

Learn How to Build Version Controlled End-to-End Data Pipelines Using Pachyderm

Data pipelines are essential for organizations that want to make data-driven decisions. They allow you to collect, process, and transform data from various sources and make it available to downstream systems. However, building and managing data pipelines can be complex and time-consuming, especially when you need to ensure data integrity and reproducibility.



Reproducible Data Science with Pachyderm: Learn how to build version-controlled, end-to-end data pipelines using Pachyderm 2.0 by Svetlana Karslioglu

★★★★★ 5 out of 5

Language	: English
File size	: 11815 KB
Text-to-Speech	: Enabled
Screen Reader	: Supported
Enhanced typesetting	: Enabled
Print length	: 364 pages
Paperback	: 200 pages
Item Weight	: 11.2 ounces
Dimensions	: 5.5 x 0.5 x 8.5 inches



Pachyderm is an open-source data management platform that makes it easy to build and manage data pipelines. It provides a version controlled environment for your data and pipelines, allowing you to track changes, roll back to previous versions, and collaborate with others on data projects.

In this article, we will provide a comprehensive guide on how to build version controlled end-to-end data pipelines using Pachyderm. We will cover the key concepts of Pachyderm, its architecture, and the steps involved in building a version controlled data pipeline using Pachyderm.

Key Concepts of Pachyderm

Before we dive into building data pipelines with Pachyderm, let's first understand some of the key concepts of Pachyderm.

- **Repos**: A repo is a collection of data and pipelines. Repos can be public or private, and they can be shared with other users.
- **Pipelines**: A pipeline is a set of steps that transform data from one format to another. Pipelines can be simple or complex, and they can be used to perform a variety of data processing tasks.
- **Datasets**: A dataset is a collection of data that is stored in a repo. Datasets can be created from a variety of sources, such as files, databases, or other repos.
- **Versions**: Every change to a repo, pipeline, or dataset is tracked as a version. This allows you to roll back to previous versions if necessary.
- **Pachyderm Client**: The Pachyderm client is a command-line tool that you can use to interact with Pachyderm. The client can be used to create and manage repos, pipelines, and datasets.

Pachyderm Architecture

Pachyderm is built on a distributed architecture that consists of the following components:

- ****Pachyderm Coordinator****: The coordinator is the central component of Pachyderm. It manages repos, pipelines, and datasets, and it orchestrates the execution of pipelines.
- ****Pachyderm Workers****: The workers are responsible for executing pipeline steps. Workers can be deployed on-premises or in the cloud.
- ****Pachyderm Storage****: Pachyderm storage is a distributed file system that is used to store data and pipeline artifacts.
- ****Pachyderm Client****: The client is used to interact with Pachyderm from the command line.

Building a Version Controlled Data Pipeline with Pachyderm

Now that we have a basic understanding of Pachyderm, let's walk through the steps involved in building a version controlled data pipeline using Pachyderm.

1. Create a Repo

The first step is to create a repo. A repo can be created using the following command:

```
pachyderm init repo my-repo
```

2. Create a Dataset

Next, we need to create a dataset. We can create a dataset from a variety of sources, including files, databases, or other repos. To create a dataset from a file, we can use the following command:

```
pachyderm create dataset my-dataset --file my-data.csv
```

3. Create a Pipeline

Once we have a dataset, we can create a pipeline to process the data. A pipeline is a set of steps that transform data from one format to another. To create a pipeline, we can use the following command:

```
pachyderm create pipeline my-pipeline
```

4. Define the Pipeline Steps

Once we have created a pipeline, we need to define the steps that will be executed in the pipeline. A pipeline step is a function that takes a dataset as input and returns a new dataset as output. To define a pipeline step, we can use the following command:

```
pachyderm add step my-pipeline my-step
```

5. Run the Pipeline

Once we have defined the pipeline steps, we can run the pipeline using the following command:

```
pachyderm run my-pipeline
```

6. Version the Pipeline

Every time we make a change to our repo, pipeline, or dataset, we should create a new version. This will allow us to roll back to previous versions if necessary. To create a new version, we can use the following command:

```
pachyderm create version my-repo my-version
```

7. Roll Back to a Previous Version

If we need to roll back to a previous version, we can use the following command:

```
pachyderm rollback my-repo my-version
```

In this article, we have provided a comprehensive guide on how to build version controlled end-to-end data pipelines using Pachyderm. We have covered the key concepts of Pachyderm, its architecture, and the steps involved in building a version controlled data pipeline using Pachyderm. By following the steps outlined in this article, you can quickly and easily build data pipelines that are reliable, reproducible, and easy to manage.



Reproducible Data Science with Pachyderm: Learn how to build version-controlled, end-to-end data pipelines using Pachyderm 2.0

by Svetlana Karslioglu

★★★★★ 5 out of 5

Language : English
File size : 11815 KB
Text-to-Speech : Enabled
Screen Reader : Supported
Enhanced typesetting : Enabled
Print length : 364 pages
Paperback : 200 pages
Item Weight : 11.2 ounces
Dimensions : 5.5 x 0.5 x 8.5 inches





Unlocking the Power of Celebrity Branding: A Comprehensive Guide by Nick Nanton

In the ever-evolving marketing landscape, celebrity branding has emerged as a potent force, captivating audiences and driving brand success. From...



The Legendary Riggins Brothers: Play-by-Play of a Football Dynasty

The Unforgettable Trio: The Impact of the Riggins Brothers on Football
The Riggins brothers, Lorenzo "Zo" and Thomas "Tom," are revered as icons in the annals...