

---

# WK 2 DE ZOOMCAMP KESTRA

---

KESTRA

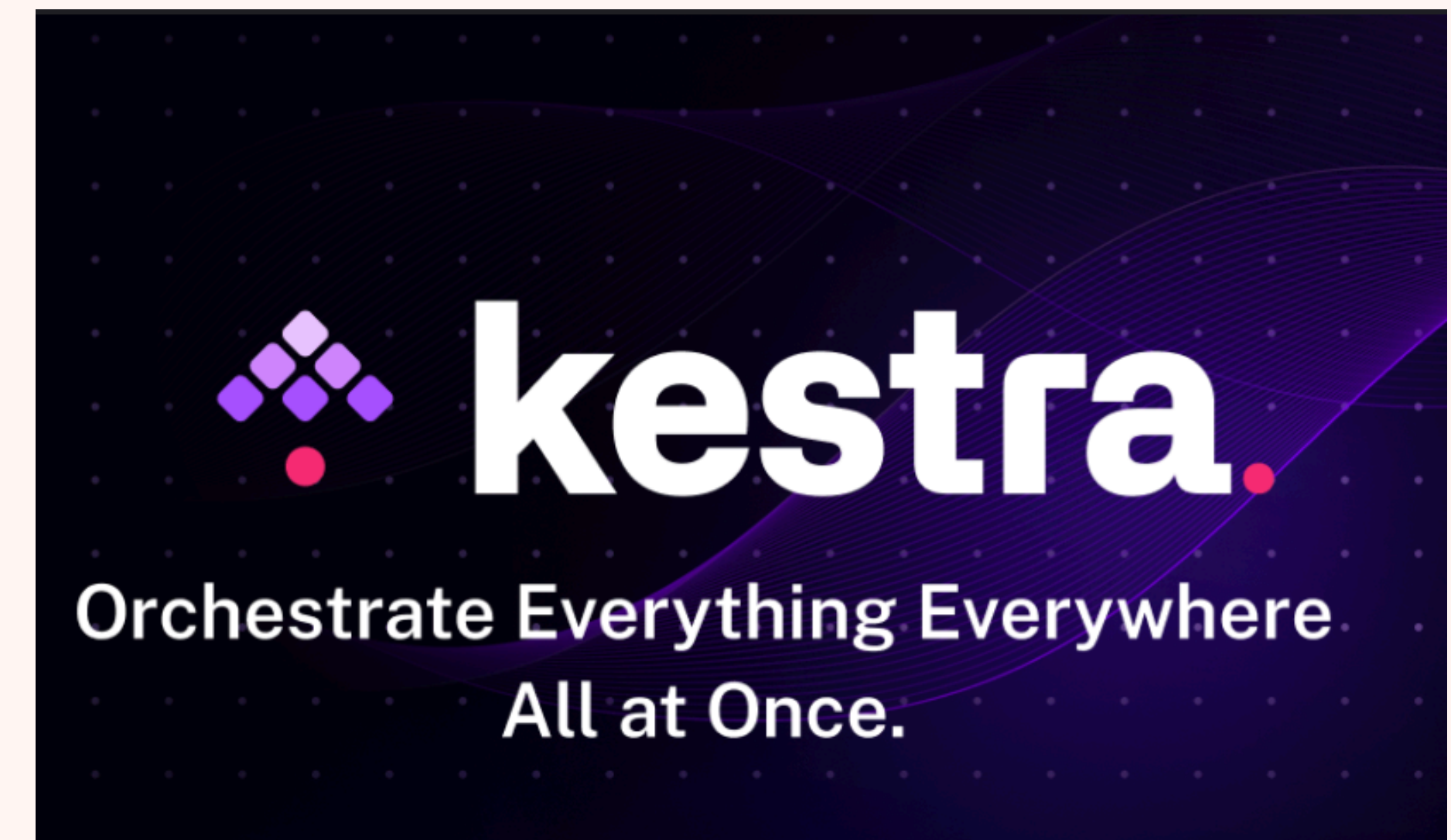
HABEEB

ZOOM

---

# INTRO

- **Kestra** is an open source data orchestration platform designed to automate complex workflows and pipelines
- It supports scheduling, dependency management, and backfill operations for seamless batch processing
- With native integrations to data sources like PostgreSQL, GCP BigQuery, and Cloud Storage, Kestra streamlines data movement and transformation
- Its extensible plugin architecture empowers users to customize and scale their data workflows effectively



Kestra

---

# RUNNING KESTRA USING DOCKER COMPOSE

- Kestra can be run or installed in multiple ways (see <https://kestra.io/docs/installation>)
- You can download the docker-compose file to run Kestra from the following link: <https://raw.githubusercontent.com/kestra-io/kestra/develop/docker-compose.yml>
- Once the docker-compose file is downloaded, run ‘docker compose up’ from your terminal. If successful, the Kestra page should be accessible on port 8080



Dashboard

Flows

Executions

Logs

Namespaces

Blueprints

Plugins

Administration

Settings

Flows

Jump to...

Ctrl/Cmd + K

Import

Source search

+ Create



scope

is one of

USER



Executions (per day)

99

Executions

70  
60  
50  
40  
30  
20  
10  
0

01/03

01/08

01/13

01/18

01/23

01/28

02/01

Date

Duration



Failed

Success

Duration  
30s  
25s  
20s  
15s  
10s  
5s  
0s

Total: 6

<input type="checkbox"/>	Id	Labels	Namespace	Last execution date	Last execution status	Execution statistics	Triggers
<input type="checkbox"/>	02_postgres_taxi_scheduled		zoomcamp	Sun, Feb 2, 2025 11:58 AM	SUCCESS		
<input type="checkbox"/>	04_gcp_kv		zoomcamp	Sun, Feb 2, 2025 12:40 PM	SUCCESS		
<input type="checkbox"/>	05_gcp_setup		zoomcamp	Sun, Feb 2, 2025 1:41 PM	SUCCESS		
<input type="checkbox"/>	06_gcp_taxi		zoomcamp	Sun, Feb 2, 2025 2:19 PM	SUCCESS		
<input type="checkbox"/>	06_gcp_taxi_scheduled		zoomcamp	Sun, Feb 2, 2025 3:59 PM	SUCCESS		
<input type="checkbox"/>	07_homework		zoomcamp	Sun, Feb 2, 2025 5:30 PM	SUCCESS		

25 per page

Total: 6



News



Docs



Help



Open an issue

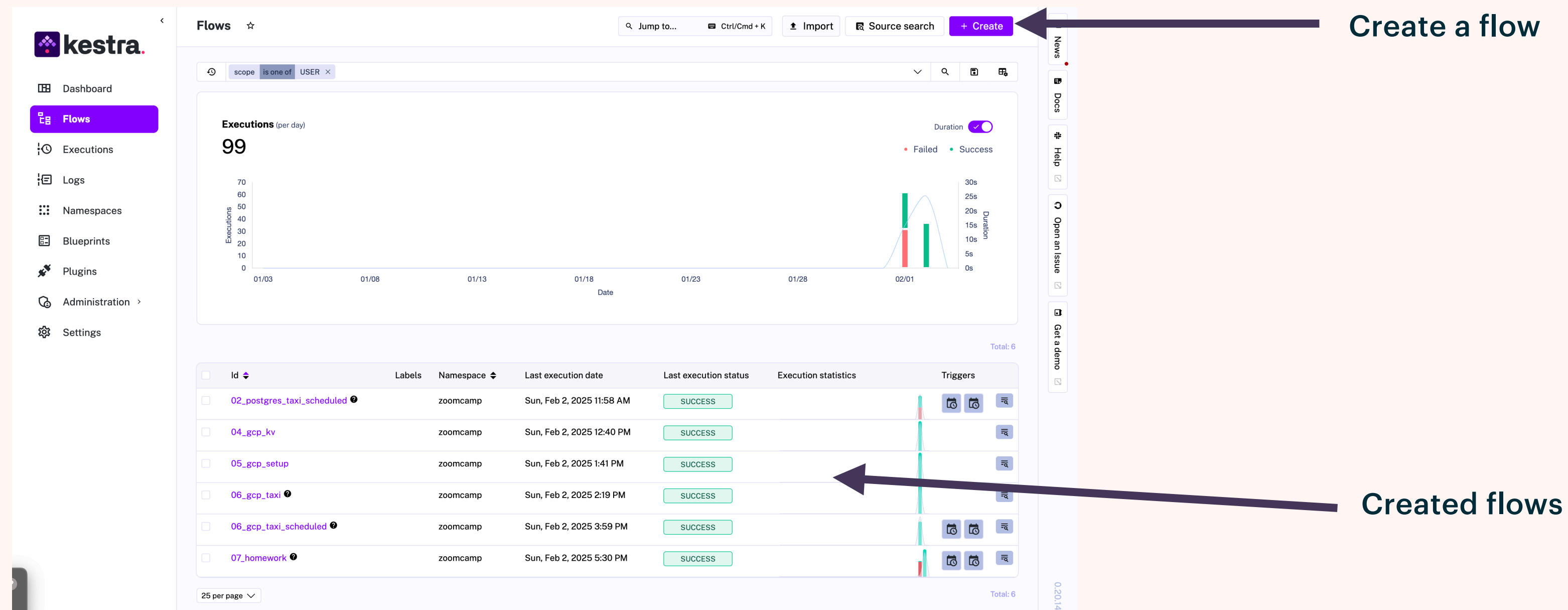


Get a demo

0.20.14

# Flows in Kestra UI

- Flows in Kestra orchestrate the sequence of tasks within a data workflow, defining clear steps for execution



Flows

---

# DATA INGESTION INTO POSTGRES USING KESTRA

---



# DATA INGESTION INTO POSTGRES USING KESTRA

- Data Ingestion with Kestra: Ingest NYC Taxi data into the Postgres database using Kestra
- Container Deployment: Run the Postgres container using Docker or Docker Compose
- Docker Compose Configuration: Since Kestra's Docker Compose already includes a Postgres instance for metadata, using a separate Compose file for external Postgres with PgAdmin can cause conflicts. I modified the Kestra Compose file to expose its Postgres on port 5433 and the external Postgres on port 5432, and created a network called "pg-network2" to connect Kestra, its metadata Postgres, and the external Postgres. (Note: I attempted to use "docker.host.internal" to avoid creating a new network, but it did not work in Codespaces.)

```
1 networks:
2   pg-network2:
3     external: true
4
5 services:
6   postgres:
7     image: postgres
8     container_name: postgres-db
9     environment:
10      - POSTGRES_USER=kestra
11      - POSTGRES_PASSWORD=k3str4
12      - POSTGRES_DB=postgres-zoomcamp
13     volumes:
14      - postgres-data:/var/lib/postgresql/data
15     ports:
16      - "5432:5432"
17     restart: unless-stopped
18     healthcheck:
19       test: ["CMD-SHELL", "pg_isready -U kestra -d postgres-zoomcamp"]
20       interval: 10s
21       timeout: 5s
22       retries: 5
23     networks:
24      - pg-network2
25
26 pgadmin:
27   image: dpage/pgadmin4
28   container_name: pg_admin
29   environment:
30     - PGADMIN_DEFAULT_EMAIL=admin@admin.com
31     - PGADMIN_DEFAULT_PASSWORD=root
32   ports:
33     - "8090:80"
34   depends_on:
35     postgres:
36       condition: service_healthy
37   restart: unless-stopped
38   networks:
39     - pg-network2
40
41 volumes:
42   postgres-data:
```

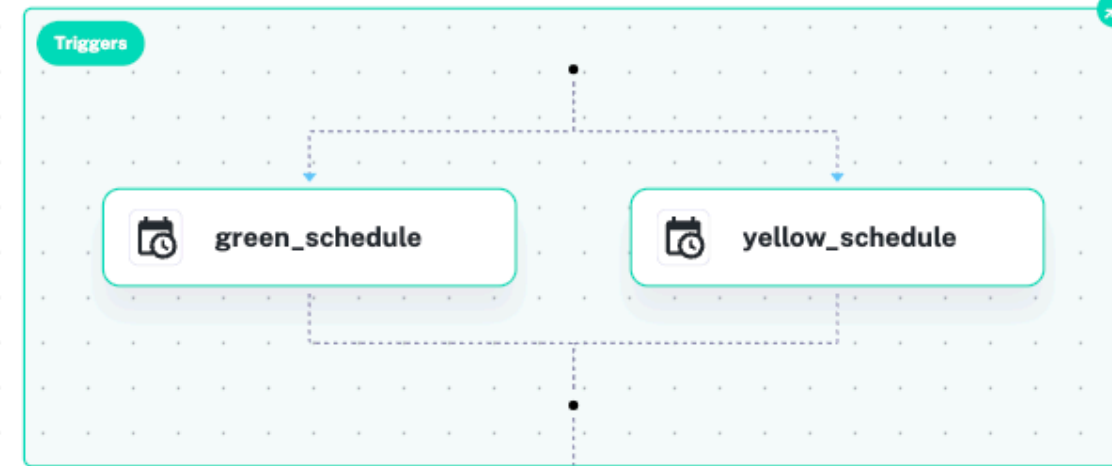
External postgres and Pgadmin

```
1 volumes:
2   postgres-data:
3     driver: local
4   kestra-data:
5     driver: local
6
7 networks:
8   pg-network2:
9     external: true
10
11 services:
12   postgres:
13     image: postgres
14     volumes:
15      - postgres-data:/var/lib/postgresql/data
16     ports:
17      - "5433:5432" # Running on a different port to avoid conflicts
18     environment:
19       POSTGRES_DB: kestra
20       POSTGRES_USER: kestra
21       POSTGRES_PASSWORD: k3str4
22     healthcheck:
23       test: ["CMD-SHELL", "pg_isready -d ${POSTGRES_DB} -U ${POSTGRES_USER}"]
24       interval: 30s
25       timeout: 10s
26       retries: 10
27     networks:
28      - pg-network2
29
```

My modified Kestra's docker compose

With this configuration, we can connect to the external Postgres from Kestra using

```
jdbc:postgresql://postgres-db:5432/postgres-zoomcamp
```



set\_label

extract



+





---

# WORKFLOW IN SLIDE 7

- We set a trigger for green\_taxi and yellow\_taxi, the yellow\_taxi is to run at 9:00AM on the first day of every month, and the green\_taxi scheduler is to run at 10:00AM on the first day of every month
  - We created 5 tasks:
    - set\_label
    - Extract
    - If\_yellow\_taxi
    - If\_green\_taxi
    - Purge
-

---

# FLOW INPUTS AND VARIABLES

- We have one input, taxi of type select which allows users to select an option between ‘yellow’ and ‘green’ in the Kestra UI
- We set 4 variables: file stores the expected filename. Staging\_table and table store the expected staging\_table name and table name respectively. Data is used to reference the output file from “extract” task.
- Inputs can be accessed using inputs.taxi
- Variables can be accessed using vars.file, vars.table and so on

```
▼ inputs:
  ▼ - id: taxi
      type: SELECT
      displayName: Select taxi type
      values: [yellow, green]
      defaults: yellow

▼ variables:
  file: "{{inputs.taxi}}_tripdata_{{trigger.date | date('yyyy-MM')}}.csv"
  staging_table: "public.{{inputs.taxi}}_tripdata_staging"
  table: "public.{{inputs.taxi}}_tripdata"
  data: "{{outputs.extract.outputFiles[inputs.taxi ~ '_tripdata_' ~ (trigger.date
```

Inputs and Variables

---

---

# TASK 1: SET LABEL

- Assigns labels to the execution for easy tracking in the UI
- Uses `render(vars.file)` to dynamically insert the selected file name
- Adds a label for the taxi type (yellow or green).

```
tasks:  
  - id: set_label  
    type: io.kestra.plugin.core.execution.Labels  
    labels:  
      file: "{{render(vars.file)}}"  
      taxi: "{{inputs.taxi}}"
```

set\_label task

---

---

# TASK 2: EXTRACT

- User selects a taxi type, year, and month in the Kestra UI
- The workflow constructs dynamic filenames and table names based on user selections.
- It downloads the selected dataset from GitHub and extracts and decompresses it as a CSV file using shell command, wget, and gunzip

```
id: extract
type: io.kestra.plugin.scripts.shell.Commands
outputFiles:
  - "*.csv"
taskRunner:
  type: io.kestra.plugin.core.runner.Process
commands:
  - wget -qO- https://github.com/DataTalksClub/nyc-  
file}}}.gz | gunzip > {{render(vars.file)}}
```

Extract

---

---

## TASK 3: IF\_YELLOW

- Checks if the input taxi type is 'yellow'
  - Executes tasks related to processing yellow taxi data
  - Creates main and staging tables if not exist
  - Delete rows in staging\_table to ensure that only the current execution data is processed
  - Loads data into staging table and updates with unique IDs
  - Merges data into the main yellow taxi table
-



---

# TASK 4: IF\_GREEN

- Checks if the input taxi type is 'green'
  - Executes tasks related to processing green taxi data
  - Creates main and staging tables if not exist
  - Delete rows in staging\_table to ensure that only the current execution data is processed
  - Loads data into staging table and updates with unique IDs
  - Merges data into the main green taxi table
-

---





# TASK 5: PURGE\_FILES

➤ This removes downloaded files from the Kestra storage

---

# EXECUTION

➤ To process data from 2019–2021, use backfill executions in the trigger tab (as it's now 2025)

	Worker Id	Next execution date	Backfill
			 
igin.core.trigger.Schedule		Sat, Mar 1, 2025 2:00 AM	 Backfill executions
igin.core.trigger.Schedule		Sat, Mar 1, 2025 3:00 AM	 Backfill executions

---

# TODO

- Modify the if\_yellow and if\_green conditions to process only months with existing data, or use ExecutionStatus/ExecutionOutputs conditions with the trigger scheduler

---

# DATA INGESTION INTO GOOGLE CLOUD

---




---


# WORKFLOW


- The Google Cloud workflow mirrors the Postgres workflow, but each month has a dynamically created staging\_table, eliminating the need to delete rows. However, Google Cloud variables must be set


# SETTING GOOGLE CLOUD VARIABLES


➤ The Google Cloud variables can be set manually from the KV Store in the namespace tab


 **kestra.**


 Dashboard


 Flows

 Executions

 Logs

 **Namespaces**

 Blueprints

 Plugins

Namespaces

zoomcamp ☆

Search

Overview

Editor

Flows

Executions

Dependencies

**KV Store**

Edit 🔒

Variables 🔒

Plugin defaults

Search

<input type="checkbox"/> Key ↕	Created date ↕	Updated date ↕
<input type="checkbox"/> GCP_CREDS	2025-02-02T19:40:15.803Z	2025-02-02T19:40:15.803Z
<input type="checkbox"/> GCP_BUCKET_NAME	2025-02-02T19:40:17.147Z	2025-02-02T19:40:17.147Z
<input type="checkbox"/> GCP_PROJECT_ID	2025-02-02T19:40:16.879Z	2025-02-02T19:40:16.879Z
<input type="checkbox"/> GCP_DATASET	2025-02-02T19:40:17.271Z	2025-02-02T19:40:17.271Z
<input type="checkbox"/> GCP_LOCATION	2025-02-02T19:40:17.007Z	2025-02-02T19:40:17.007Z

KV Store

---

# SETTING GOOGLE CLOUD VARIABLES

- The Google Cloud variables can also be set by creating a KV Set-based flow

```
- id: gcp_project_id
  type: io.kestra.plugin.core.kv.Set
  key: GCP_PROJECT_ID
  kvType: STRING
  value: de-zoomcamp-449719 #TODO
```