# Data lineage via Apache Airflow

# Short bio Matteo

## Matteo Senardi

\* 2020-ongoing Head of Data in Docsity

\* 2019-2020 Senior Data Engineer in Mediaset

\* 2016-2019 Innovation Engineer in Webranking

\*  email: matteo.s@docsity.com

SCAN ME
https://github.com/pualien

# Short bio Irene



## Irene Soligno

* 2020-ongoing Data Scientist in Docsity

* 2019-2020 Data Scientist in Enginium

* 2015-2019 PhD and research fellow in POLITO
(environmental engineering)

*  email: irene.s@docsity.com

* ORCID: https://orcid.org/0000-0001-9884-5316

# docsity

Founded in Italy in 2011 - Now present in 70 countries, with offices in Turin, Rome and Sao Paulo (Brazil).

Docsity is an EdTech company operating on an international scale that distributes and produces digital education content to help university students prepare for their exams.

# docsity

Docsity also collaborates with over 150 universities and business schools from around the world, by promoting their educational programs through the Docsity student community (Business To Business - B2B).

Docsity has today 500K new registered users each month and 12M visitors per month.

Full-service approach: implement and optimize the whole student funnel.

# WHY AIRFLOW?
## part I

# Let's imagine - a very simple use case

## 01

**PROBLEM**

Decide if batch job should run today and how long to backfill

## 02

**PROCESS**

Data sources,
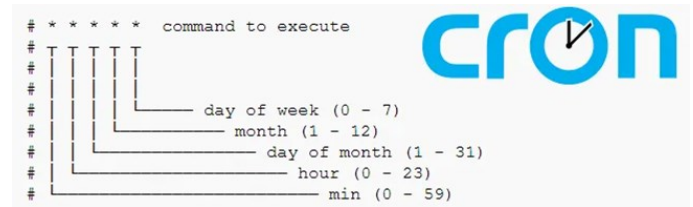that need to be passed through data lake
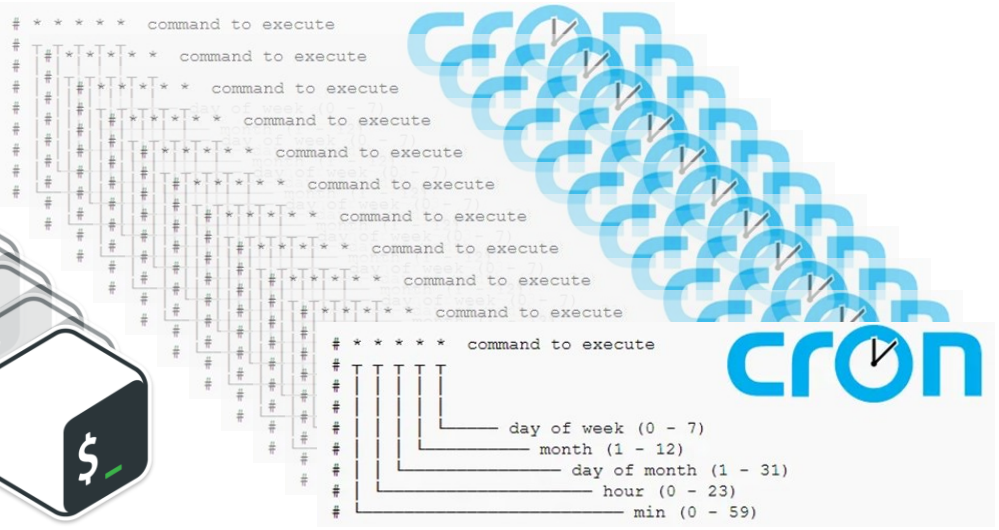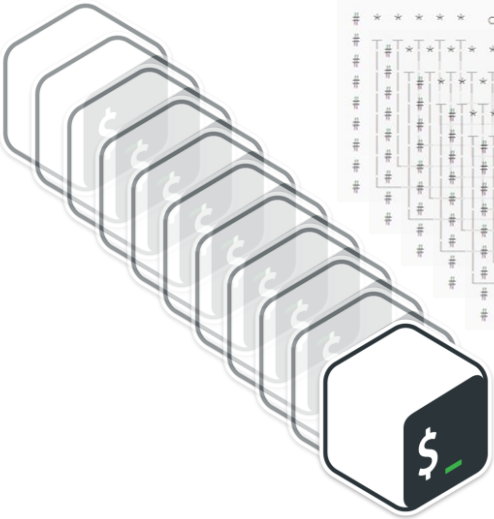
## 03

**TARGET**

When all files are in place, submit the pre-defined job.

# Let's imagine - a very simple use case



## Scripting + Cron could do!

# Let's imagine - a very simple use case

What if **hundreds** of workflows to be managed?

# Let's imagine - a very ~~simple~~ production use case

▶ Which features **are required** to manage production workflows?

**MANAGE**

Manage scripts and crontab for hundred of workflows

**ENVIRONMENT**

Jobs could have different requirements and environments

**MONITOR**

Track performance and completion of each step

**EXECUTION**

Consider performance in parallel

**CONNECTION**

Operational DBs, APIs, Cloud Services, ssh tunneling, all with their own configuration

**RETRY**

Re-run a specific step in case of failures or whatever

Apache Airflow

born in 2014 from airbnb

The technology to **automate** & **orchestrate** data pipelines

~1700 GITHUB CONTRIBUTORS

PYTHON BASED

MONITORING

DAGS AS CODE

FRIENDLY
INTERFACE

SECURITY BY DESIGN

DISTRIBUTED EXECUTION

DATA
LINEAGE

AIRFLOW DEPLOYMENT

# Airflow Docsity solution selection

**MANAGED**

Google Cloud Composer

Amazon Managed Workflows
for Apache Airflow

ASTRONOMER

**CUSTOM CONTAINER**

docker

Google Compute Engine

AWS ECS

# Airflow Docsity solution selection

## MANAGED

* KEYS IN HAND

* MAINTENANCE BY PROVIDER

* MANAGED SCALABILITY

* PREDEFINED AIRFLOW VERSIONS

## CUSTOM CONTAINER

* REPRODUCIBILITY

* TESTING

* ISOLATION

* PORTABILITY

* PRICING

* CUSTOM AIRFLOW VERSIONS

# Airflow Docsity solution

## MANAGED

* KEYS IN HAND

* MAINTENANCE BY PROVIDER

* MANAGED SCALABILITY

* PREDEFINED AIRFLOW VERSIONS

## CUSTOM CONTAINER

* REPRODUCIBILITY

* TESTING

* ISOLATION

* PORTABILITY

* PRICING

* CUSTOM AIRFLOW VERSIONS

# Dockerfile

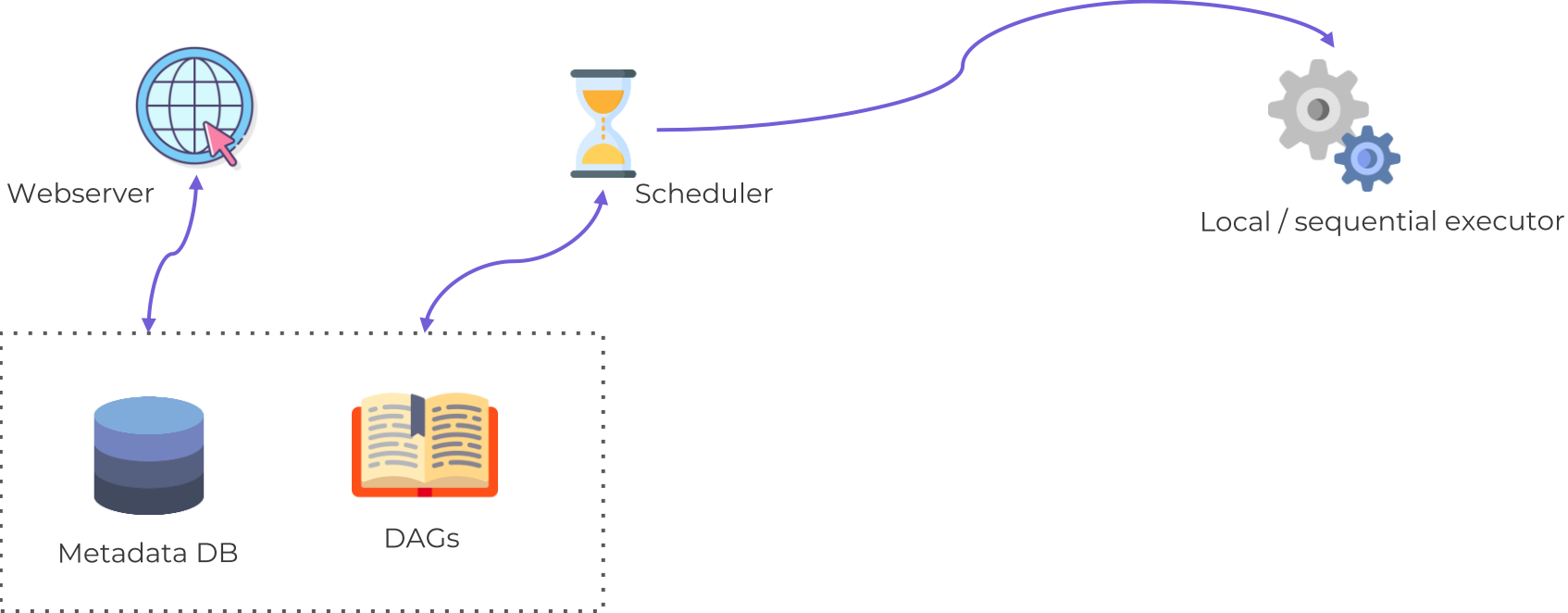Python base image
size: 340MB
time: ~6 minutes

Dockerfile

```
FROM python:3.9.9-slim-buster

ARG AIRFLOW_VERSION=2.2.3

...

ENV BUILD_COMMIT_ID=$BUILD_COMMIT_ID
```

Python slim-buster
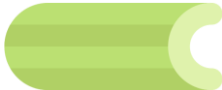size: 60MB
time: 4 minutes

# Airflow basic architecture

Webserver

Scheduler

Local / sequential executor

Metadata DB

DAGs

# Airflow distributed process



Celery
distributed task queue

Webserver

Scheduler

Broker

Workers

Metadata DB

DAGs

# Airflow distributed engine



Kubernetes

Webserver

Scheduler

Broker

Nodes / Workers

Metadata DB

DAGs

# Production architecture



Webserver

Scheduler

Local executor

Metadata DB

DAGs

Compute layer

Google Big Query

Google Cloud Storage

Cloud Dataproc

# CI/CD pipeline

Google Cloud Build

**PREPARE & TRACK TO PRODUCTION**

**DECOUPLED STAGES**

**FASTER RELEASES**

**REPEATABLE**

**MORE ROBUST RELEASES**

**FAILING FAST**

**BETTER VISIBILITY ON CHANGE**

**CENTRALISED ARTIFACTS**

via **cloudbuild.yaml**

# CI/CD pipeline

Github → Cloud Build → Container Registry → Compute Engine

# CI/CD flow



**Github**

*push on production branch
to trigger Google Cloud Build

**Cloud Build**

*build and push Docker image
*restarts Compute Engine

**Container Registry**

*new Airflow image tagged
with $COMMIT_SHA

**Compute Engine**

*production container
updated during startup-script

**Cloud Function**

*subscribed to build
update notifications

**slack**

*build status log

# CI/CD pipeline

**cloudbuild.yaml**

```yaml
steps:
- name: 'gcr.io/cloud-builders/docker'
  args: ['build', '-t', 'gcr.io/$PROJECT_ID/airflow:$COMMIT_SHA',
         '-t', 'gcr.io/$PROJECT_ID/airflow',
         '--build-arg', 'BUILD_COMMIT_ID=$COMMIT_SHA',
         'src/airflow-prod']
- name: gcr.io/cloud-builders/gcloud
  args: [ compute, instances, stop, **-airflow-engine, --zone=europe-***-*]
- name: gcr.io/cloud-builders/gcloud
  args: [ compute, instances, start, **-airflow-engine, --zone=europe-***-*]


substitutions:
  _AIRFLOW_VERSION: 2.0.2

images:
- 'gcr.io/$PROJECT_ID/airflow:latest'
- 'gcr.io/$PROJECT_ID/airflow:$COMMIT_SHA'
```

Github

Cloud Build

# WHY AIRFLOW?
## part II
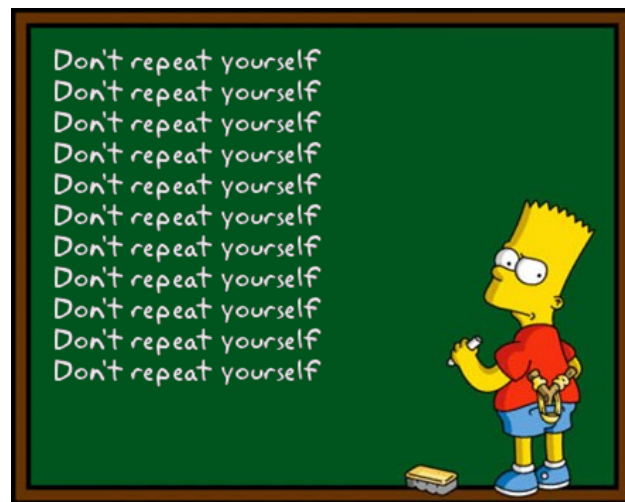
# Lot of errors and time spent



Dashboarding + Analysis

MariaDB

Operational DB

manual code execution to refresh data

Manual refresh

# Pandas + SQLAlchemy



MariaDB

SSH

pandas

SQLAlchemy

Google Sheets

via df2gspread and
gspread-pandas

# Pandas + SQLAlchemy is not DRY

## ~15 lines of code per query/table

```python
import pandas as pd
from sqlalchemy import create_engine
from sshtunnel import SSHTunnelForwarder


server = SSHTunnelForwarder(
    'myuser',
    ssh_username='myuser',
    ssh_pkey='/home/myplace/.ssh/id_rsa',
    remote_bind_address=('0.0.0.0',
                         '37017')
)
```



Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself
Don't repeat yourself

# sqlalchemy-connector



```
pip install sqlalchemy-connector
```



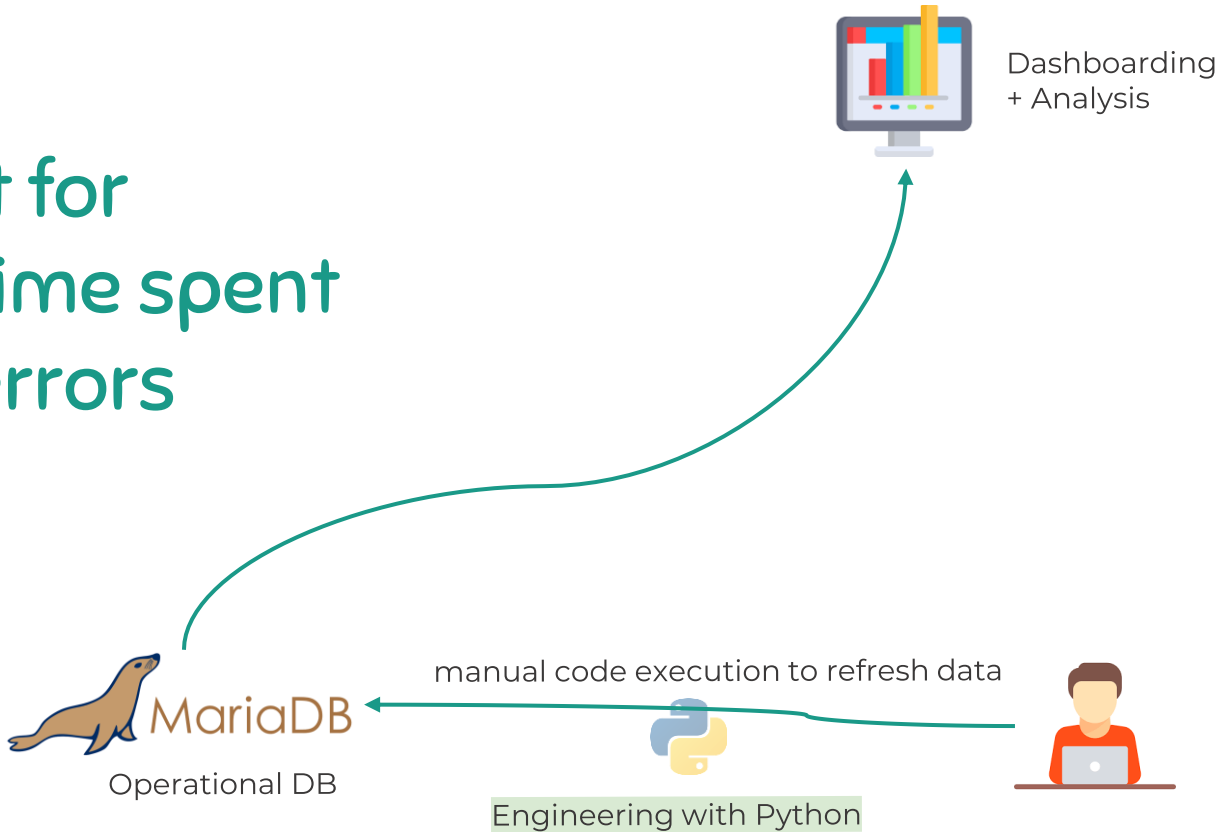INSTALL ME
https://pypi.org/project/sqlalchemy-connector/
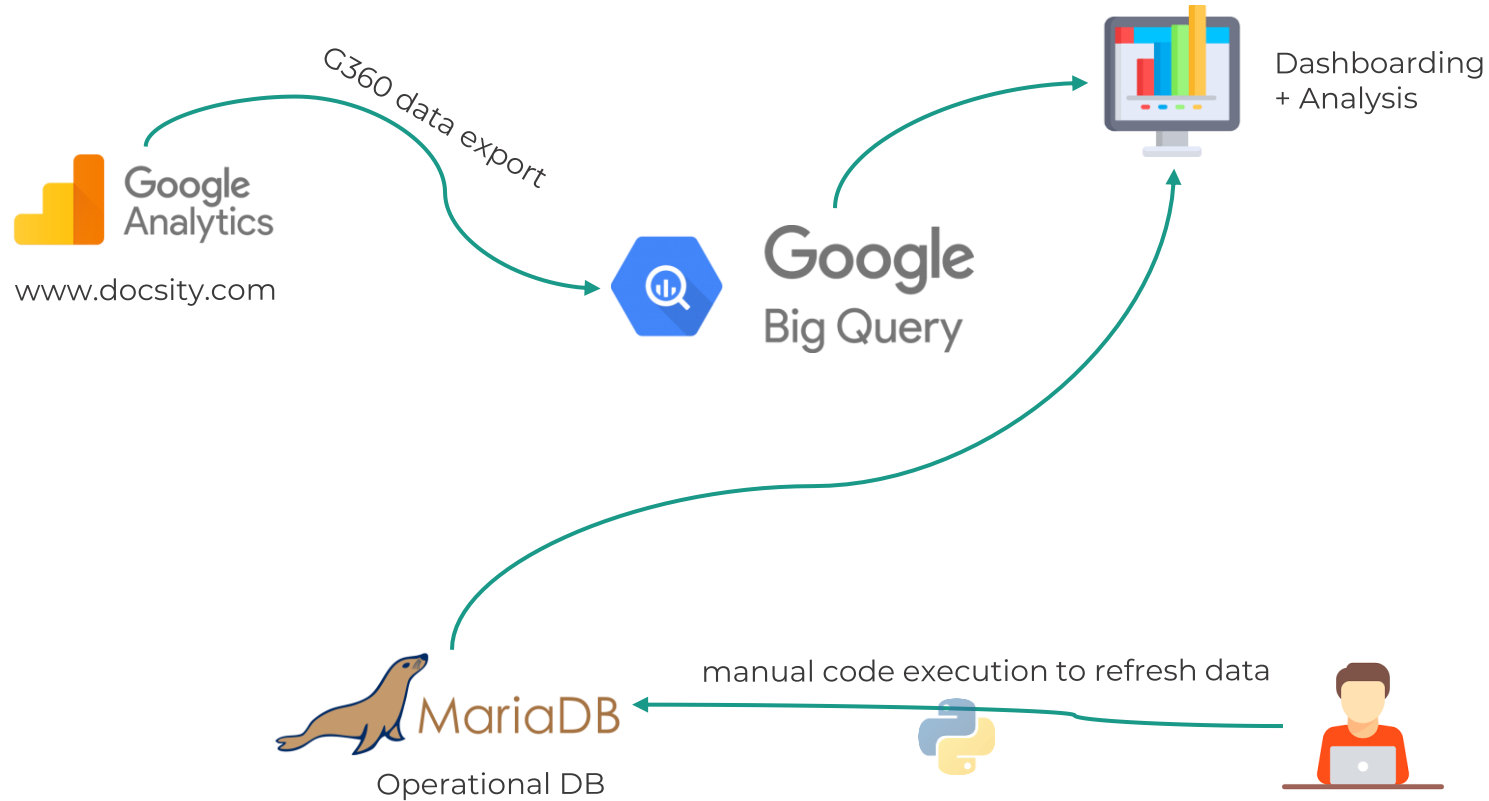
## only 4 lines of code!

```python
from alchemy_connector import SQLAlchemySession


session = SQLAlchemySession(
    host='db.example.com',
    port='21',
    user='myuser',
 key='/home/myplace/.ssh/id_rsa',
    to_port='37017',
    to_host='0.0.0.0'
)
```

# Boost for
# less time spent
# and errors

Dashboarding
+ Analysis

manual code execution to refresh data

MariaDB

Operational DB

Engineering with Python

# 2020 Q1 data lake start



G360 data export

Dashboarding + Analysis

www.docsity.com

Operational DB

manual code execution to refresh data

# Automate BigQuery tables from Cloud Storage



Pandas + Google Cloud
Storage API

Google BigQuery API

# gcloud-connectors



```
pip install gcloud-connectors
```



INSTALL ME
https://github.com/pualien/py-gcloud-connectors

## write

```python
from gcloud_connectors.gstorage import GStorageConnector


gstorage_service = GStorageConnector(confs_path=None)


... load df


gstorage_service.pd_to_gstorage(
    df=df, bucket_name=bucket_name,
```
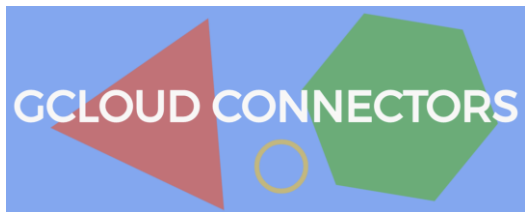
Google Cloud Storage

33

# gcloud-connectors



```
pip install gcloud-connectors
```

Google Big Query

read

```python
import os
from gcloud_connectors.bigquery import BigQueryConnector




project_id = os.environ['PROJECT_ID']
bq_service = BigQueryConnector(project_id=project_id)
```
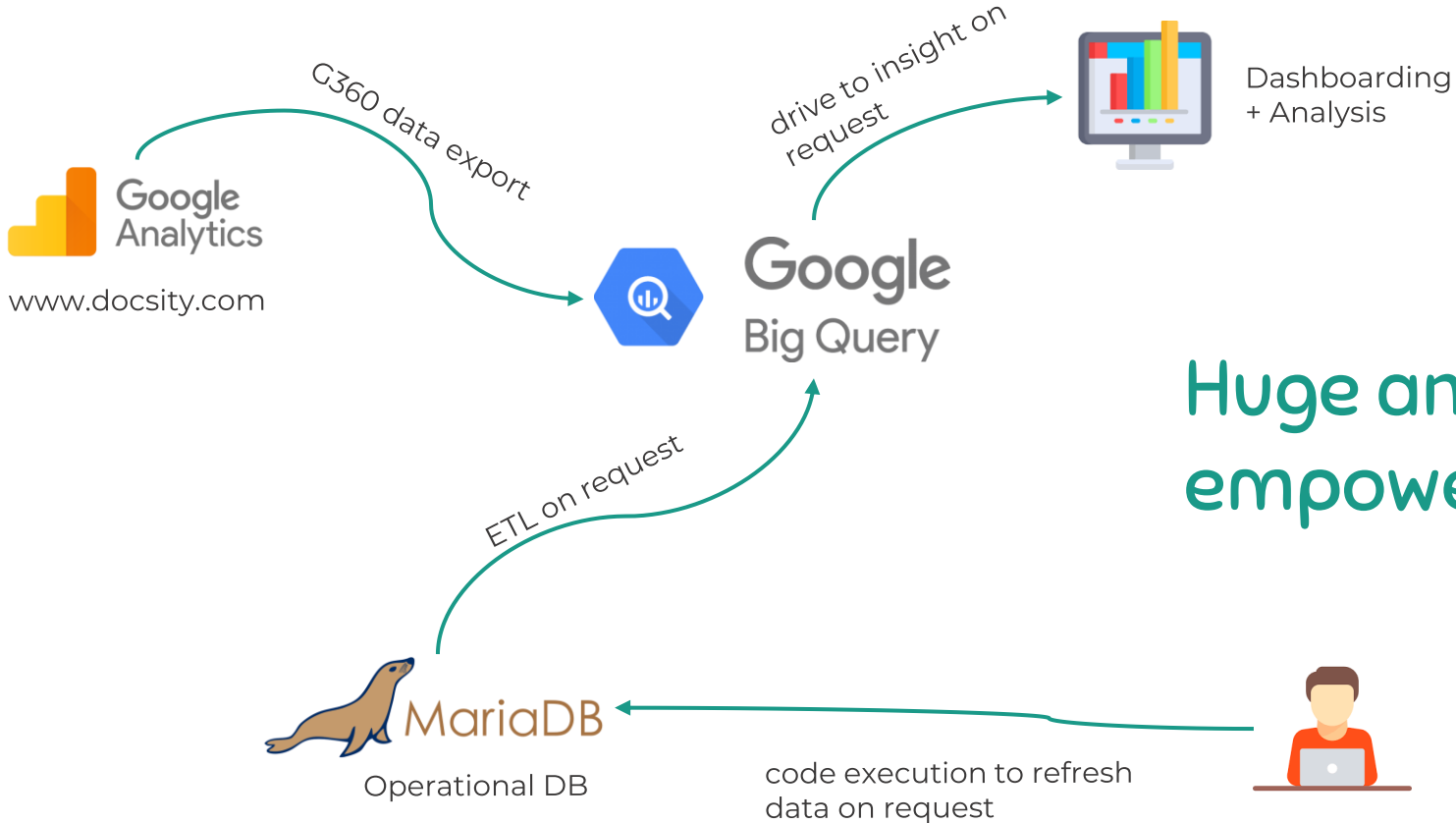
# Extended gcloud-connectors

# 2020 Q2 data lake enhancement



G360 data export

Google Analytics
www.docsity.com

Google Big Query

drive to insight on request

Dashboarding + Analysis

ETL on request

MariaDB
Operational DB

code execution to refresh data on request

Huge analytics empowerment

# Terraform infrastructure as code

**VERSIONING**

**AUTOMATE DEPLOY/RECOVERY**

**EASY TO INTEGRATE WITH CI**
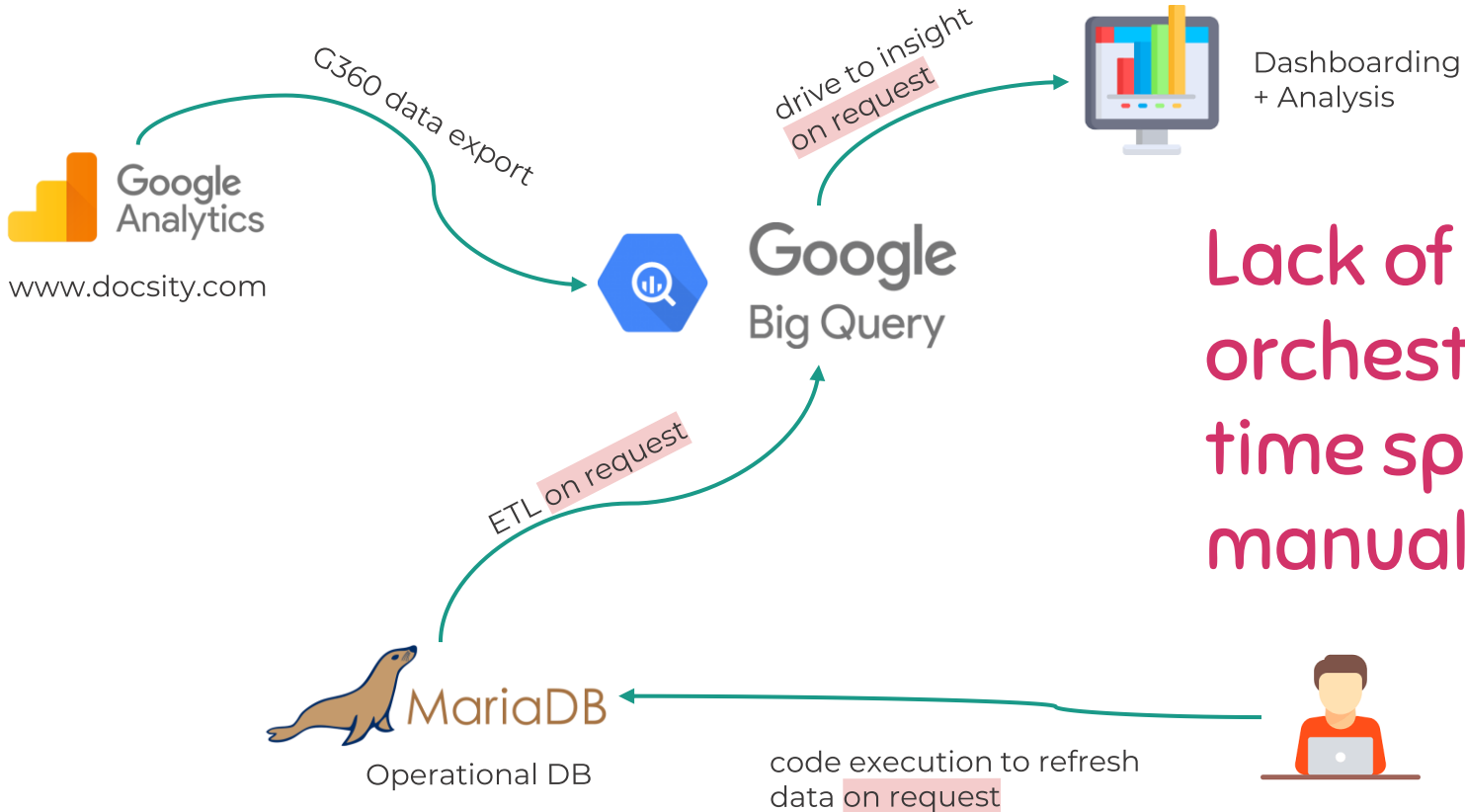
**ONE CLICK TO DESTROY**

**NO NEED TO REPAIR, JUST REDEPLOY**

```
resource "google_bigquery_table" "db_docsity_table" {
    dataset_id = google_bigquery_dataset.db_dataset.friendly_name
    table_id = "db_docsity"
    schema = <<EOF
[ { "name": "id", "type": "FLOAT", "mode": "NULLABLE" } ...]
}
```

BigQuery table

# 2020 Q2 data lake enhancement



Google Analytics

www.docsity.com

G360 data export

Google Big Query

drive to insight on request

Dashboarding + Analysis

ETL on request

MariaDB

Operational DB

code execution to refresh data on request

Lack of orchestration, time spent for manual refresh

# 2020 Q2 data lake orchestration



Google Analytics

www.docsity.com

G360 data export

Google Big Query

drive to insight

Dashboarding + Analysis

Pandas ETL

This is why Airflow matters

MariaDB
Operational DB

automation

Apache Airflow

39

# AIRFLOW SECURITY

# Security for Google Compute Engine



Compute Engine

service account key
& predefined role
with granular permissions
managed by Terraform

IAM & Admin

Apache Airflow

APIs

Google Analytics

Google Cloud Storage

Google Sheets

Google Search Console

Google Drive

Google Big Query

# Airflow encrypted variables & connections

Encrypted credentials pulled from Compute Engine metadata



Airflow  DAGs  Security ▾  Browse ▾  Admin ▾  Docs ▾                                    09:27 UTC ▾  MS ▾

Choose file  No chosen  ☁ Import Variable

## List Variable

Search ▾

+ | Actions ▾ | ←                                                                        Record Count: 4

| | Key | Val | Is Encrypted |
|---|---|---|---|
| ☐ ✎ 🗑 | BUILD_COMMIT_ID | 42354ywirgdkshfhrekth3oit | True |
| ☐ ✎ 🗑 | ENV | PROD | True |
| ☐ ✎ 🗑 | MARIADB_PASSWORD_CREDS | ******** | True |
| ☐ ✎ 🗑 | MONGODB_PASSWORD_CREDS | ******** | True |

# Airflow + Google Secret Manager



Google Secret Manager

credentials pulled from
Google Secret Manager

WHAT TIME TO UPDATE?

# Schedule time?



AWS RDS prod DB

snapshot to sync

AWS Lambda

AWS RDS data DB

right time
to pull data?

# airflow-add-ons



MariaDB

AWS RDS prod DB

MariaDB

AWS RDS data DB

snapshot to sync

AWS Lambda

semaphore file written to
understand when to pull data

AWS S3

**class** ReturnS3KeySensor

Apache Airflow

# airflow-add-ons

```
pip install airflow-add-ons
```



INSTALL ME
https://pypi.org/project/airflow-add-ons/



via **custom Airflow Operators & Sensors**

# 2020 Q3 data lake orchestration



Google Analytics

www.docsity.com

G360 data export

Google Big Query

drive to insight

Looker — Dashboarding + Analysis

Pandas ETL

MariaDB — Operational DB

automation + best time to refresh

Apache Airflow

# pymongo-ssh

```
pip install pymongo-ssh
```

INSTALL ME
https://pypi.org/project/pymongo-ssh/

## DRY

```python
from pymongo_ssh import MongoSession


session = MongoSession(
    host='db.example.com',
    port='21',
    user='myuser',
    key='/home/myplace/.ssh/id_rsa',
    to_port='37017',
    to_host='0.0.0.0'
)
```
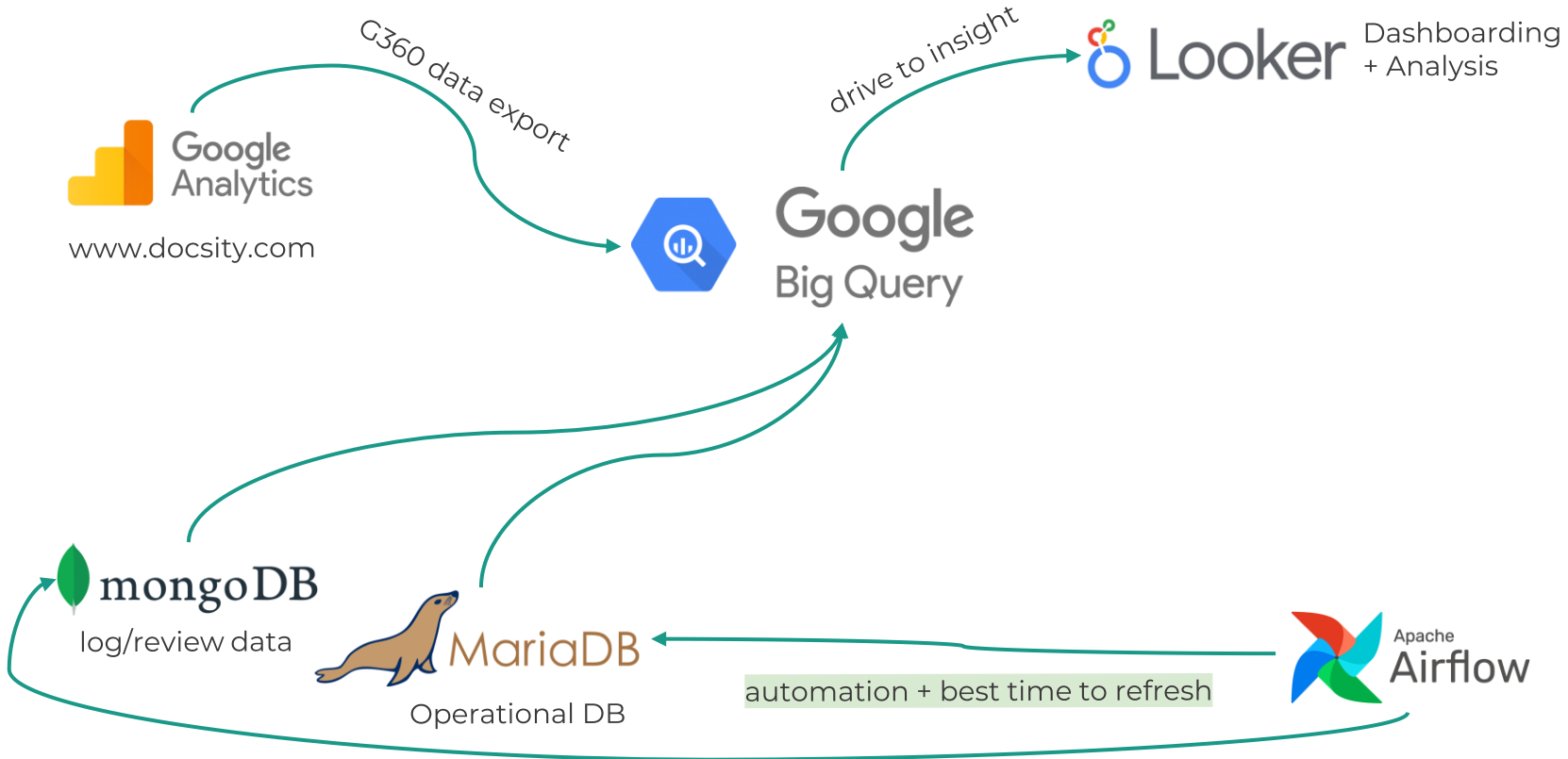
# 2020 Q4 data lake orchestration



Google Analytics
www.docsity.com

G360 data export

Google Big Query

drive to insight

Looker
Dashboarding + Analysis

mongoDB
log/review data

MariaDB
Operational DB

automation + best time to refresh

Apache Airflow

# 2021 Q1 data lake Orchestration



Google Analytics
www.docsity.com

G360 data export

Google Big Query

drive to insight

Looker — Dashboarding + Analysis

mongoDB
log/review data

MariaDB
Operational DB

CRUX report

automation + best time to refresh

Apache Airflow

# 2021 Q2 data lake Orchestration



Google Analytics — www.docsity.com

G360 data export

drive to insight

**Looker** — Dashboarding + Analysis

Google Big Query

mongoDB — log/review data

MariaDB — Operational DB

CRUX report

unbounce

Google Ad Manager — Campaign data

automation + best time to refresh

Apache Airflow

# 2021 Q3 data lake Orchestration



Google Analytics

www.docsity.com

G360 data export

Google Big Query

drive to insight

Looker — Dashboarding + Analysis

unbounce · Facebook Business Manager · Google Ad Manager · customer.io

Campaign data

mongoDB

log/review data

MariaDB

Operational DB

CRUX report

automation + best time to refresh

Apache Airflow

# PRE Airflow



Dashboarding + Analysis

manual code execution to refresh data

MariaDB
Operational DB

Manual refresh

SIMPLE

ERRORS

TIME SPENT

LACK OF CONTROL

# Airflow

# Thanks!

**Does anyone have questions?**

matteo.s@docsity.com
https://github.com/pualien
irene.s@docsity.com
docsity.com