# Pegasus: An In-Depth Technical Guide to Single-Cell Analysis

**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 technical guide provides a comprehensive overview of the **Pegasus** Python package, a powerful and scalable tool for single-cell RNA sequencing (scRNA-seq) data analysis. **Pegasus**, developed as part of the Cumulus project, offers a rich set of functionalities for processing, analyzing, and visualizing large-scale single-cell datasets.[1] This document details the core workflow, experimental protocols, and data presentation, enabling users to effectively leverage **Pegasus** for their research and development needs.
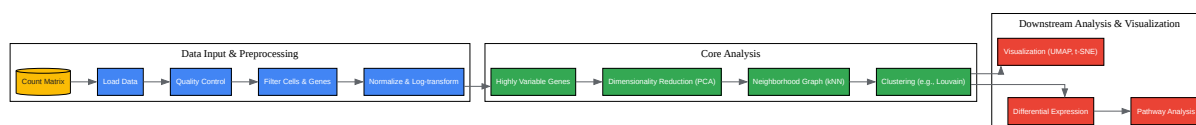
## Introduction to **Pegasus**

**Pegasus** is a command-line tool and a Python package designed for the analysis of transcriptomes from millions of single cells.[2] It is built upon the popular AnnData data structure, ensuring interoperability with the broader scverse ecosystem. **Pegasus** provides a comprehensive suite of tools covering the entire scRNA-seq analysis pipeline, from initial data loading and quality control to advanced analyses like differential gene expression and gene set enrichment.

## The **Pegasus** Workflow

The standard **Pegasus** workflow encompasses several key stages, each with dedicated functions to ensure robust and reproducible analysis. The typical progression involves data loading, quality control and filtering, normalization, identification of highly variable genes,

Tech Support

dimensionality reduction, cell clustering, and differential gene expression analysis to identify cluster-specific markers.



Click to download full resolution via product page

A high-level overview of the standard **Pegasus** single-cell analysis workflow.

# Experimental Protocols & Quantitative Data

This section provides detailed methodologies for the core steps in the **Pegasus** workflow, accompanied by tables summarizing key quantitative parameters.

## Data Loading

**Pegasus** supports various input formats, including 10x Genomics' Cell Ranger output, MTX, CSV, and TSV files. The **pegasus**.read_input function is the primary entry point for loading data into an AnnData object.

Experimental Protocol: Data Loading

- Purpose: To load the gene expression count matrix and associated metadata into memory.

- Methodology: Utilize the **pegasus**.read_input() function, specifying the file path and format. For 10x Genomics data, provide the path to the directory containing the matrix.mtx.gz, barcodes.tsv.gz, and features.tsv.gz files.

- Example Code:

# Quality Control and Filtering

Quality control (QC) is a critical step to remove low-quality cells and genes that could otherwise introduce noise into downstream analyses. **Pegasus** provides the pg.qc_metrics and pg.filter_data functions for this purpose.

Experimental Protocol: Quality Control and Filtering

- Purpose: To calculate QC metrics and filter out cells and genes based on these metrics.

- Methodology:

  - Calculate QC metrics using pg.qc_metrics(). This function computes metrics such as the number of genes detected per cell (n_genes), the total number of UMIs per cell (n_counts), and the percentage of mitochondrial gene expression (percent_mito).

  - Filter the data using pg.filter_data(). This function applies user-defined thresholds to remove cells and genes that do not meet the quality criteria.

- Example Code:

Table 1: Recommended Filtering Parameters

| Parameter | pegasus.qc_metrics argument | Description | Recommended Range |
|---|---|---|---|
| Minimum Genes per Cell | min_genes | The minimum number of genes detected in a cell. | 200 - 1000 |
| Maximum Genes per Cell | max_genes | The maximum number of genes detected in a cell to filter out potential doublets. | 3000 - 8000 |
| Mitochondrial Percentage | percent_mito | The maximum percentage of mitochondrial gene content. | 5 - 20 |
| Minimum Cells per Gene | (within pg.filter_data) | The minimum number of cells a gene must be expressed in to be retained. | 3 - 10 |

# Normalization and Highly Variable Gene Selection

Normalization adjusts for differences in sequencing depth between cells. Subsequently, identifying highly variable genes (HVGs) focuses the analysis on biologically meaningful variation.

Experimental Protocol: Normalization and HVG Selection

- Purpose: To normalize the data and identify genes with high variance across cells.

- Methodology:

  - Normalize the data using pg.log_norm(). This function performs total-count normalization and log-transforms the data.

  - Identify HVGs using pg.highly_variable_features(). **Pegasus** offers methods similar to Seurat for HVG selection.

- Example Code:

Table 2: Highly Variable Gene Selection Parameters

| Parameter | pegasus.highly_variable_features argument | Description | Default Value |
| --- | --- | --- | --- |
| Flavor | flavor | The method for HVG selection. | "seurat_v3" |
| Number of Top Genes | n_top_genes | The number of highly variable genes to select. | 2000 |

# Dimensionality Reduction and Clustering

Principal Component Analysis (PCA) is used to reduce the dimensionality of the data, followed by graph-based clustering to group cells with similar expression profiles.

Experimental Protocol: PCA and Clustering

- Purpose: To reduce the dimensionality of the data and identify cell clusters.

- Methodology:

  - Perform PCA on the highly variable genes using pg.pca().

  - Construct a k-nearest neighbor (kNN) graph using pg.neighbors().

  - Perform clustering on the kNN graph using algorithms like Louvain or Leiden (pg.louvain() or pg.leiden()).

- Example Code:

Table 3: PCA and Clustering Parameters

| Parameter | Function | Description | Default Value |
| :--- | :--- | :--- |
| Number of Principal Components | pg.pca | The number of principal components to compute. | 50 |
| Number of Neighbors | pg.neighbors | The number of nearest neighbors to use for building the kNN graph. | 15 |
| Resolution | pg.louvain / pg.leiden | The resolution parameter for clustering, which influences the number of clusters. | 1.0 |

# Differential Gene Expression and Visualization

Differential expression (DE) analysis identifies genes that are significantly upregulated in each cluster compared to all other cells. The results are often visualized using UMAP or t-SNE plots.

Experimental Protocol: DE Analysis and Visualization

- Purpose: To find marker genes for each cluster and visualize the cell populations.

- Methodology:

  - Perform DE analysis using pg.de_analysis(), specifying the cluster annotation.

  - Generate a UMAP embedding using pg.umap().

  - Visualize the clusters and gene expression on the UMAP plot using pg.scatter().

- Example Code:

# Signaling Pathway Analysis

**Pegasus** facilitates the analysis of signaling pathways and other gene sets through its gene set enrichment analysis (GSEA) and signature score calculation functionalities.

# Gene Set Enrichment Analysis (GSEA)

The **pegasus**.gsea() function allows for the identification of enriched pathways in the differentially expressed genes of each cluster.

Experimental Protocol: Gene Set Enrichment Analysis

- Purpose: To identify biological pathways that are significantly enriched in each cell cluster.
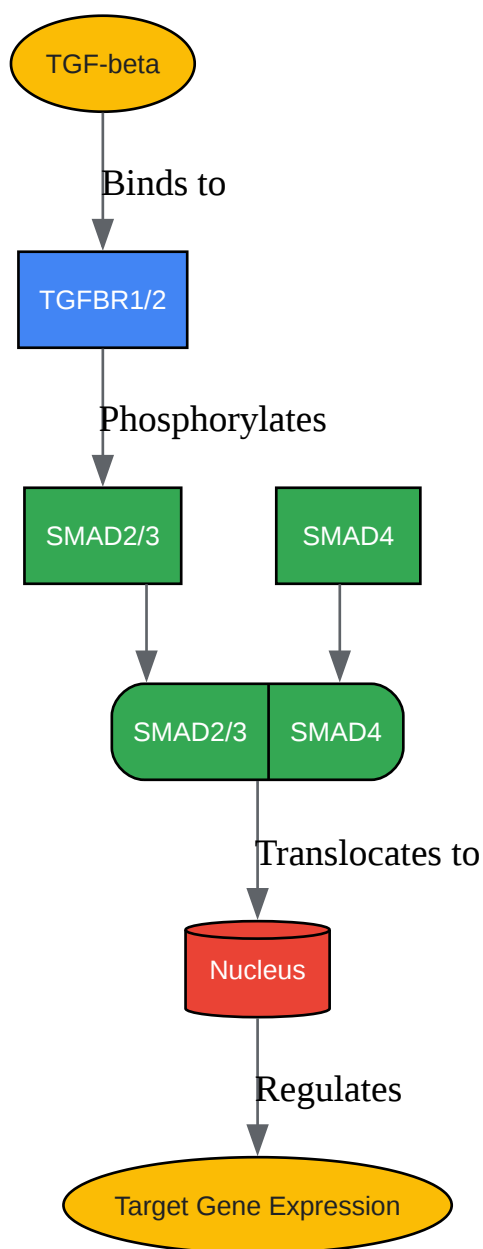
- Methodology:

  - Perform differential expression analysis as described in section 3.5.

  - Run pg.gsea(), providing the DE results and a gene set file in GMT format (e.g., from MSigDB).

- Example Code:

# Signature Score Calculation for a Signaling Pathway

The **pegasus**.calc_signature_score() function can be used to calculate a score for a given gene set (e.g., a signaling pathway) for each cell. This allows for the visualization of pathway activity across the dataset.

Hypothetical Example: Analysis of the TGF-β Signaling Pathway

The TGF-β signaling pathway plays a crucial role in various cellular processes. We can define a gene set representing this pathway and analyze its activity.

Tech Support

Click to download full resolution via product page

A simplified diagram of the TGF-β signaling pathway.

Experimental Protocol: TGF-β Pathway Activity Score

- Purpose: To quantify the activity of the TGF-β signaling pathway in each cell.

- Methodology:

  - Define a list of genes belonging to the TGF-β pathway.

- Use **pegasus**.calc_signature_score() to calculate a score for this gene set.

- Visualize the signature score on a UMAP plot using pg.scatter().

- Example Code:

# Conclusion

**Pegasus** provides a robust and user-friendly framework for the analysis of large-scale single-cell RNA sequencing data. Its comprehensive functionalities, scalability, and integration with the Python ecosystem make it an invaluable tool for researchers and scientists in both academic and industrial settings. This guide has outlined the core workflow and provided detailed protocols to enable users to effectively apply **Pegasus** to their own single-cell datasets. For more detailed information, users are encouraged to consult the official **Pegasus** documentation.

---

**Need Custom Synthesis?**

*BenchChem offers custom synthesis for rare earth carbides and specific isotopiclabeling.*

*Email: info@benchchem.com or Request Quote Online.*

---

# References

- 1. Cumulus provides cloud-based data analysis for large-scale single-cell and single-nucleus RNA-seq - PMC [pmc.ncbi.nlm.nih.gov]

- 2. GitHub - lilab-bcb/pegasus: A tool for analyzing trascriptomes of millions of single cells. [github.com]

- To cite this document: BenchChem. [Pegasus: An In-Depth Technical Guide to Single-Cell Analysis]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b039198#pegasus-single-cell-analysis-python-package-tutorial]

---

**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