasks across the available compute resources.
- Tool:mProject

4. Background Rectification:

- Action: The background levels of the re-projected images are matched to a common level to ensure a seamless mosaic.
- Tool:mBgModel, mBgExec

5. Co-addition:

- Action: The background-corrected, re-projected images are co-added to create the final mosaic.
- Tool:mAdd

6. Image Generation (Optional):

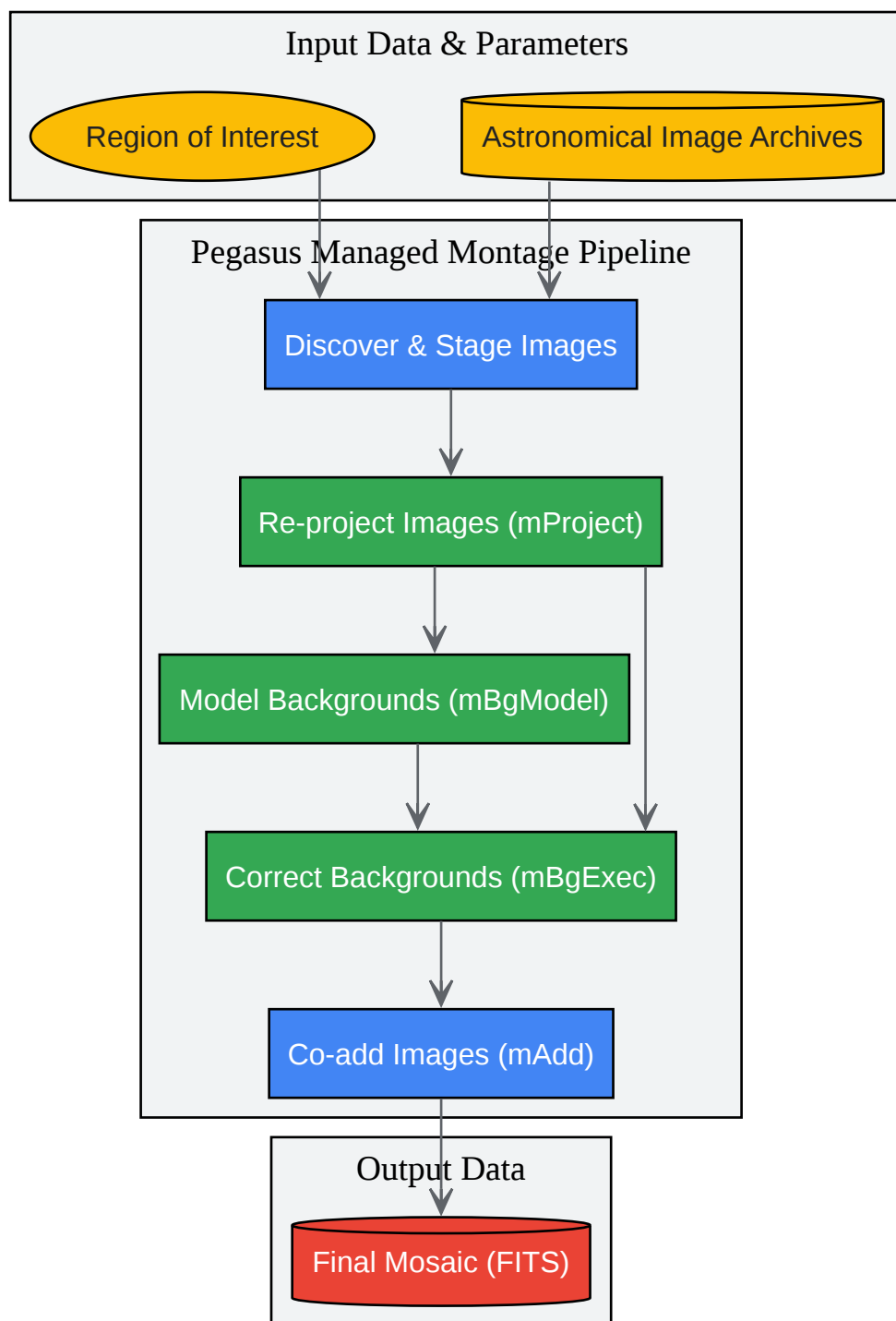
- Action: The final mosaic can be converted to a more common image format like JPEG for visualization.
- Tool:mJPEG

Quantitative Data

While specific performance metrics can vary greatly depending on the infrastructure and the size of the mosaic, the following table provides a conceptual overview of the scalability of **Pegasus**-managed Montage workflows.

Workflow Scale	Number of Input Images	Number of Tasks	Total Data Processed	Estimated Wall Time (on a 100-core cluster)
Small	100s	1,000s	10s of GB	< 1 hour
Medium	1,000s	10,000s	100s of GB	Several hours
Large	10,000s+	100,000s+	Terabytes	Days

Workflow Visualization



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

Caption: Astronomical image mosaicking workflow with Montage.

Application Note 3: A Representative Drug Target Identification Workflow

While there are no specific published examples of **Pegasus** WMS in a drug development pipeline, its capabilities are well-suited for automating the bioinformatics-intensive stages of early drug discovery, such as drug target identification.^{[14][15][16]} This application note presents a representative workflow for identifying potential drug targets from genomic and transcriptomic data, structured for execution with **Pegasus**.

Experimental Protocol: In-Silico Drug Target Identification

This protocol outlines a computational workflow to identify genes that are differentially expressed in a disease state and are predicted to be "druggable".

1. Data Acquisition and Pre-processing:

- Input: RNA-Seq data (FASTQ files) from disease and control samples.
- Action: Raw sequencing reads are pre-processed to remove low-quality reads and adapters.
- Tool: FastQC, Trimmomatic

2. Gene Expression Quantification:

- Action: The cleaned reads are aligned to a reference genome, and the expression level of each gene is quantified.
- Tool: STAR (aligner), RSEM (quantification)

3. Differential Expression Analysis:

- Action: Statistical analysis is performed to identify genes that are significantly up- or down-regulated in the disease samples compared to the controls.
- Tool: DESeq2 (R package)

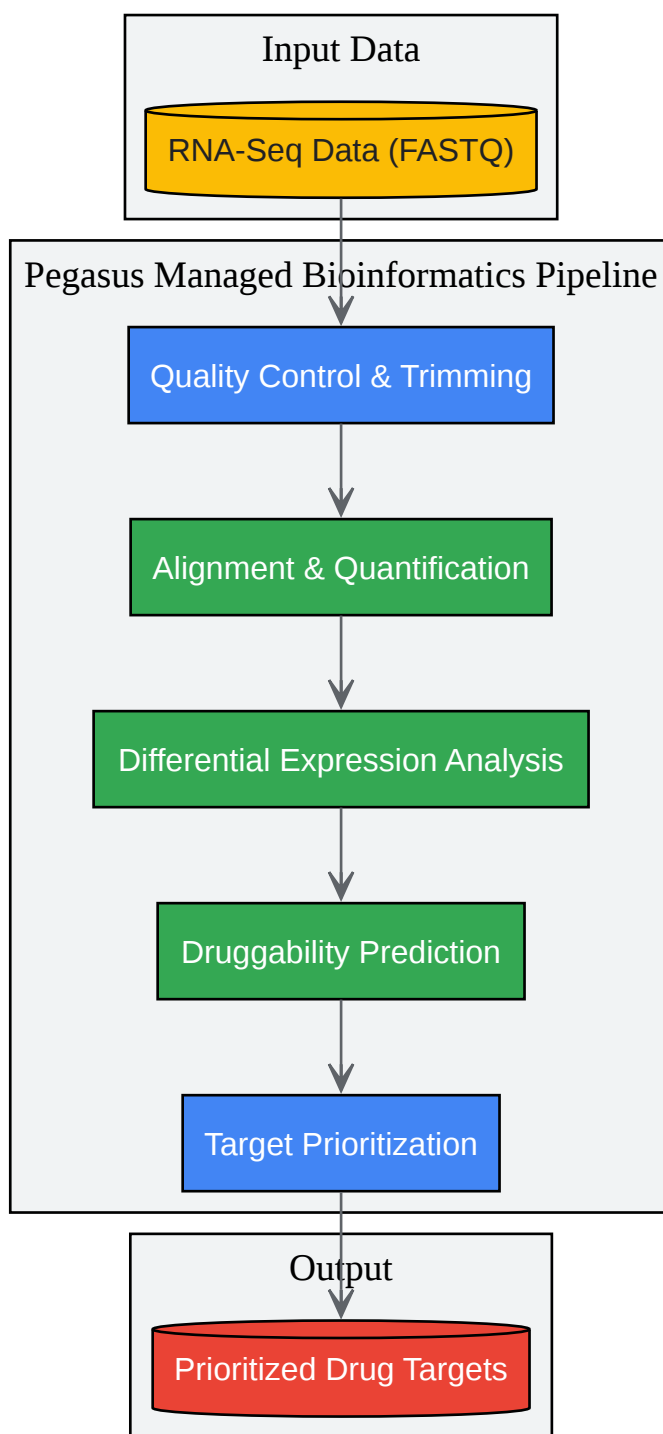
4. Druggability Prediction:

- Action: The differentially expressed genes are annotated with information from various databases to predict their potential as drug targets. This can include checking if they belong to gene families known to be druggable (e.g., kinases, GPCRs) or if they have known binding pockets.
- Tool: Custom scripts integrating data from databases like DrugBank, ChEMBL, and the Human Protein Atlas.

5. Target Prioritization:

- Action: The list of potential targets is filtered and ranked based on criteria such as the magnitude of differential expression, druggability score, and known association with the disease pathway.
- Tool: Custom analysis scripts.

Logical Relationship Visualization



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

Caption: A representative drug target identification workflow.

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. Pegasus WMS – Automate, recover, and debug scientific computations [pegasus.isi.edu]
- 3. GitHub - pegasus-isi/pegasus: Pegasus Workflow Management System - Automate, recover, and debug scientific computations. [github.com]
- 4. marketing.globuscs.info [marketing.globuscs.info]
- 5. rafaelsilva.com [rafaelsilva.com]
- 6. pegasus.isi.edu [pegasus.isi.edu]
- 7. research.cs.wisc.edu [research.cs.wisc.edu]
- 8. danielskatz.org [danielskatz.org]
- 9. research.cs.wisc.edu [research.cs.wisc.edu]
- 10. DNA Sequencing – Pegasus WMS [pegasus.isi.edu]
- 11. GitHub - gatk-workflows/gatk4-data-processing: Workflows for processing high-throughput sequencing data for variant discovery with GATK4 and related tools [github.com]
- 12. researchgate.net [researchgate.net]
- 13. PegasusHub [pegasushub.io]
- 14. Bioinformatics and Drug Discovery - PMC [pmc.ncbi.nlm.nih.gov]
- 15. Drug Discovery Workflow - What is it? [vipergen.com]
- 16. Target Identification and Validation in Drug Development | Technology Networks [technologynetworks.com]
- To cite this document: BenchChem. [Automating Data Processing Pipelines with Pegasus WMS: Application Notes and Protocols]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b039198#automating-data-processing-pipelines-with-pegasus-wms]

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