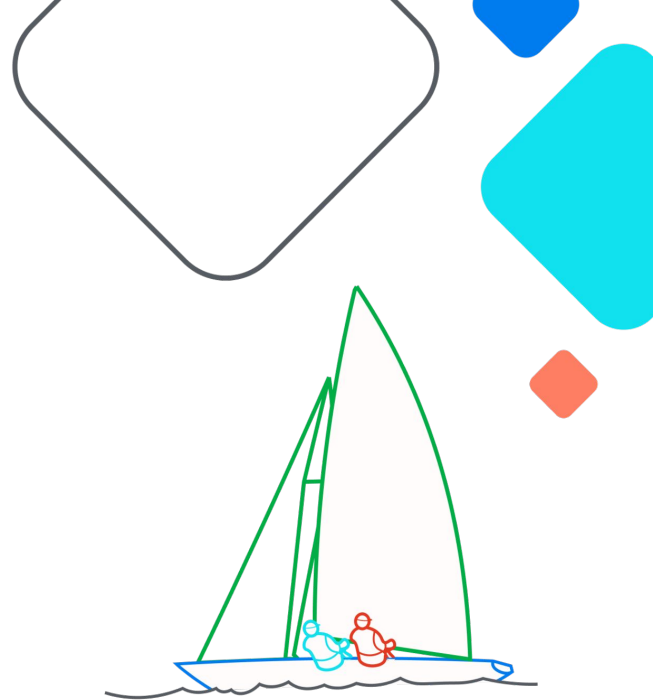# DAG Authoring without PhD

Filip Knapik, Rafał Biegacz

**Airflow Summit**
Let's flow together

September 19-21, 2023,
Toronto, Canada

# About us

**Filip Knapik**

Group Product Manager
Cloud Composer

Working with Airflow for ~4 years
Ex-Product Manager for Google Workflows

**Rafal Biegacz**

Senior Eng Manager
Cloud Composer

Working with Airflow for ~4 years
Member of Airflow Summit Organizing Team

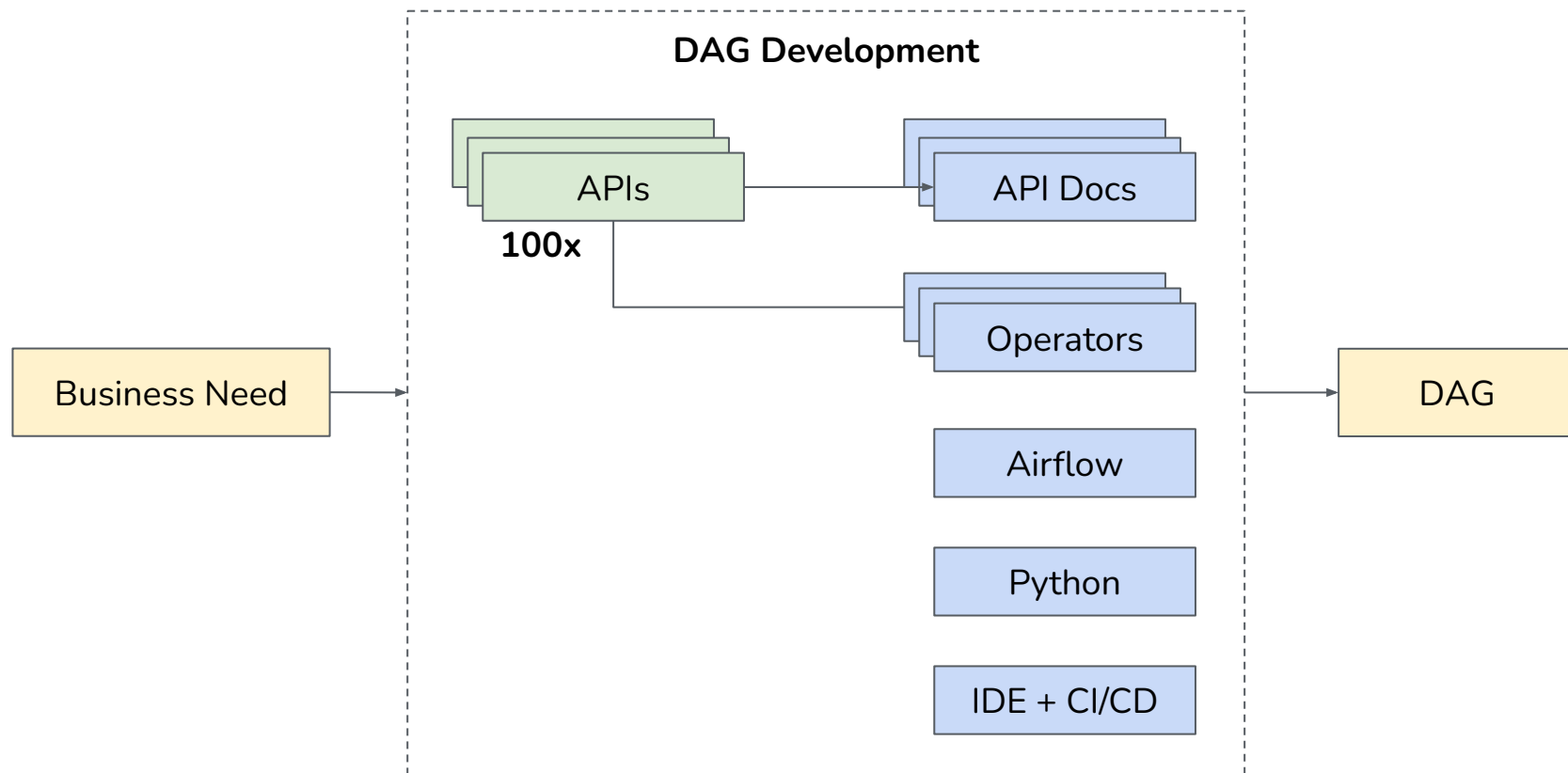Google Cloud

# Contacting Us / Where to Find Us

During the conference **visit us at Google Booth** to learn more about:
- Serverless Composer
- Disaster Recovery & Data Lineage support in Composer
- Support for Public Sector and Assured Workloads

Fill in the form: **bit.ly/airflow-summit-2023-composer** if you would like to meet, request more information or you are interested in getting a voucher for GCP credits.
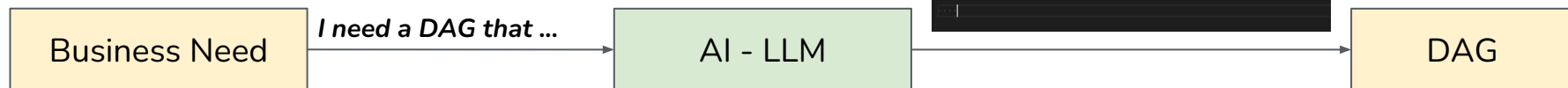
Google Cloud

# DAG Authoring is non-trivial

**DAG Development**

**Business Need** → | APIs (100x) | → API Docs
| Operators
| Airflow
| Python
| IDE + CI/CD |

→ **DAG**

Google Cloud

# Can it be any easier?

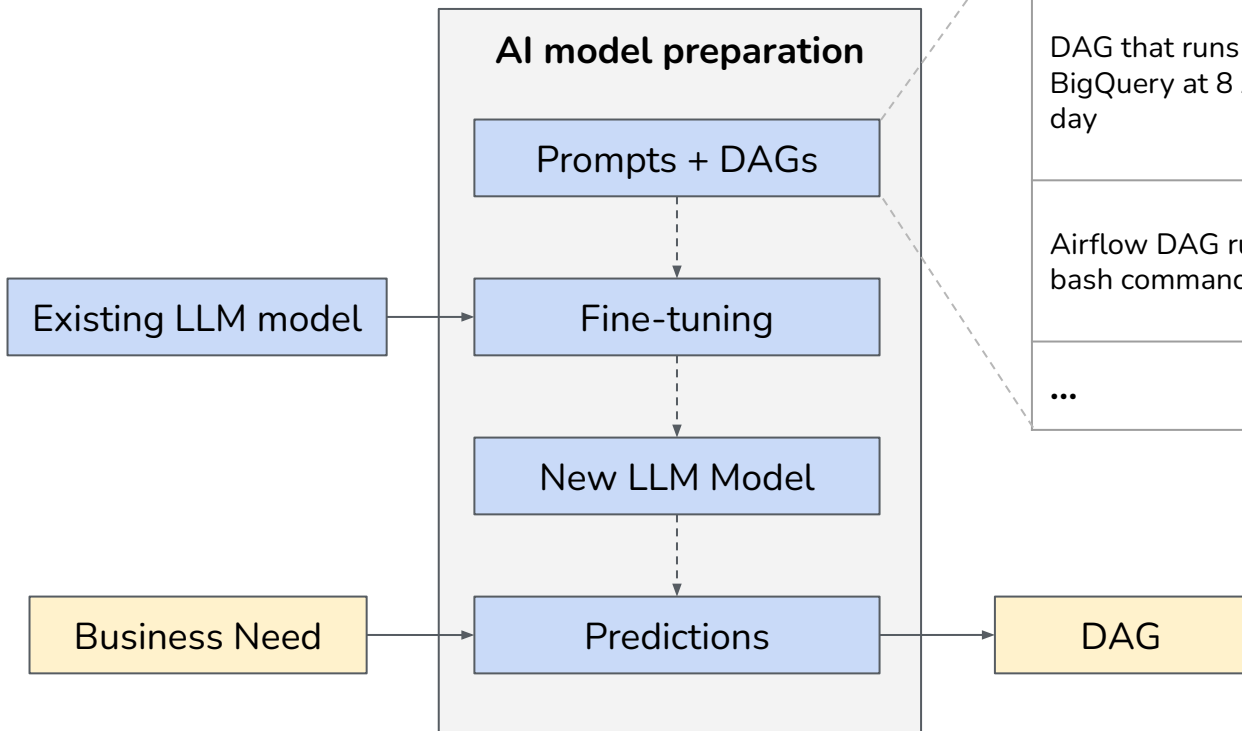LLM-based Generative AI model acting as a translation layer

```
from airflow import DAG
from airflow.providers.google.cloud.operators.bigquery import BigQueryOperator
from airflow.utils.dates import days_ago

with DAG(
    "bigquery-query",
    schedule_interval="0 8 * * *",
) as dag:

    run_query = BigQueryOperator(
        task_id="run_query",
        sql="SELECT * FROM `my_dataset.my_table`",
        use_legacy_sql=False,
    )
```

| Business Need | → *I need a DAG that ...* → | AI - LLM | → | DAG |

Google Cloud

# How to get there?



AI model preparation

Prompts + DAGs

Existing LLM model → Fine-tuning

New LLM Model

Business Need → Predictions → DAG

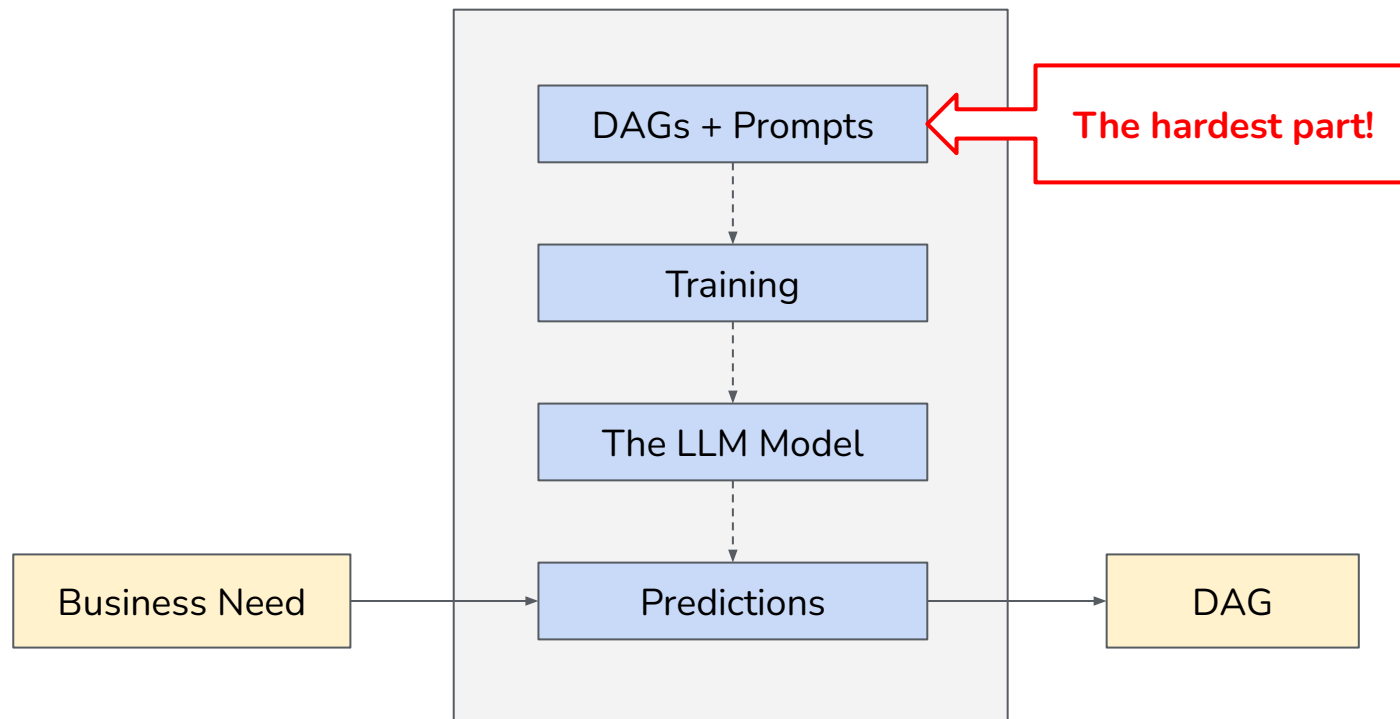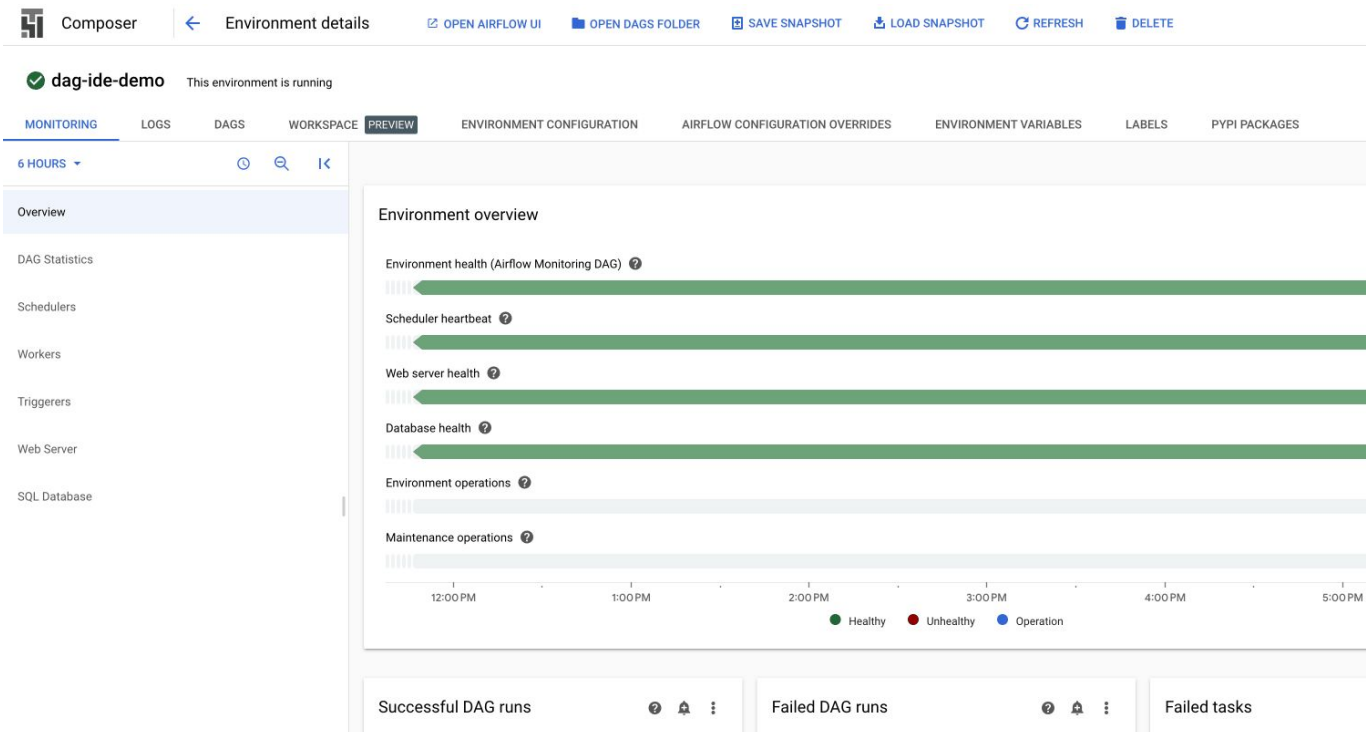| Prompt | DAG |
|---|---|
| DAG that runs a query in BigQuery at 8 AM every day | ```<br>from airflow import DAG<br>from airflow.providers.google.cloud.operators.bigquery import BigQueryOper<br>from airflow.utils.dates import days_ago<br><br>with DAG(<br>    "bigquery-query",<br>    schedule_interval="0 8 * * *",<br>) as dag:<br><br>    run_query = BigQueryOperator(<br>        task_id="run_query",<br>        sql="SELECT * FROM `my_dataset.my_table`",<br>        use_legacy_sql=False,<br>    )<br>``` |
| Airflow DAG running 10 bash commands in parallel | ```<br>from airflow import models<br>from airflow.operators.bash_operator import BashOperator<br><br>with models.DAG(<br>    "Bash_10_parallel",<br>    schedule_interval="0 0 * * *", # Override to match your needs<br>) as dag:<br><br>    for k in range(10):<br>        newStep = BashOperator(<br>``` |
| ... | |

# The challenge

# Where we are

**Prompt:** *Write Airflow DAG that that runs BigQuery query using BigQueryInsertJobOperator to get number of GitHub commits in May 2022. Please, use public dataset for github data published by Google Cloud Platform*
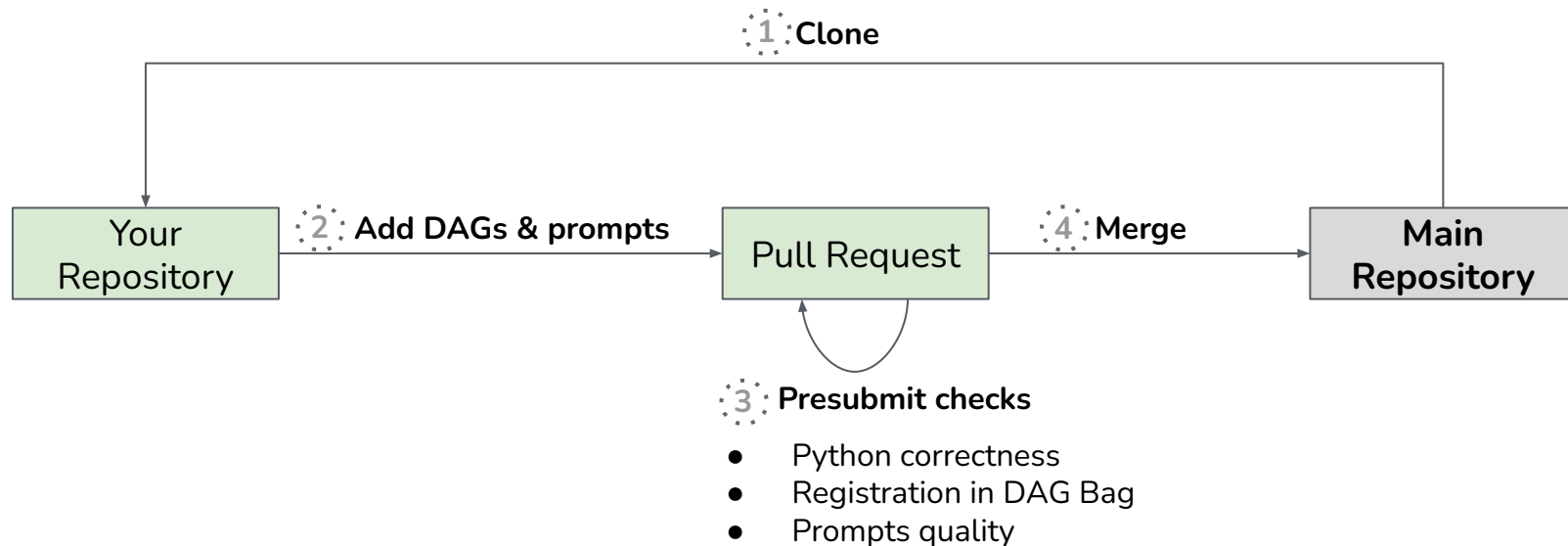
# Our proposal

**Establish an open source Airflow DAGs + Prompts repository to enable LLM models creation & experimentation by the community**

- Google establishes open source repository
- Google donates its Airflow DAGs training data
- Community joins in and contributes with its own DAGs examples
- Community can build its own DAG code generation models and tools
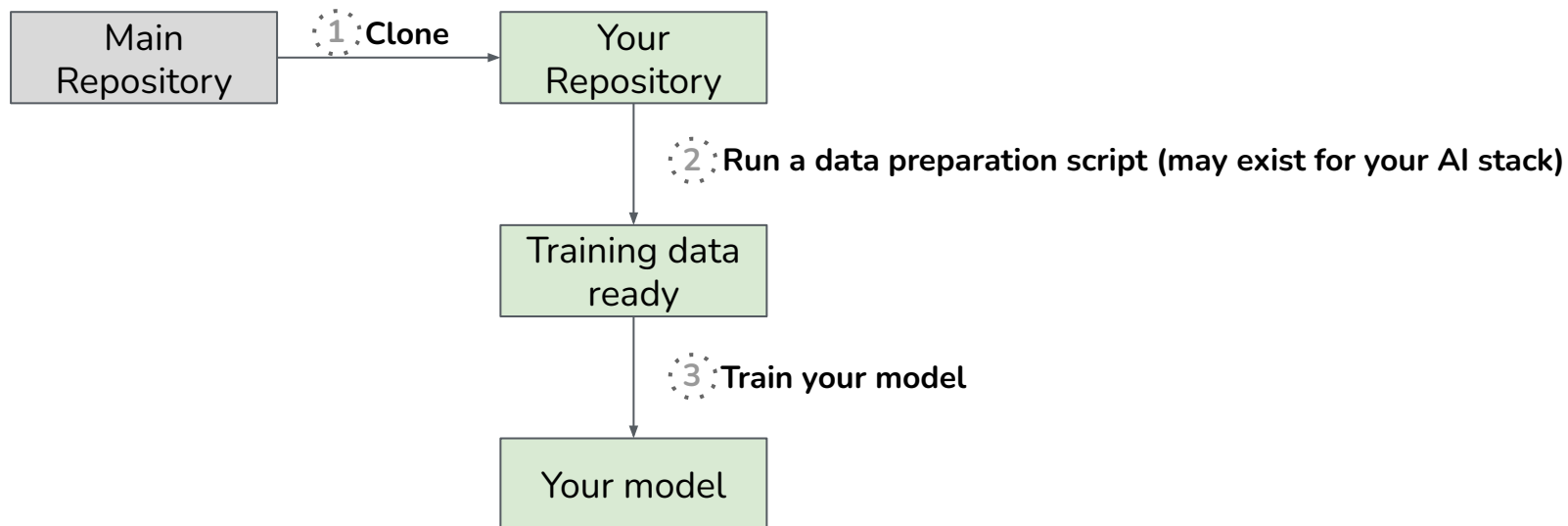
Google Cloud

# Contribution process

# Contribution guidelines

- The more (DAG, prompt) pairs are contributed, the better
- Ensure that DAGs and prompts are:
    - Free of any PII references
    - Following latest Airflow development practices
    - Using particular versions of python modules and provider packages
    - As diverse as possible
- Use examples that you can contribute under Apache 2.0 license

Google Cloud

# How to use it?

```
┌─────────────┐   1  Clone    ┌─────────────┐
│    Main     │ ───────────►  │    Your     │
│ Repository  │               │ Repository  │
└─────────────┘               └─────────────┘
                                     │
                                     │  2  Run a data preparation script (may exist for your AI stack)
                                     ▼
                              ┌─────────────┐
                              │ Training data│
                              │    ready    │
                              └─────────────┘
                                     │
                                     │  3  Train your model
                                     ▼
                              ┌─────────────┐
                              │ Your model  │
                              └─────────────┘
```

Google Cloud

# Next Steps

1. Composer team will send a call for action email on Airflow Dev list
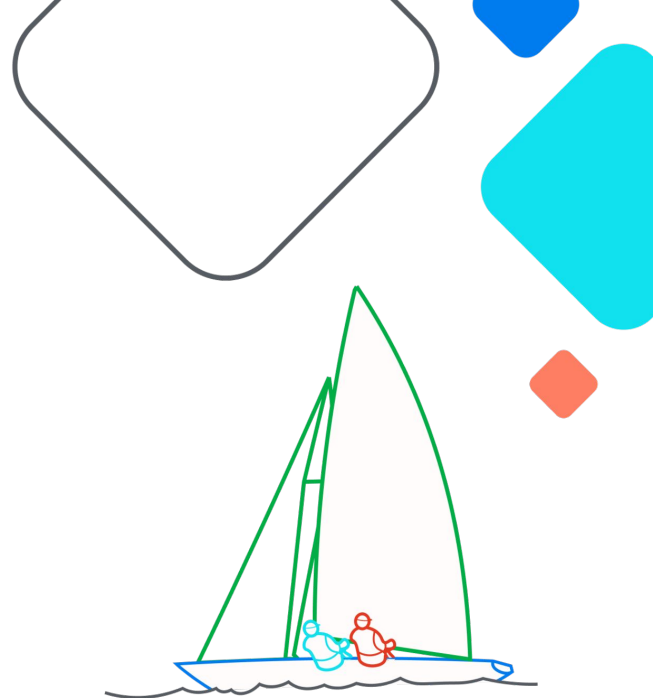2. All community members are encouraged to join us!

# Questions?

Optionally share some contact info like email, blog or social media handles