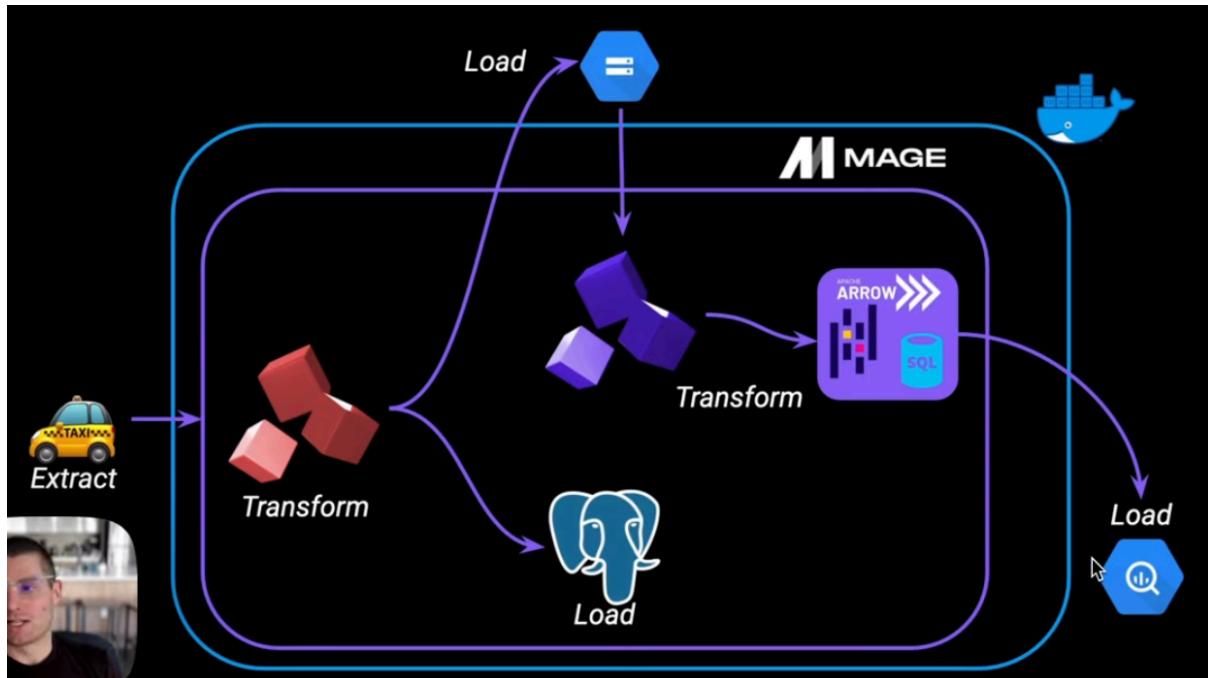


(32) DE Zoomcamp 2.2.1 - What is Orchestration? - YouTube

Final goal



What is Orchestration?

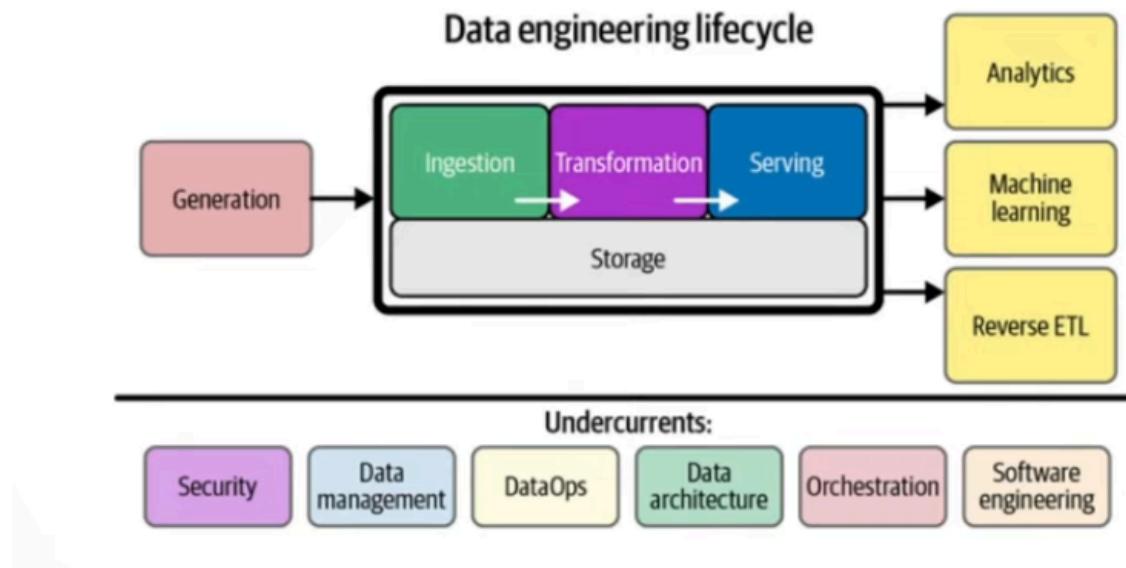
A large part of data engineering is **extracting, transforming, and loading** data between sources.

Orchestration is a process of dependency management, facilitated through **automation**.

The data orchestrator manages scheduling, triggering, monitoring, and even resource allocation.

Step = task = block

Worklog = dag = directed acyclic graphs



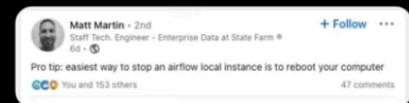
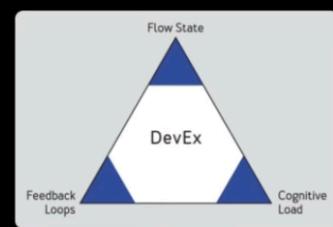
A good orchestrator handles...

- Workflow management
- Automation
- Error handling
- Recovery
- Monitoring, alerting
- Resource optimization
- Observability
- Debugging
- Compliance/Auditing

A good orchestrator prioritizes....

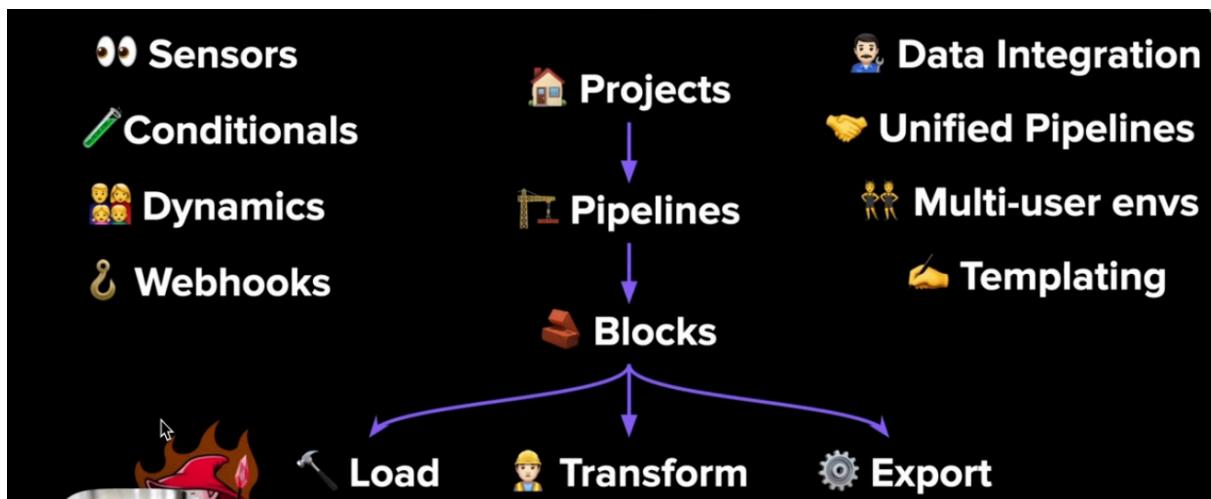
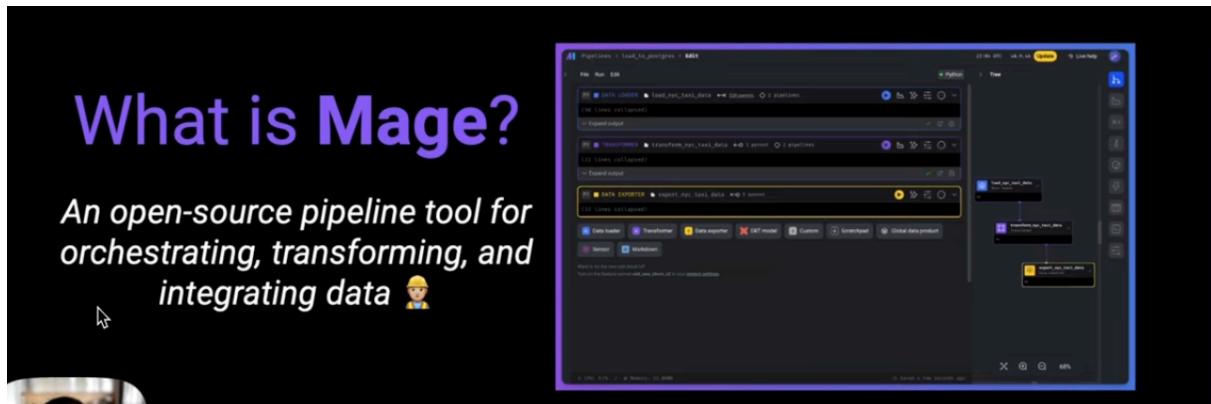
The developer experience

- Flow state 🌊
 - "I need to switch between 7 tools/services."
- Feedback Loops ⏪
 - "I spent 5 hours locally testing this DAG."
- Cognitive Load 🧠



How much do you need to know to do your job?

(32) DE Zoomcamp 2.2.2 - What is Mage? - YouTube



Mage accelerates pipeline development

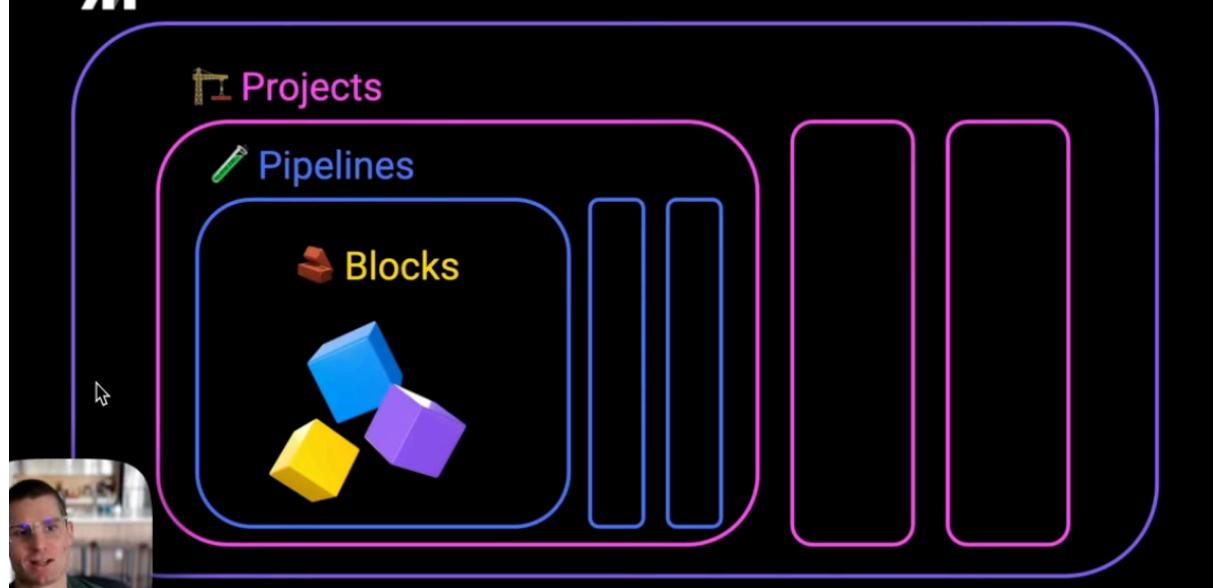
- Hybrid environment
 - Use our **GUI** for interactive development (or don't, I like VSCode)
 - Use **blocks** as testable, reusable pieces of code.
- Improved DevEx
 - Code and test in parallel.
 - Reduce your dependencies, switch tools less, be efficient.



Engineering best-practices built-in



- In-line testing and debugging
 - Familiar, notebook-style format
- Fully-featured observability
 - Transformation *in one place*: dbt models, streaming, & more.
- DRY principles
 - No more DAGs with duplicate functions and weird imports
 - DEaaS (sorry, I had to 😊)



Projects



- A project forms the basis for all the work you can do in Mage—you can think of it like a GitHub repo.
- It contains the code for all of your pipelines, blocks, and other assets.
- ↗ A Mage instance has one or more projects

Pipelines



- A pipeline is a workflow that executes some data operation—maybe extracting, transforming, and loading data from an API. They're also called DAGs on other platforms
- In Mage, pipelines can contain *Blocks* (written in SQL, Python, ↗ or R) and charts.
- Each pipeline is represented by a YAML file in the “pipelines” folder of your project.



Blocks



- A block is a file that can be executed independently or within a pipeline.
- Together, blocks form Directed Acyclic Graphs (DAGs), which we call pipelines.
- ↗ A block won't start running in a pipeline until all its upstream dependencies are met.



Blocks continued

- Blocks are reusable, atomic pieces of code that perform certain actions.
- Changing one block will change it everywhere it's used, but **don't worry**, it's easy to detach blocks to separate instances if necessary.
- Blocks can be used to perform a variety of actions, from simple data transformations to complex machine learning models.



blocks should return dataframes

Anatomy of a Block

The screenshot shows a Python code editor with a dark theme. A code block is highlighted, and a cursor is visible over the code. The code defines a class `DATA_LOADER` with a method `load_nyc_taxi_data`. The method uses `requests` to download a CSV file and `pd.read_csv` to parse it into a DataFrame. The DataFrame is then converted to a dictionary of types and returned. A test function `test_output` is also defined to check if the output is None.

```
1 import io
2 import pandas as pd
3 import requests
4 if 'data_loader' not in globals():
5     from mage_ai.data_preparation.decorators import data_loader
6 if 'test' not in globals():
7     from mage_ai.data_preparation.decorators import test
8
9 @data_loader
10 def load_nyc_taxi_data(*args, **kwargs):
11     url = 'https://github.com/DataTalksClub/nyc-tlc-data/releases/download/yellow/yellow_tripdata_2021-01.csv.gz'
12
13     taxi_dtypes = {-
14         # native date parsing
15         'tpep_pickup_datetime': 'tpep_dropoff_datetime': -
16     }
17
18     return pd.read_csv(-
19
20     @test
21     def test_output(output, *args) -> None:
22         """
23             Write code for testing the output of the block.
24
25             If output is not None, 'The output is undefined'.
26
27     
```

tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	RatecodeID	store_and_fwd_flag	PULocationID
1609461010000	1609461372000	1	2.1	1	N	142

Imports

Decorator

Function*

Assertion

*returns df

(32) DE Zoomcamp 2.2.2 - Configure Mage - YouTube

[mage-ai/mage-zoomcamp: This repository will contain all of the resources for the Mage component of the Data Engineering Zoomcamp:](#)

<https://github.com/DataTalksClub/data-engineering-zoomcamp/tree/main>

1. git clone <https://github.com/mage-ai/mage-zoomcamp.git> mage-zoomcamp
2. cd the folder
3. Create .env from dev.env
4. docker compose build
5. docker pull mageai/mageai:latest (optional)
6. docker compose up
7. <http://localhost:6789/>

DE Zoomcamp 2.2.2 - A Simple Pipeline (youtube.com)

If two blocks are connected DFs will be passed between them. The result of one is the input of the other

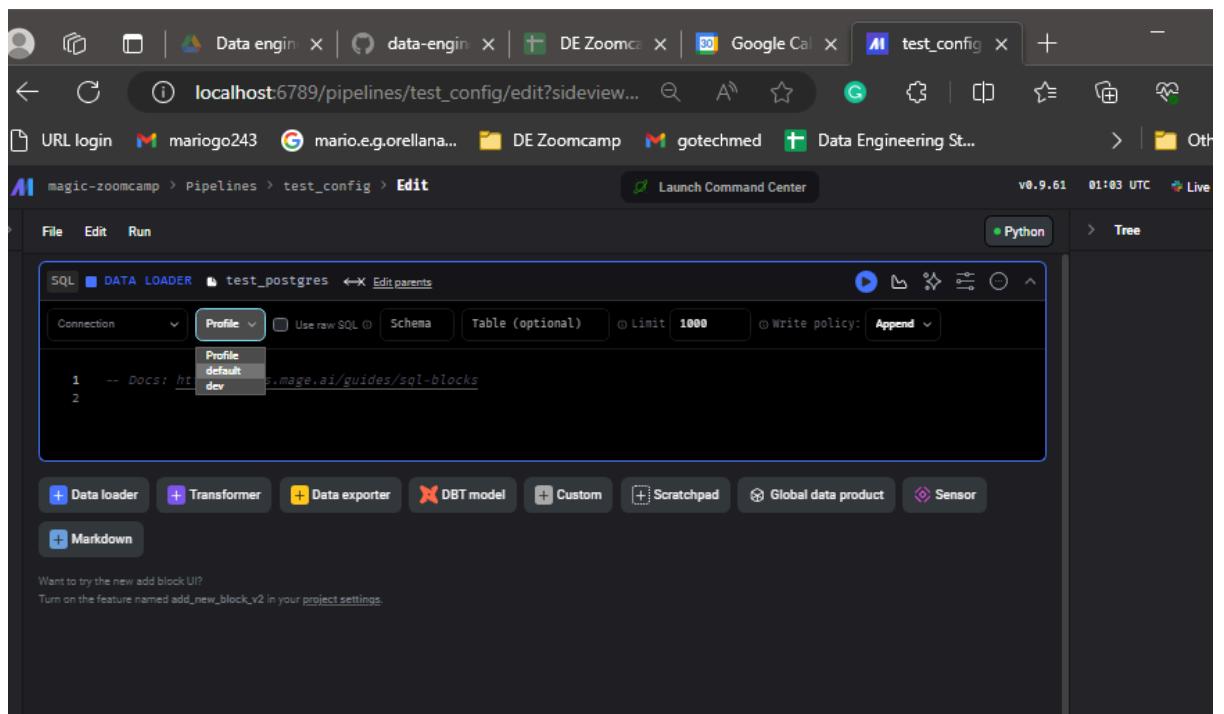
[Pipelines | Mage](#)

DE Zoomcamp 2.2.3 - Configuring Postgres (youtube.com)

1. Setup io_config.yaml is where we manage the connection, we can define a new profile
2. Let's set up a new profile called dev, let's set up environment variables, the way to do this in mage is with jinja notation "{{ {} }}"

```
o      POSTGRES_DBNAME: "{{ env_var('POSTGRES_DBNAME') }}"  
o  dev:  
o      # PostgreSQL  
o      POSTGRES_CONNECT_TIMEOUT: 10  
o      POSTGRES_DBNAME: "{{ env_var('POSTGRES_DBNAME') }}"  
o      POSTGRES_SCHEMA: "{{ env_var('POSTGRES_SCHEMA') }}"  
o      POSTGRES_USER: "{{ env_var('POSTGRES_USER') }}"  
o      POSTGRES_PASSWORD: "{{ env_var('POSTGRES_PASSWORD') }}"  
o      POSTGRES_HOST: "{{ env_var('POSTGRES_HOST') }}"  
o      POSTGRES_PORT: "{{ env_var('POSTGRES_PORT') }}"  
o
```

3. Create new pipeline and change name in edit pipeline settings
4. Add a new sql data loader
5. define profile



6. RUN SELECT 1; to test the connection

DE Zoomcamp 2.2.3 - ETL: API to Postgres (youtube.com)

[data-engineering-zoomcamp/Module 2/my_image_files at master · MarioOrellana58/data-engineering-zoomcamp \(github.com\)](https://github.com/MarioOrellana58/data-engineering-zoomcamp)

1. Create a new pipeline
2. Add a new data loader as python api
3. Add a new transformer python generic
4. Add a new data exporter python postgres
5. Add a new data loader for sql postgres to inspect the data, select postgres connection and mark raw sql checkbox

DE Zoomcamp 2.2.4 - Configuring GCP (youtube.com)

1. create a new google cloud storage bucket
2. Go to IAM > service accounts
3. Go to the new account and create a new key, paste the json into mage directory
4. Make sure that docker-compose.yml has this, this will make that all the files of my localhost directory will be mapped into that directory inside docker

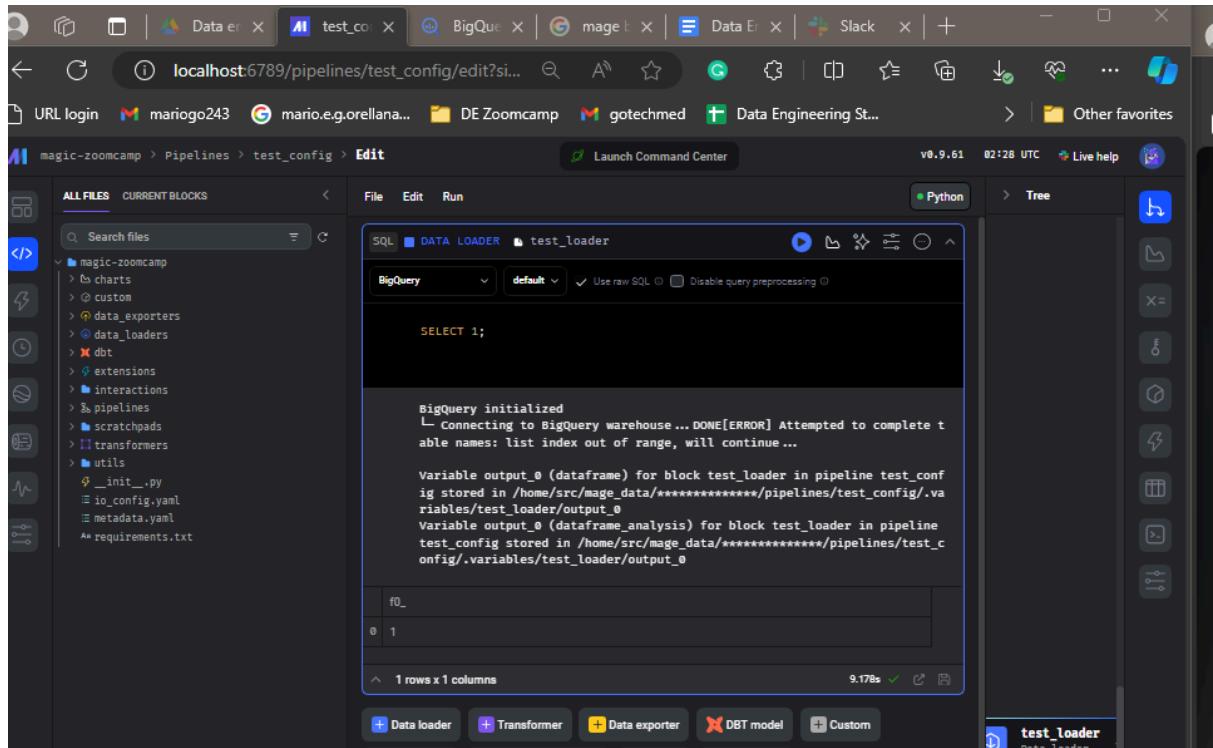
```
volumes:  
  - ./home/src/
```

5. Open mage, go to files then to io_config.yaml
 - a. Here you will find two ways to authenticate with Google, one is to paste all the payload from the key into the file, delete this one

```
b. GOOGLE_SERVICE_ACC_KEY:  
c.   type: service_account  
d.   project_id: project-id  
e.   private_key_id: key-id  
f.   private_key: "-----BEGIN PRIVATE  
KEY-----\nyour_private_key\n-----END_PRIVATE_KEY"  
g.   client_email: your_service_account_email  
h.   auth_uri: "https://accounts.google.com/o/oauth2/auth"  
i.   token_uri: "https://accounts.google.com/o/oauth2/token"  
j.   auth_provider_x509_cert_url:  
    "https://www.googleapis.com/oauth2/v1/certs"  
k.   client_x509_cert_url:  
    "https://www.googleapis.com/robot/v1/metadata/x509/your_serv  
ice_account_email"
```

- I. Use GOOGLE_SERVICE_ACC_KEY_FILEPATH instead
 - i. Go to terminal, run ls -la to check the files
 - ii. Copy the key name
 - iii. Update io config
6. Save the file
7. Use the test_config pipeline, change the connection to bigquery, then default profile, running it threw indexerror, here's the solution [Data Engineering Zoomcamp FAQ -](#)

Documentos de Google



The screenshot shows the Mage UI interface for pipeline configuration. On the left, there's a sidebar with a file tree. The main area is titled 'Edit' and contains a 'DATA LOADER' block named 'test_loader'. The code editor shows a simple SQL query: 'SELECT 1;'. The output window displays logs related to connecting to BigQuery and initializing variables. The bottom status bar shows a execution time of 9.178s.

8. Go to example_pipeline go to the last step, execute with the previous steps
9. Upload titanic_clean.csv to gcp bucket
10. go to test config pipeline, delete the bigquery block
11. Create a new data loader, python, google cloud storage
12. Update bucket name and object key (filename)
13. Run it

DE Zoomcamp 2.2.4 - ETL: API to GCS (youtube.com)

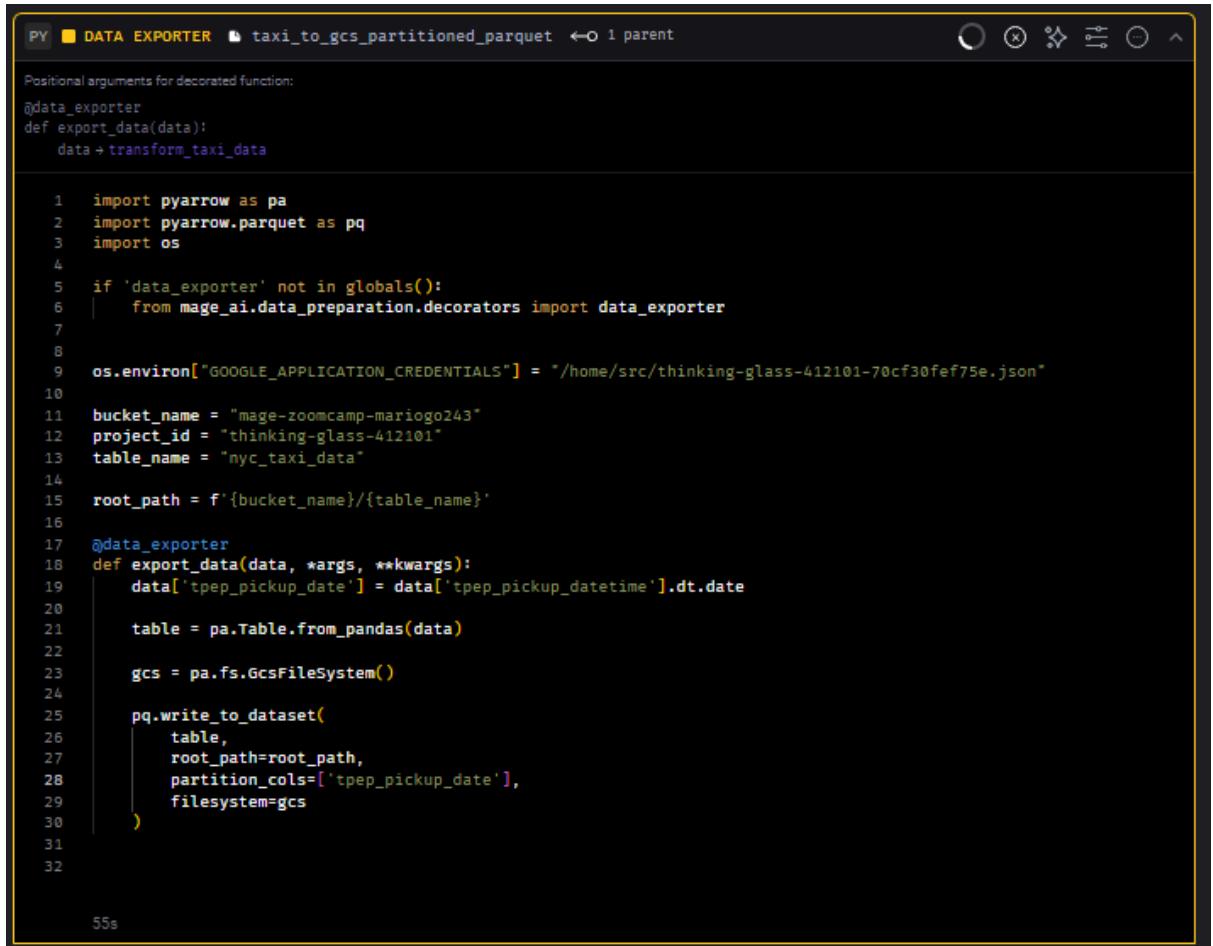
[data-engineering-zoomcamp/Module 2/my_mage_files/3. Load Data GCS Partitioned.py at master · MarioOrellana58/data-engineering-zoomcamp \(github.com\)](https://github.com/MarioOrellana58/data-engineering-zoomcamp/blob/master/Module%202/my_mage_files/3.%20Load%20Data%20GCS%20Partitioned.py)

Write to a single parquet file

1. New standard pipeline
2. Expand data_loaders folder in the left, then drag and drop load_api_data
3. Expand transformers, drag and drop transform taxi data
4. In the blocks view connect the two blocks dragging the lower dot of data loader to the top dot of transformer
5. Create a new data exporter python gcs
6. Run all the blocks

Write to a partitioned parquet file

1. Add a new data exporter python generic
2. Create a new date column
3. Partition data by date



```
PY DATA EXPORTER taxi_to_gcs_partitioned_parquet ← 1 parent
Positional arguments for decorated function:
@data_exporter
def export_data(data):
    data = transform_taxi_data

1  import pyarrow as pa
2  import pyarrow.parquet as pq
3  import os
4
5  if 'data_exporter' not in globals():
6      from mage_ai.data_preparation.decorators import data_exporter
7
8
9  os.environ["GOOGLE_APPLICATION_CREDENTIALS"] = "/home/src/thinking-glass-412101-70cf30fef75e.json"
10
11 bucket_name = "mage-zoomcamp-mariogo243"
12 project_id = "thinking-glass-412101"
13 table_name = "nyc_taxi_data"
14
15 root_path = f'{bucket_name}/{table_name}'
16
17 @data_exporter
18 def export_data(data, *args, **kwargs):
19     data['tpep_pickup_date'] = data['tpep_pickup_datetime'].dt.date
20
21     table = pa.Table.from_pandas(data)
22
23     gcs = pa.fs.GcsFileSystem()
24
25     pq.write_to_dataset(
26         table,
27         root_path=root_path,
28         partition_cols=['tpep_pickup_date'],
29         filesystem=gcs
30     )
31
32
```

55s

(35) DE Zoomcamp 2.2.5 - ETL: GCS to BigQuery - YouTube

1. New pipeline
2. Drag the loader python gcs
3. Transform the DF
4. Export, if this sql export fails try python export

The screenshot shows a Data Exporter interface with the following configuration:

- Target: BigQuery
- Dataset: defa
- Table: ny_taxi
- Schema: yellow_cab_data
- Limit: 1000
- Write policy: Append

The interface displays the following SQL query:

```
1 SELECT * FROM {{ df_1 }}
```

```
dfs = []
for object_key in object_keys:
    df =
pq.read_table(f'gs://{bucket_name}/{object_key}').to_pandas()
    dfs.append(df)

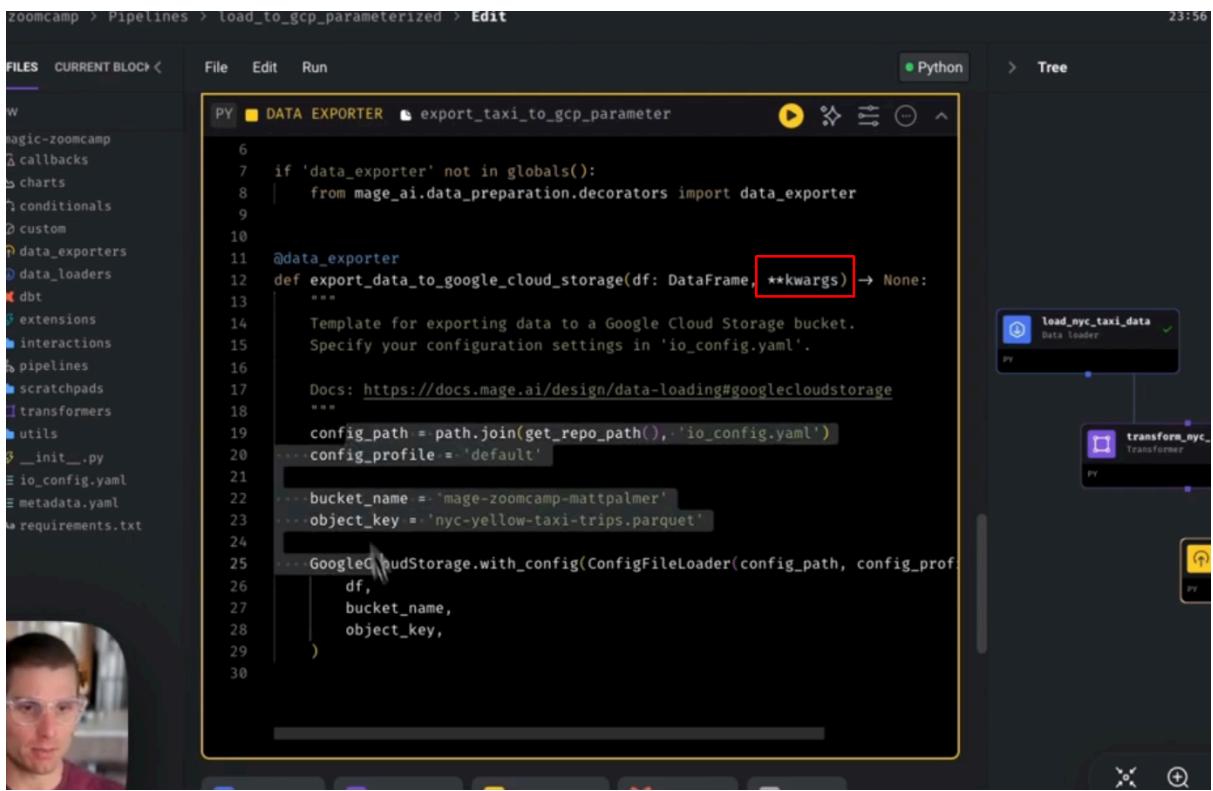
result_df = pd.concat(dfs, ignore_index=True)
```

DE Zoomcamp 2.2.6 - Parameterized Execution (youtube.com)

Mage has runtime variables, global variables, many variables... Look into the documentation.

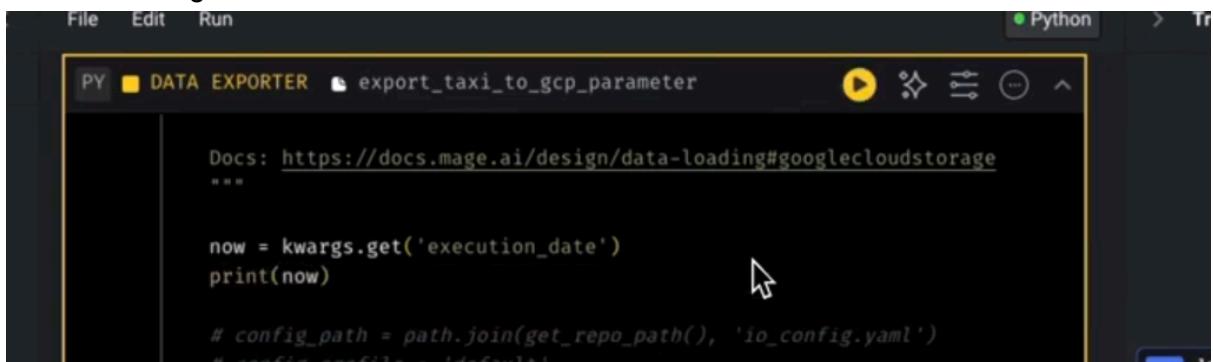
This video show how to create a different file for each day

1. Clone the previous pipeline
2. If you're going to edit a block create a new one and copy the code, if you edit an existing block it will affect all the pipelines where it's in use
3. This saves the runtime variables



```
zoomcamp > Pipelines > load_to_gcp_parameterized > Edit
FILEs CURRENT BLOCK < File Edit Run Python > Tree
PY DATA EXPORTER export_taxi_to_gcp_parameter
6
7 if 'data_exporter' not in globals():
8     from mage_ai.data_preparation.decorators import data_exporter
9
10
11 @data_exporter
12 def export_data_to_google_cloud_storage(df: DataFrame, **kwargs) -> None:
13     """
14         Template for exporting data to a Google Cloud Storage bucket.
15         Specify your configuration settings in 'io_config.yaml'.
16
17         Docs: https://docs.mage.ai/design/data-loading#googlecloudstorage
18     """
19     config_path = path.join(get_repo_path(), 'io_config.yaml')
20     config_profile = 'default'
21
22     bucket_name = 'mage-zoomcamp-mattpalmer'
23     object_key = 'nyc-yellow-taxi-trips.parquet'
24
25     GoogleCloudStorage.with_config(ConfigFileLoader(config_path, config_profile),
26         df,
27         bucket_name,
28         object_key,
29     )
30
```

4. This is how to get it

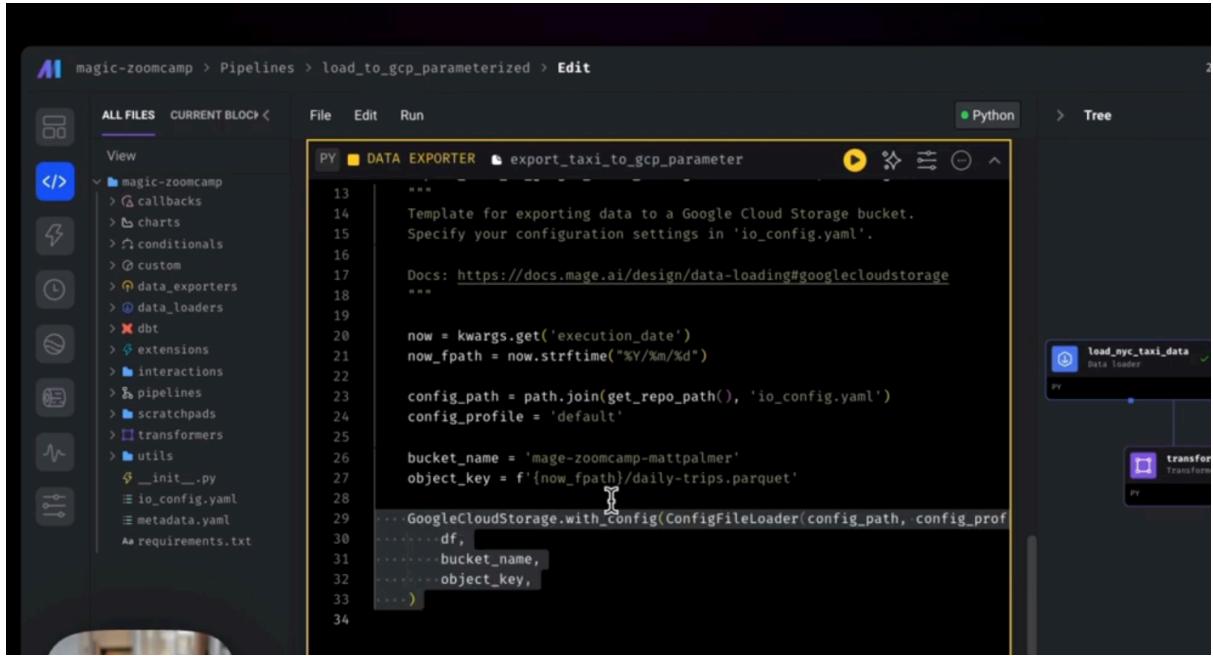


```
File Edit Run Python > Tree
PY DATA EXPORTER export_taxi_to_gcp_parameter
Docs: https://docs.mage.ai/design/data-loading#googlecloudstorage
"""

now = kwargs.get('execution_date')
print(now)

# config_path = path.join(get_repo_path(), 'io_config.yaml')
```

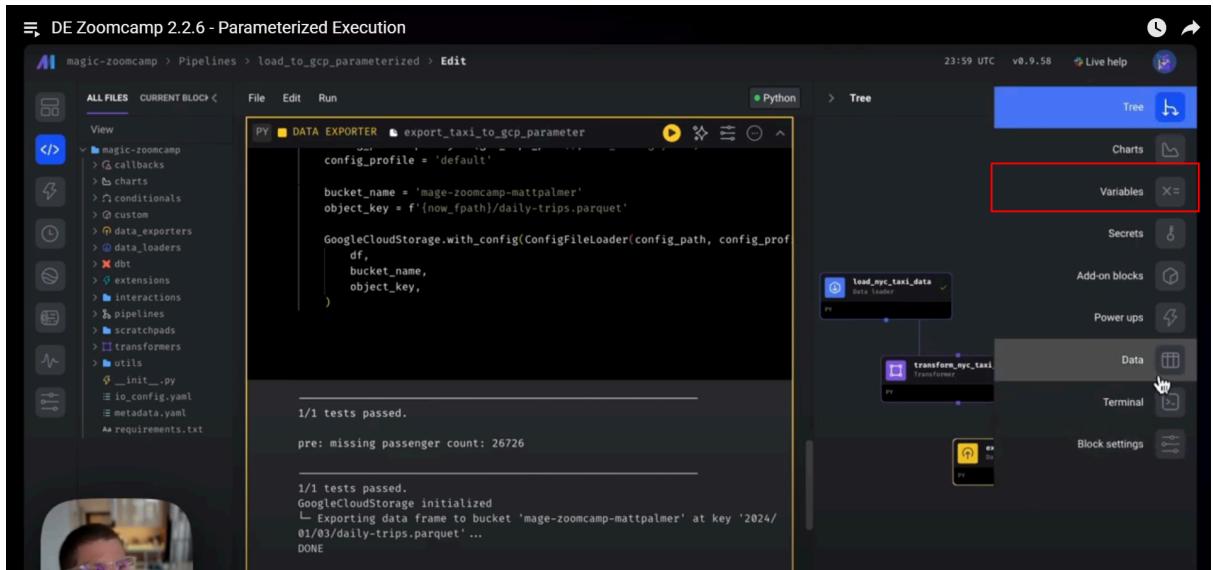
5. Create a file path for each file



The screenshot shows the Mage AI IDE interface. The left sidebar lists project files: callbacks, charts, conditionals, custom, data_exporters, data_loaders, dbt, extensions, interactions, pipelines, scratchpads, transformers, utils, __init__.py, io_config.yaml, metadata.yaml, and requirements.txt. The main editor window displays Python code for a 'DATA EXPORTER' named 'export_taxi_to_gcp_parameter'. The code uses the GoogleCloudStorage.with_config() method to write a DataFrame to a Google Cloud Storage bucket. A tooltip for 'GoogleCloudStorage.with_config()' is visible, pointing to the line 'GoogleCloudStorage.with_config(ConfigFileLoader(config_path, config_profile))'. The right sidebar shows a pipeline diagram with nodes for 'load_nyc_taxi_data' (Data Loader) and 'transform_nyc_taxi' (Transformer).

```
13     """
14     Template for exporting data to a Google Cloud Storage bucket.
15     Specify your configuration settings in 'io_config.yaml'.
16
17     Docs: https://docs.mage.ai/design/data-loading#googlecloudstorage
18     """
19
20     now = kwargs.get('execution_date')
21     now_fpath = now.strftime("%Y/%m/%d")
22
23     config_path = path.join(get_repo_path(), 'io_config.yaml')
24     config_profile = 'default'
25
26     bucket_name = 'mage-zoomcamp-mattpalmer'
27     object_key = f'{now_fpath}/daily-trips.parquet'
28
29     GoogleCloudStorage.with_config(ConfigFileLoader(config_path, config_profile))
30         df,
31         bucket_name,
32         object_key,
33     )
34
```

6. Here's how to set global variables



The screenshot shows the Mage AI IDE interface. The left sidebar lists project files: callbacks, charts, conditionals, custom, data_exporters, data_loaders, dbt, extensions, interactions, pipelines, scratchpads, transformers, utils, __init__.py, io_config.yaml, metadata.yaml, and requirements.txt. The main editor window displays Python code for a 'DATA EXPORTER' named 'export_taxi_to_gcp_parameter'. The code uses the GoogleCloudStorage.with_config() method to write a DataFrame to a Google Cloud Storage bucket. The right sidebar includes sections for 'Variables' (highlighted with a red box), 'Secrets', 'Add-on blocks', 'Power ups', 'Data', 'Terminal', and 'Block settings'. The pipeline diagram on the right shows nodes for 'load_nyc_taxi_data' (Data Loader) and 'transform_nyc_taxi' (Transformer). The terminal output at the bottom shows test results and the execution of the export command.

```
1/1 tests passed.

pre: missing passenger count: 26726

1/1 tests passed.
GoogleCloudStorage initialized
└ Exporting data frame to bucket 'mage-zoomcamp-mattpalmer' at key '2024/01/03/daily-trips.parquet' ...
DONE
```

The screenshot shows the DataCamp Platform interface for a Python pipeline. At the top, it displays "00:00 UTC v0.9.58" and a "Live help" button. On the left, there's a sidebar with a "Python" tab selected, followed by a navigation arrow and the word "Variables". A large text area in the center says "Press Enter or Return to save changes." Below this, a note states: "Global variables will be passed into all non-scratchpad blocks as keyword arguments (Python), interpolated variables (SQL), or vector elements (R). To load a global variable, use the following syntax:". It then provides examples for Python, SQL, and R. On the right side, there's a vertical toolbar with various icons for file operations like copy, paste, delete, and search. At the bottom, there's a "Trigger Runtime Variables" section with a note about default runtime variables based on trigger type, and a "Schedule" section showing an "event" trigger.

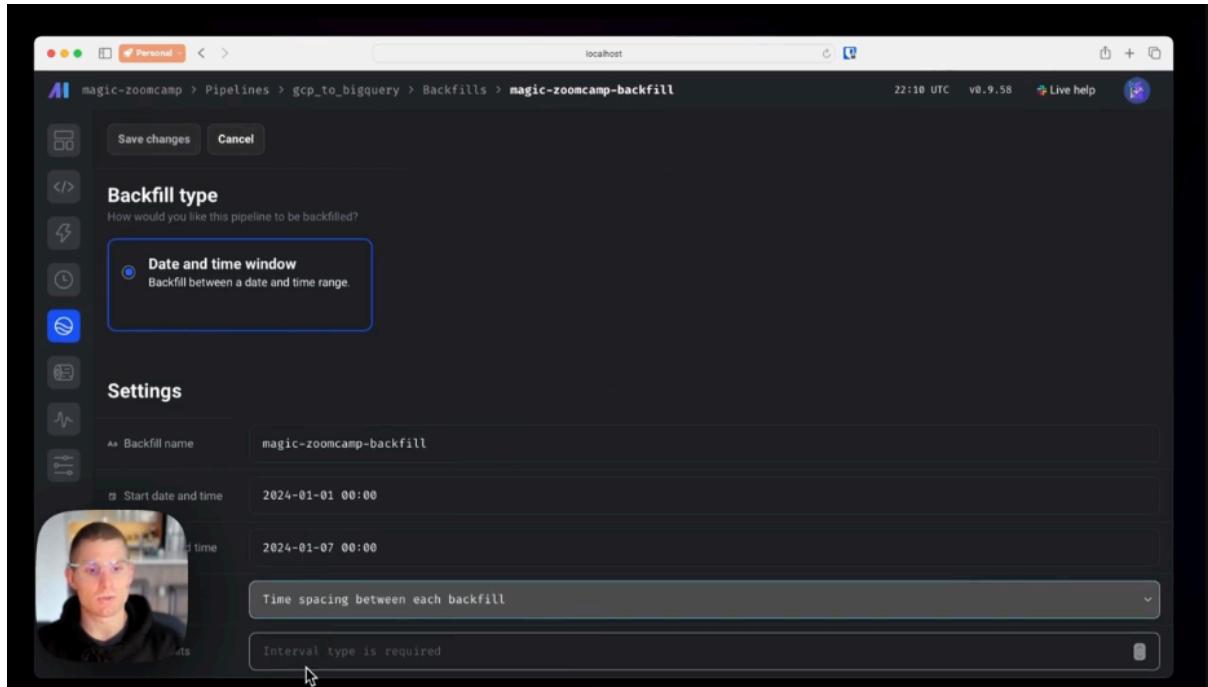
- 7.
8. We can also add runtime variables in the trigger

The screenshot shows the "Triggers" settings for a pipeline named "victorious destiny". In the "Settings" section, the "Trigger name" is set to "victorious destiny" and the "Trigger description" is "Description". Under "Frequency", the "Daily" option is selected. There's a note about enabling landing time. In the "Run settings" section, there are options for "Configure trigger SLA", "Keep running pipeline even if blocks fail", and "Skip run if previous run still in progress". The "Runtime variables" section indicates that no runtime variables are currently defined for this pipeline.

DE Zoomcamp 2.2.6 - Backfills (youtube.com)

Backfill is to restore missing data for example

1. Click any pipeline
2. Go to backfills tab
3. Select the data frame of the backfill (last day is inclusive)



Deploy MAGE to Google Cloud

Videos:

- [!\[\]\(c35b00deb15b9877bd4bf11a0ea5ff77_img.jpg\) DE Zoomcamp 2.2.7 - Deployment Prerequisites](#)
- [!\[\]\(ce81ecc9a883bc9780490b6bbaa8cd9f_img.jpg\) DE Zoomcamp 2.2.7 - Google Cloud Permissions](#)
- [!\[\]\(a9d31affde24f6c67af062c826defab2_img.jpg\) DE Zoomcamp 2.2.7 - Deploying to Google Cloud Part 1](#)
- [!\[\]\(9b2dfd55c9797b68203e02f02d2b522a_img.jpg\) DE Zoomcamp 2.2.7 - Deploying to Google Cloud Part 2](#)

1. install terraform
2. install gcloud cli [Install the gcloud CLI | Google Cloud CLI Documentation](#)

3. Change gcloud permissions

Edit access to "terraform-demo"

Principal [?](#)
mage-zoomcamp-mariogo243@thinking-glass-412101.iam.gserviceaccount.com

Project terraform-demo

Summary of changes
Role removed
Owner
Roles added
Artifact Registry Reader
Artifact Registry Writer
Cloud Run Developer
Cloud SQL Admin
Service Account Token C

Assign roles

Roles are composed of sets of permissions and determine what the principal can do with this resource. [Learn more](#)

Role Artifact Registry Reader **IAM condition (optional)** [?](#) [+ ADD IAM CONDITION](#) [TEST CHANGES](#) [?](#)

Access to read repository items.

Role Artifact Registry Writer **IAM condition (optional)** [?](#) [+ ADD IAM CONDITION](#)

Access to read and write repository items.

Role Cloud Run Developer **IAM condition (optional)** [?](#) [+ ADD IAM CONDITION](#)

Read and write access to all Cloud Run resources.

Role Cloud SQL Admin **IAM condition (optional)** [?](#) [+ ADD IAM CONDITION](#)

Full control of Cloud SQL resources.

Role Service Account Token Creator **IAM condition (optional)** [?](#) [+ ADD IAM CONDITION](#)

Impersonate service accounts (create OAuth2 access tokens, sign blobs or JWTs, etc.).

[+ ADD ANOTHER ROLE](#)

SAVE **TEST CHANGES** [?](#) **CANCEL**

4. Git clone [mage-ai/mage-ai-terraform-templates: Terraform templates for deploying mage-ai to AWS, GCP and Azure \(github.com\)](#)
5. cd into gcp folder
6. gcloud auth application-default login
7. Enable google cloud firestore api
8. terraform init, terraform plan and terraform apply
9. Go to google cloud run, allow all ips and access the new service URL

d