

# Configuring Pegasus for High-Throughput Drug Discovery in the Cloud

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

## Compound of Interest

Compound Name: Pegasus

Cat. No.: B039198

[Get Quote](#)

Application Notes and Protocols for Researchers, Scientists, and Drug Development Professionals

This document provides detailed application notes and protocols for configuring and utilizing the **Pegasus** Workflow Management System (WMS) in cloud computing environments for drug discovery research. **Pegasus** is an open-source platform that enables the automation and execution of complex scientific workflows across a variety of computational infrastructures, including commercial and academic clouds.[1][2] By abstracting the workflow from the underlying execution environment, **Pegasus** allows researchers to define complex computational pipelines that are portable, scalable, and resilient to failures.[1][3] These capabilities are particularly advantageous for computationally intensive tasks common in drug discovery, such as virtual screening and molecular dynamics simulations.

## Introduction to Pegasus in Cloud Environments

**Pegasus** facilitates the execution of scientific workflows on Infrastructure-as-a-Service (IaaS) clouds, such as Amazon Web Services (AWS), Google Cloud Platform (GCP), and Microsoft Azure.[4] It achieves this by creating a virtual cluster on the cloud, which consists of virtual machines configured with the necessary software, such as the HTCondor high-throughput computing system.[1] This approach provides researchers with a familiar cluster environment while leveraging the on-demand scalability and resource flexibility of the cloud.

A key aspect of using **Pegasus** in the cloud is its robust data management capabilities.

**Pegasus** can be configured to work with various data storage solutions, including cloud-native

object storage services like Amazon S3 and distributed file systems like GlusterFS.<sup>[4]</sup> It automatically manages the staging of input data required for workflow tasks and the transfer of output data back to a designated storage location.<sup>[1]</sup>

## Configuring Pegasus on a Cloud Platform

Configuring **Pegasus** for a cloud environment involves several key steps, from setting up the cloud resources to configuring the **Pegasus** workflow management system. The following protocol outlines a general approach for configuring **Pegasus** on a cloud platform, using AWS as an example.

### Protocol: Setting up a Virtual Cluster on AWS for Pegasus

Objective: To create a virtual cluster on Amazon EC2 that can be used to execute **Pegasus** workflows.

Materials:

- An Amazon Web Services (AWS) account.
- A submit host (a local machine or a small, persistent EC2 instance) with **Pegasus** and HTCondor installed.
- A virtual machine (VM) image (Amazon Machine Image - AMI) with HTCondor and the necessary scientific software pre-installed.

Methodology:

- Prepare the Submit Host:
  - Install and configure the **Pegasus** WMS and HTCondor on your designated submit host. This machine will be used to plan and submit your workflows.
  - Configure the AWS Command Line Interface (CLI) with your AWS credentials.
- Create a Custom AMI:

- Launch a base Amazon Linux or Ubuntu EC2 instance.
- Install HTCondor and configure it to join the Condor pool managed by your submit host.
- Install the scientific applications required for your workflow (e.g., AutoDock Vina for virtual screening).
- Create an Amazon Machine Image (AMI) from this configured instance. This AMI will be used to launch worker nodes in your virtual cluster.
- Configure **Pegasus** for AWS:
  - On the submit host, configure the **Pegasus** site catalog to describe the AWS resources. This includes specifying the AMI ID of your custom AMI, the desired instance type, and the security group.
  - Configure the replica catalog to specify the location of your input data. For cloud environments, it is recommended to store input data in an object store like Amazon S3.
  - Configure the transformation catalog to define the logical names of your executables and where they are located on the worker nodes.
- Define the Workflow:
  - Define your scientific workflow as a Directed Acyclic Graph (DAG) using the **Pegasus** Python API or another supported format.<sup>[5]</sup> This abstract workflow will describe the computational tasks and their dependencies.
- Plan and Execute the Workflow:
  - Use the **pegasus-plan** command to map the abstract workflow to the AWS resources defined in your site catalog. **Pegasus** will generate an executable workflow that includes jobs for data staging, computation, and data registration.<sup>[2]</sup>
  - Use the **pegasus-run** command to submit the executable workflow to HTCondor for execution on your virtual cluster.

## Application: High-Throughput Virtual Screening for Drug Discovery

Virtual screening is a computational technique used in drug discovery to search large libraries of small molecules to identify those that are most likely to bind to a drug target, typically a protein receptor or enzyme. This process can be computationally intensive, making it an ideal candidate for execution on the cloud using **Pegasus**.

### Experimental Protocol: Virtual Screening with Pegasus and AutoDock Vina on AWS

Objective: To perform a high-throughput virtual screening of a compound library against a protein target using a **Pegasus** workflow on AWS.

Methodology:

- Prepare the Input Files:
  - Receptor: Prepare the 3D structure of the target protein in PDBQT format. This is the format required by AutoDock Vina.
  - Compound Library: Obtain a library of small molecules in a format that can be converted to PDBQT, such as SMILES or SDF.
  - Configuration File: Create a configuration file for AutoDock Vina that specifies the search space (the region of the receptor to be docked against) and other docking parameters.
  - Upload all input files to an Amazon S3 bucket.
- Define the **Pegasus** Workflow:
  - The workflow will consist of the following main steps:
    - A "split" job that divides the large compound library into smaller chunks.
    - Multiple "docking" jobs that run in parallel, each processing one chunk of the compound library. Each docking job will use AutoDock Vina to dock the compounds to the receptor.

- A "merge" job that gathers the results from all the docking jobs and combines them into a single output file.
- A "rank" job that sorts the docked compounds based on their binding affinity scores to identify the top candidates.
- Execute and Monitor the Workflow:
  - Plan and run the workflow using the **pegasus-plan** and **pegasus-run** commands as described in the previous protocol.
  - Monitor the progress of the workflow using **pegasus-status** and other monitoring tools provided by **Pegasus**.

## Quantitative Data and Performance

The performance and cost of running **Pegasus** workflows in the cloud can vary depending on the cloud provider, the types of virtual machines used, and the data storage solution. The following tables provide an illustrative comparison of different configurations.

Table 1: Illustrative Performance of a Virtual Screening Workflow

| Cloud Provider | VM Instance Type | Number of VMs | Workflow Wall Time (hours) |
|----------------|------------------|---------------|----------------------------|
| AWS            | c5.2xlarge       | 10            | 5.2                        |
| GCP            | n2-standard-8    | 10            | 4.9                        |
| Azure          | Standard_F8s_v2  | 10            | 5.5                        |

Note: The data in this table is illustrative and will vary based on the specific workflow, dataset size, and other factors.

Table 2: Illustrative Cost Comparison for a 100-Hour Virtual Screening Workflow

| Cloud Provider | VM Instance Type (On-Demand) | Cost per Hour per VM | Total Estimated Cost |
|----------------|------------------------------|----------------------|----------------------|
| AWS            | c5.2xlarge                   | \$0.34               | \$340                |
| GCP            | n2-standard-8                | \$0.38               | \$380                |
| Azure          | Standard_F8s_v2              | \$0.39               | \$390                |

Note: Cloud provider pricing is subject to change. This table does not include costs for data storage and transfer. Significant discounts can be achieved using spot instances or reserved instances.[\[6\]](#)[\[7\]](#)[\[8\]](#)

Table 3: Data Staging Performance Comparison

| Storage Solution | Throughput for Large Files | Latency for Small Files | Cost   |
|------------------|----------------------------|-------------------------|--------|
| Amazon S3        | High                       | Higher                  | Lower  |
| GlusterFS on EBS | Moderate                   | Lower                   | Higher |

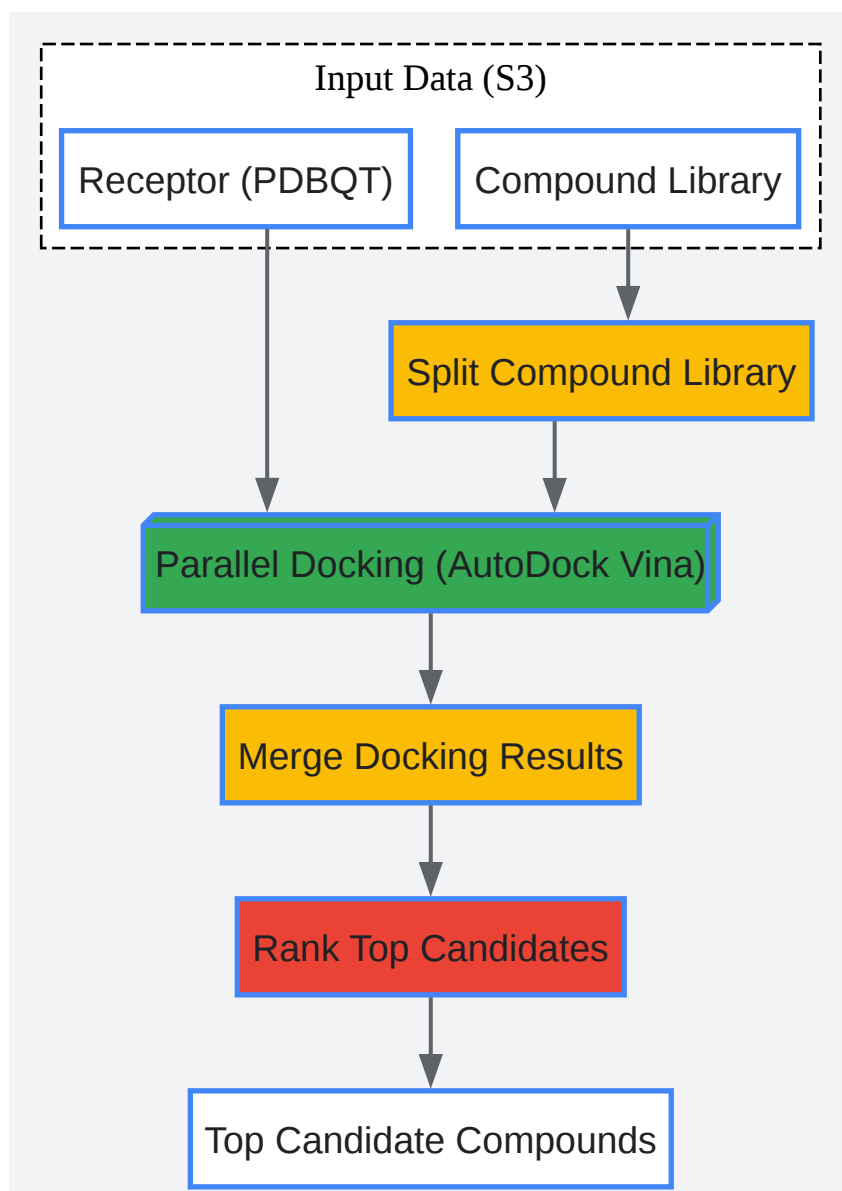
Note: The choice of storage solution depends on the specific I/O patterns of the workflow. Object stores like S3 are generally more cost-effective and scalable for large datasets.[\[4\]](#)

## Visualizing Workflows and Signaling Pathways

Visual representations are crucial for understanding complex workflows and biological pathways. **Pegasus** workflows can be visualized as Directed Acyclic Graphs (DAGs), and signaling pathways relevant to drug discovery can be modeled to identify potential targets.

## Virtual Screening Experimental Workflow

The following diagram illustrates the logical flow of the virtual screening workflow described in the protocol.

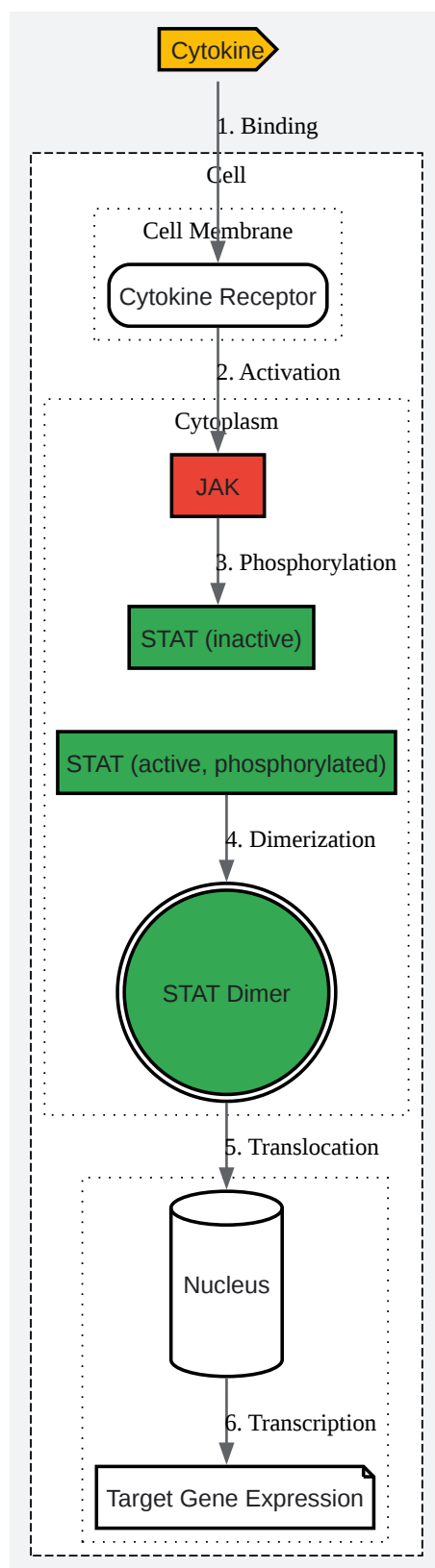


[Click to download full resolution via product page](#)

A high-level workflow for virtual screening using **Pegasus**.

## JAK-STAT Signaling Pathway

The Janus kinase (JAK) and signal transducer and activator of transcription (STAT) signaling pathway is a critical pathway in the regulation of the immune system.<sup>[9][10][11]</sup> Its dysregulation is implicated in various diseases, making it a significant target for drug discovery.<sup>[12]</sup>



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

The canonical JAK-STAT signaling pathway.



## Conclusion

**Pegasus** provides a powerful and flexible framework for orchestrating complex drug discovery workflows in cloud computing environments. By leveraging the scalability and on-demand resources of the cloud, researchers can significantly accelerate their research and development efforts. The ability to define portable and reproducible workflows also enhances collaboration and ensures the integrity of scientific results. While the initial setup and configuration require some effort, the long-term benefits of using a robust workflow management system like **Pegasus** for drug discovery research are substantial.

### Need Custom Synthesis?

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

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

## References

- 1. rafaelsilva.com [rafaelsilva.com]
- 2. arokem.github.io [arokem.github.io]
- 3. 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]
- 4. Pegasus in the Cloud – Pegasus WMS [pegasus.isi.edu]
- 5. youtube.com [youtube.com]
- 6. AWS vs. GCP: A Comprehensive Pricing Breakdown for 2025 - Hykell [hykell.com]
- 7. Cloud Pricing Comparison: AWS vs. Azure vs. Google in 2025 [cast.ai]
- 8. cloudzero.com [cloudzero.com]
- 9. Small molecule drug discovery targeting the JAK-STAT pathway | CoLab [colab.ws]
- 10. Small molecule drug discovery targeting the JAK-STAT pathway - PubMed [pubmed.ncbi.nlm.nih.gov]
- 11. bocsci.com [bocsci.com]
- 12. researchgate.net [researchgate.net]

- To cite this document: BenchChem. [Configuring Pegasus for High-Throughput Drug Discovery in the Cloud]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b039198#configuring-pegasus-for-cloud-computing-environments]

---

### 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](mailto:info@benchchem.com)