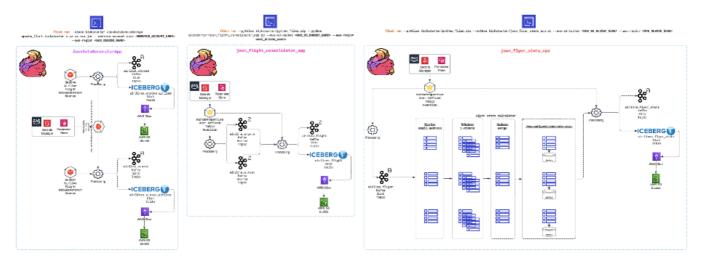


## Apache Flink v2.1 Kickstarter

This project showcases the exceptional capabilities of Apache Flink, known for its high performance, scalability, and advanced stream processing features, which is a core component of the signalRoom technology stack. Staying true to our mission of empowering developers by sharing practical knowledge, we demonstrate the use of Flink through real-world examples based on the blog series "Building Apache Flink Applications in Java."



#### Flink App

#### Description

Data Generator App This Flink App generates realistic flight data for the fictional airlines Sunset Air and Sky One Airlines, seamlessly integrating modern data streaming technologies. Flight events are published to dedicated Kafka topics (airline.sunset and airline.skyone), enabling real-time processing and analysis. Simultaneously, the synthetic data flows into the apache\_kickstarter.airlines.skyone\_airline and apache\_kickstarter.airlines.sunset\_airline Apache Iceberg Tables, respectively, stored in an AWS S3 bucket, providing a robust, scalable foundation for historical data analytics. The entire solution is implemented in Java, overcoming limitations in PyFlink, which currently lacks a Python-based Apache Flink data generator source. This design ensures a powerful, flexible data pipeline, ideal for both real-time and batch use cases.

Flight Consolidator App Imports flight data from airline.sunset and airline.skyone Kafka topics and standardizes them into a unified airline.flight Kafka topic, and the apache\_kickstarter.airlines.flight Apache Iceberg Table. Implemented in both Java and Python.

Flyer Stats App Processes data from the airline.flight Kafka topic to aggregate flyer statistics into the airline.flyer\_stats Kafka topic, and the apache\_kickstarter.airlines.flyer\_stats Apache Iceberg Table. Implemented in both Java and Python.

Created by Wade Waldron, Staff Software Practice Lead at Confluent Inc., these applications are designed to handle enterprise-scale use cases to demonstrate Flink's potential for real-world applications. This will start with securely fetching Kafka Cluster and Schema Registry Cluster API keys via *AWS Secrets Manager*, and retrieving essential Kafka client configuration properties through the *AWS Systems Manager Parameter*.

Beyond simply transforming and enriching data, we stream it into Kafka topics and *Apache Iceberg Tables*, designed for high-performance, large-scale analytics with support for ACID transactions, offering highly scalable, real-time data processing and durable storage.

Our journey with Flink doesn't end with Java; Python-based Flink applications also shine in this project, leveraging Flink SQL, Table API, and DataFrame API, each offering unique advantages: Flink SQL for declarative queries, Table API for a mix of SQL-like syntax and programmatic control, and DataFrame API for efficient data manipulation similar to pandas or R. This provides a seamless, flexible approach to solving complex streaming challenges, ultimately unlocking an efficient and scalable architecture for modern data needs.

Explore the repo, run the Docker containers (compatible with both Mac M chips and non-Mac M chips), and see how Apache Flink can be utilized to build robust, enterprise-grade streaming applications.

#### **Table of Contents**

- 1.0 Are these examples better described as Flink Jobs or Flink Applications—and why does the distinction matter?
- 2.0 Let's get started!
  - 2.1 DevOps in Action: Running Terraform Locally
    - 2.1.1 Run locally
  - 2.2 DevOps in Action: Running Terraform in the cloud
    - 2.2.1 Run from the cloud
  - 2.3 Visualizing the Terraform Configuration
- 3.0 Hands-On Kickoff: Practical Examples for Rapid Learning
  - o 3.1 Flink Applications Powered by Java on a locally running Apache Flink Cluster in Docker
  - 3.2 Flink Applications Powered by Python on a locally running Apache Flink Cluster in Docker
  - 3.3 Flink Applications Powered by Python on Confluent Cloud for Apache Flink
- 4.0 Powering Analytics with Apache Iceberg tables in Snowflake
- 5.0 Resources
- 6.0 Important Note(s)

# 1.0 Are these examples better described as Flink Jobs or Flink Applications—and why does the distinction matter?

"What's in a name? That which we call a rose by any other name would smell just as sweet."

-- William Shakespeare

Did you know Flink jobs are often called Flink applications? It's not just a fancy name upgrade—there's a good reason behind it! Unlike a single task or computation, a Flink job is like a mini software ecosystem. Calling it an 'application' captures its true essence: a complex, integrated workflow that solves big data challenges, just like any full-fledged software app.

By using 'Flink applications,' we're emphasizing their depth and sophistication. It's not just about crunching numbers; it's about building something comprehensive, dynamic, and powerful—just like any great app you use every day. (Curious for more details? Check out the rationale here!)

### 2.0 Let's get started!

Ready to supercharge your data streaming skills? As of October 4, 2024, Apache Flink on Confluent Cloud brings incredible possibilities and a few challenges! Currently, it doesn't support the DataStream API, and its Table API is limited (click here for a list of what is supported) to creating Java or Python-based Flink applications. Don't worry, though—we've got you covered!

In this repo, you'll find an example code where we use the Java DataStream API to create a *custom source* beyond Flink's typical capabilities. Imagine making external calls to AWS Secrets Manager to securely retrieve Kafka Cluster API Key secrets or pulling Kafka consumer and producer client configuration properties from the AWS Systems Manager Parameter Store—all directly integrated into Flink!

Prefer Python? We've got you there too. With the Table API, you can use a *User-Defined Table Function* to do the same. Plus, we don't just sink data into Kafka topics; this repo showcases how to store it in *Apache Iceberg tables*, unlocking even more power for your data.

To help you start quickly, the repo comes with **Docker containers** for Mac M chip and non-Mac M chip machines, letting you run an Apache Flink 1.20.0 cluster locally. For the Kafka and Schema Registry resources, we're fully leveraging Confluent Cloud. Dive in and explore—the future of data streaming is here!

#### These are the steps

1. Take care of the cloud and local environment prequisities listed below:

You need to have the following cloud accounts:

- AWS Account with SSO configured
- Confluent Cloud Account
- Docker Account
- GitHub Account with OIDC configured for AWS
- Snowflake Account
- Terraform Cloud Account

You need to have the following installed on your local machine:

- o AWS CLI version 2
- Confluent CLI version 3 or higher
- Docker Desktop
- Java JDK (Java Development Kit) 17
- Python 3.11
- Terraform CLI version 1.9.3 or higher

#### 2. Clone the repo:

```
git clone https://github.com/j3-signalroom/apache_flink-
kickstarter.git
```

3. Set up your Terraform Cloud environment locally or leverage GitHub Actions to create the complete setup effortlessly. Here's what you can expect:

- A Confluent Cloud environment featuring a Kafka Cluster, fully equipped with pre-configured example Kafka topics—ready to power your data streaming needs.
- AWS Secrets Manager securely storing API Key Secrets for the Kafka Cluster, along with AWS Systems Manager Parameter Store containing Kafka Consumer and Producer properties for easy integration.
- An AWS S3 bucket with a dedicated warehouse folder, serving as the landing zone for Apache lceberg Tables populated by two Python-based Flink apps, bringing your data streaming architecture to life.
- 4. Run Apache Flink locally on your Mac, or use the provided Docker containers from the project to launch Apache Flink and Apache Iceberg seamlessly on your machine.

With these steps, you'll have everything set up to run enterprise-grade data streaming applications in no time!

#### 2.1 DevOps in Action: Running Terraform Locally

Install the Terraform CLI on your local machine, and make sure you have an HCP Terraform account to run the Terraform configuration. Learn how to set up Terraform Cloud for local use by clicking here.

#### 2.1.1 Run locally

```
scripts/deploy-terraform.sh <create | delete> --profile=<SSO_PROFILE_NAME>
                                                    --confluent-api-key=
<CONFLUENT_API_KEY>
                                                    --confluent-api-secret=
<CONFLUENT_API_SECRET>
                                                    --snowflake-warehouse=
<SNOWFLAKE_WAREHOUSE>
                                                    --service-account-user=
<SERVICE_ACCOUNT_USER>
                                                    --day-count=<DAY COUNT>
                                                    --auto-offset-reset=
<earliest | latest>
                                                    --number-of-api-keys-
to-retain=<NUMBER_OF_API_KEYS_TO_RETAIN>
                                                    --admin-service-user-
secrets-root-path=<ADMIN_SERVICE_USER_SECRETS_ROOT_PATH>
```

Argument placeholder	Replace with
----------------------	--------------

Argument placeholder	Replace with
<sso_profile_name></sso_profile_name>	your AWS SSO profile name for your AWS infrastructue that host your AWS Secrets Manager.
<confluent_api_key></confluent_api_key>	your organization's Confluent Cloud API Key (also referred as Cloud API ID).
<confluent_api_secret></confluent_api_secret>	your organization's Confluent Cloud API Secret.
<snowflake_warehouse></snowflake_warehouse>	the Snowflake warehouse (or "virtual warehouse") you choose to run the resources in Snowflake.
<service_account_user></service_account_user>	the Snowflake service account user who is to be assigned the RSA key pairs for its authentication.
<day_count></day_count>	how many day(s) should the API Key be rotated for.
<auto_offset_reset></auto_offset_reset>	Use earliest, when you want to read the first event in a Kafka topic. Otherwise, specify latest.
<number_of_api_keys_to_retain></number_of_api_keys_to_retain>	Specifies the number of API keys to create and retain.
<admin_service_user_secrets_root_path></admin_service_user_secrets_root_path>	the root path in AWS Secrets Manager for admin service user secrets.

To learn more about this script, click here.

#### 2.2 DevOps in Action: Running Terraform in the cloud

In order to run the Terraform configuration from GitHub, the Terraform Cloud API token and Confluent Cloud API Key are required as GitHub Secret variables. Learn how to do to get the Terraform Cloud API token and Confluent Cloud API key here.

#### 2.2.1 Run from the cloud

Follow these steps:

- a. **Deploy the Repository**: Ensure that you have cloned or forked the repository to your GitHub account.
- b. **Set Required Secrets and Variables**: Before running any of the GitHub workflows provided in the repository, you must define at least the AWS\_DEV\_ACCOUNT\_ID variable (which should contain your AWS Account ID for your development environment). To do this:
  - Go to the **Settings** of your cloned or forked repository in GitHub.

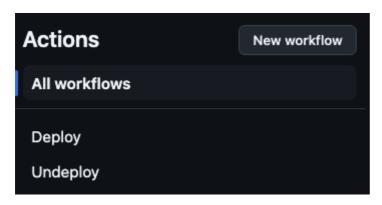
- Navigate to Secrets and Variables > Actions.
- Add the AWS\_DEV\_ACCOUNT\_ID and any other required variables or secrets.

#### c. Navigate to the Actions Page:

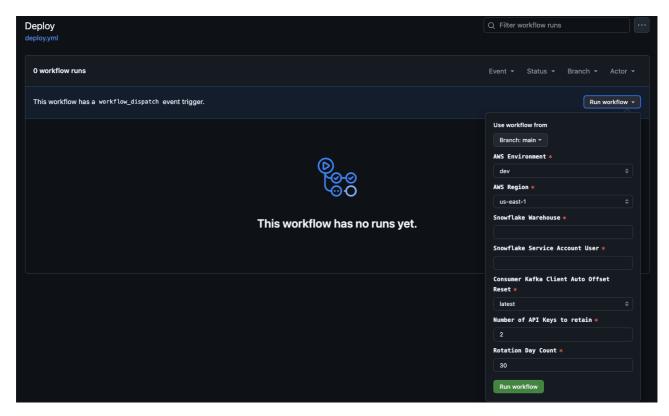
• From the cloned or forked repository on GitHub, click on the **Actions** tab.

#### d. Select and Run the Deploy Workflow:

• Find the **Deploy workflow** link on the left side of the Actions page and click on it.



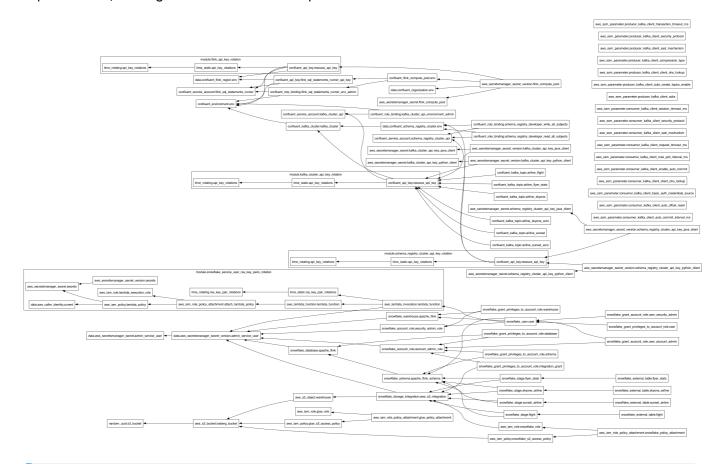
- On the **Deploy workflow** page, click the **Run workflow** button.
- A workflow dialog box will appear. Fill in the necessary details and click **Run workflow** to initiate the Terraform deployment process.



By following these steps, you will run the Terraform configuration directly from GitHub, leveraging GitHub Actions for automation and deployment.

#### 2.3 Visualizing the Terraform Configuration

Below is the Terraform visualization of the Terraform configuration. It shows the resources and their dependencies, making the infrastructure setup easier to understand.



To fully view the image, open it in another tab on your browser to zoom in.

When you update the Terraform Configuration, to update the Terraform visualization, use the terraform graph command with Graphviz to generate a visual representation of the resources and their dependencies. To do this, run the following command:

terraform graph | dot -Tpng > .blog/images/terraform-visualization.png

## 3.0 Hands-On Kickoff: Practical Examples for Rapid Learning

- 3.1 Flink Applications Powered by Java on a locally running Apache Flink Cluster in Docker Let's go!
- 3.2 Flink Applications Powered by Python on a locally running Apache Flink Cluster in Docker Let's go!
- 3.3 Flink Applications Powered by Python on Confluent Cloud for Apache Flink Let's qo!
- 4.0 Powering Analytics with Apache Iceberg tables in Snowflake

After running the Flink applications, you can now leverage the data stored in the Apache Iceberg tables in Snowflake for analytics. Click here to learn more about using **Non-Managed Apache Iceberg tables** in Snowflake. This blog post will help you understand the trade-offs, operational characteristics, and integration patterns of using Apache Iceberg tables in Snowflake. These terraform configurations:

- aws-iam-snowflake-setup creates the necessary Snowflake IAM roles and policies in AWS to get started.
- snowflake-non-managed-apache-iceberg-table-privileges-setup grants the necessary privileges to the Snowflake resources to get started.
- snowflake-non-managed-apache-iceberg-table-setup creates the necessary resources in Snowflake to get started.

#### 5.0 Resources

Apache Flink's Core is Dataflow Programming

What is Apache Flink? — Architecture

Apache Flink Use Cases

Building Apache Flink Applications in Java

J3's techStack Lexicon

Unlocking Schema Registry Access: Granting Confluent Service Account Permissions with Terraform

## 6.0 Important Note(s)

Known Issue(s)