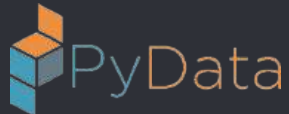


- Building a Data Platform from scratch

Rodel van Rooijen

PyData Amsterdam 2024



About me



Over 8 years of experience in building data infrastructure and data products.

● ING

○ Data Scientist

● Adyen

○ Senior ML Scientist

○ Tech Lead (Manager)

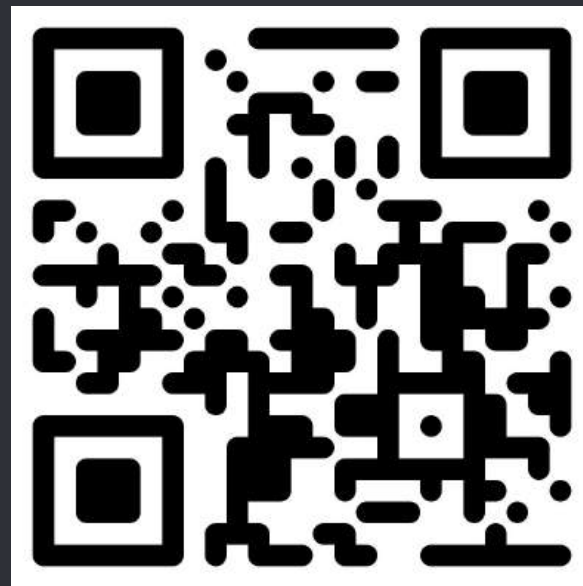
○ Engineering Lead (Director)

● Solvimon

○ Founding Engineer, Data

● Achmea

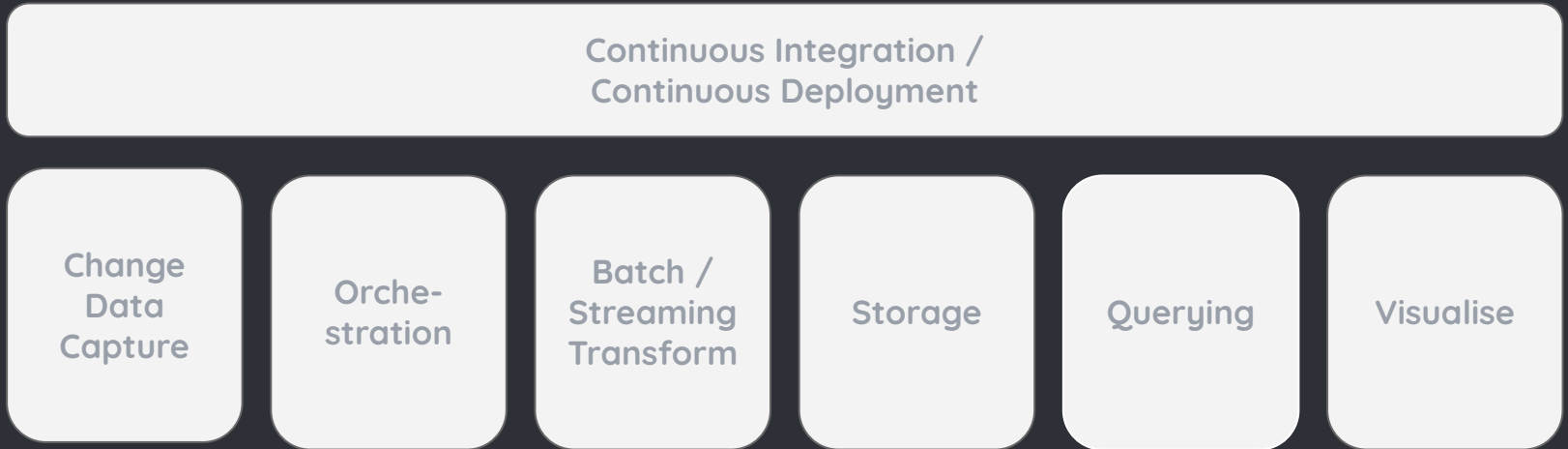
○ Senior ML Engineer



<https://rodel.dev/>



- Dissecting a data platform





Continuous Integration /
Continuous Deployment



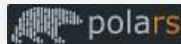
Change
Data
Capture



Kubeflow



Orche-
stration



Batch /
Streaming
Transform



databricks



Storage



PostgreSQL



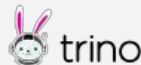
snowflake



databricks



Storage



Querying



Visualise

**NOT SURE WHICH TOOL I
SHOULD USE**

SO MANY CHOICES

memegenerator.net

Where to start?



Cloud
Platform

1. Pick cloud platform (or use existing).



My team's
knowledge

2. Start with what you know.



Use
open-source

3. Use open-source solutions where possible.

● What are the (hosting) options?

Self-hosting

- ++ Full control
- ++ High degree of flexibility
- + Lowest immediate cost
- Long implementation time
- Maintenance & expertise

Managed (open-source)

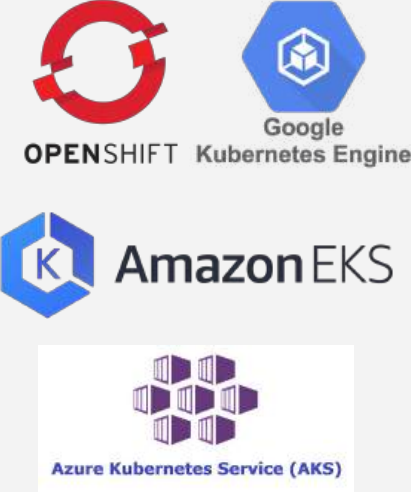
- ++ Less expertise required
- + Decent flexibility
- + Low implementation time
- Less control
- Can be more expensive

Proprietary

- + Low implementation time
- +/- Can be more cost efficient
- Least amount of flexibility
- Least level of control




How to self-host?



Logos for OpenShift, Google Kubernetes Engine, Amazon EKS, and Azure Kubernetes Service (AKS).

Managed Kubernetes

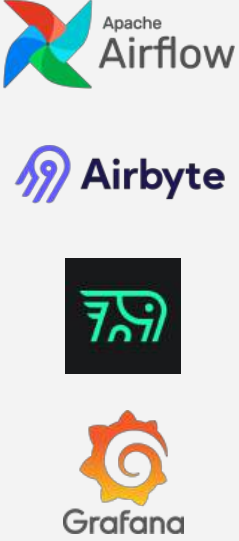
1. Pick a managed Kubernetes platform



Github repo

Set-up infrastructure

2. Set-up cluster with auto-scaling

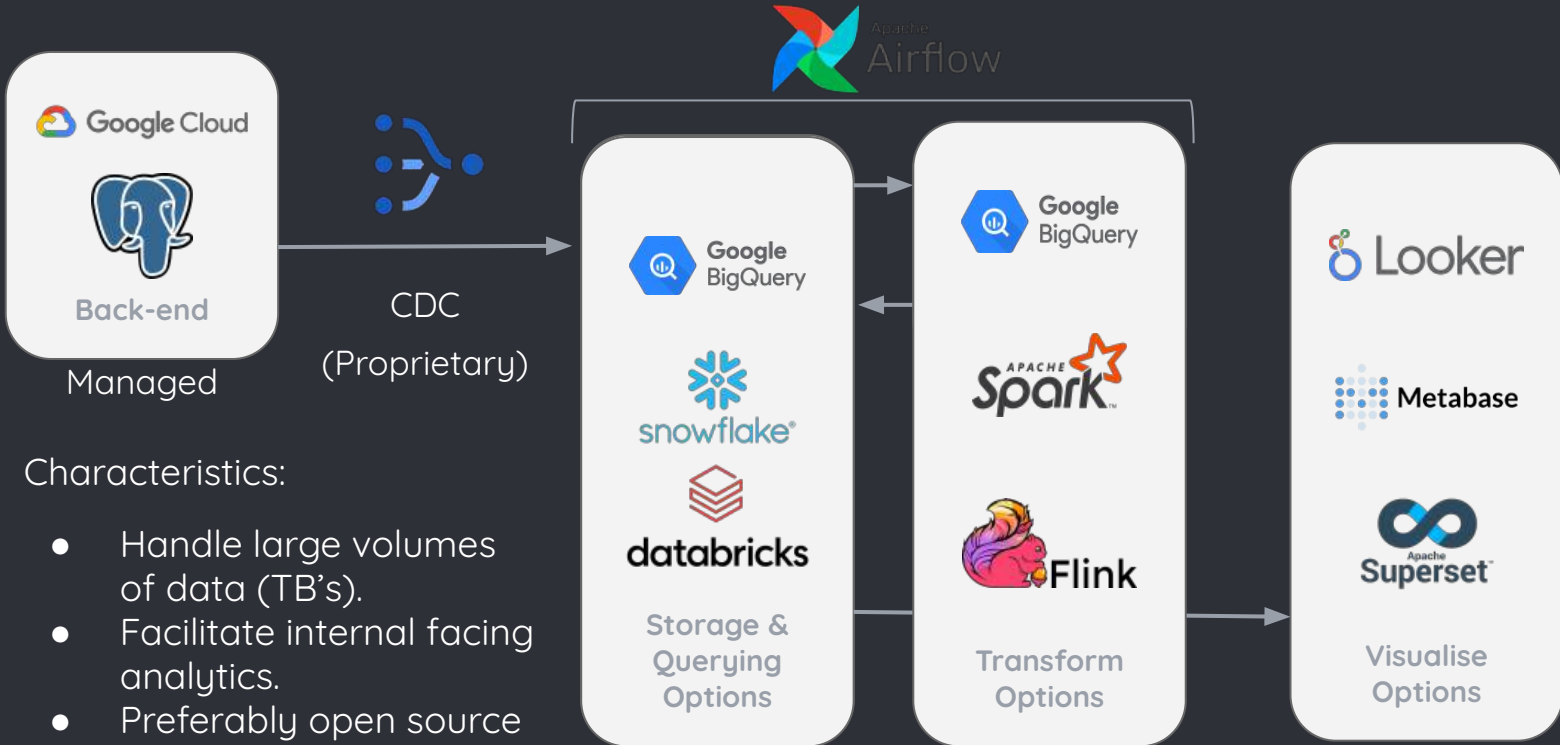


Logos for Apache Airflow, Airbyte, and Grafana.

Deploy using helm charts

3. Use pre-configured helm charts and customise values

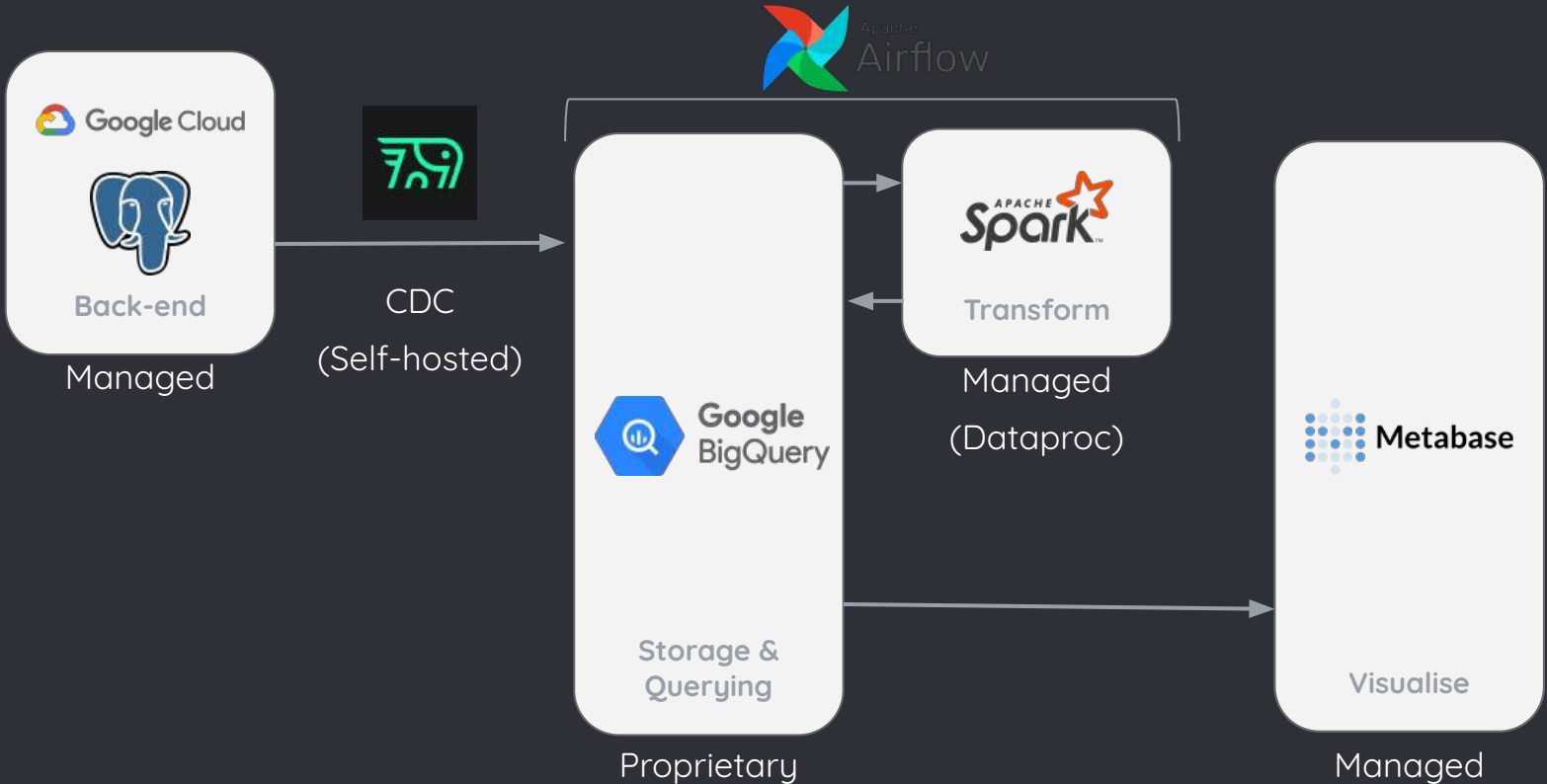
Building the data platform (starting options)



Characteristics:

- Handle large volumes of data (TB's).
- Facilitate internal facing analytics.
- Preferably open source and managed tools.
- Keep costs under control.

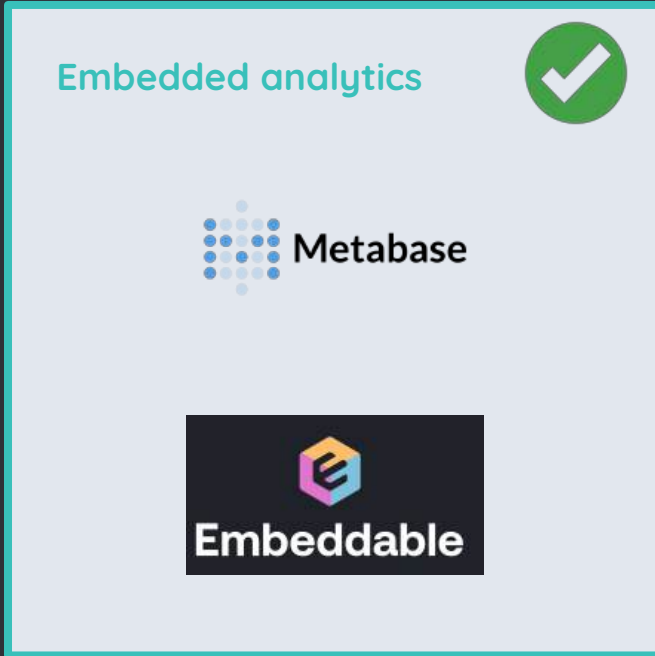
- Building the data platform (conclusion)





Medium blogpost

- Value-add data products



Medium blogpost

Data imports & exports



Authentication

Send information to front-end

Back-end

FastAPI

- List imports/exports
- Initiate new import/export set-up
- List schedules of syncs

Cloud Run

Perform sync

Airbyte

Store credentials

Apache Airflow

Set-up schedule

Data Platform

Azure aws Google Cloud

Customer's



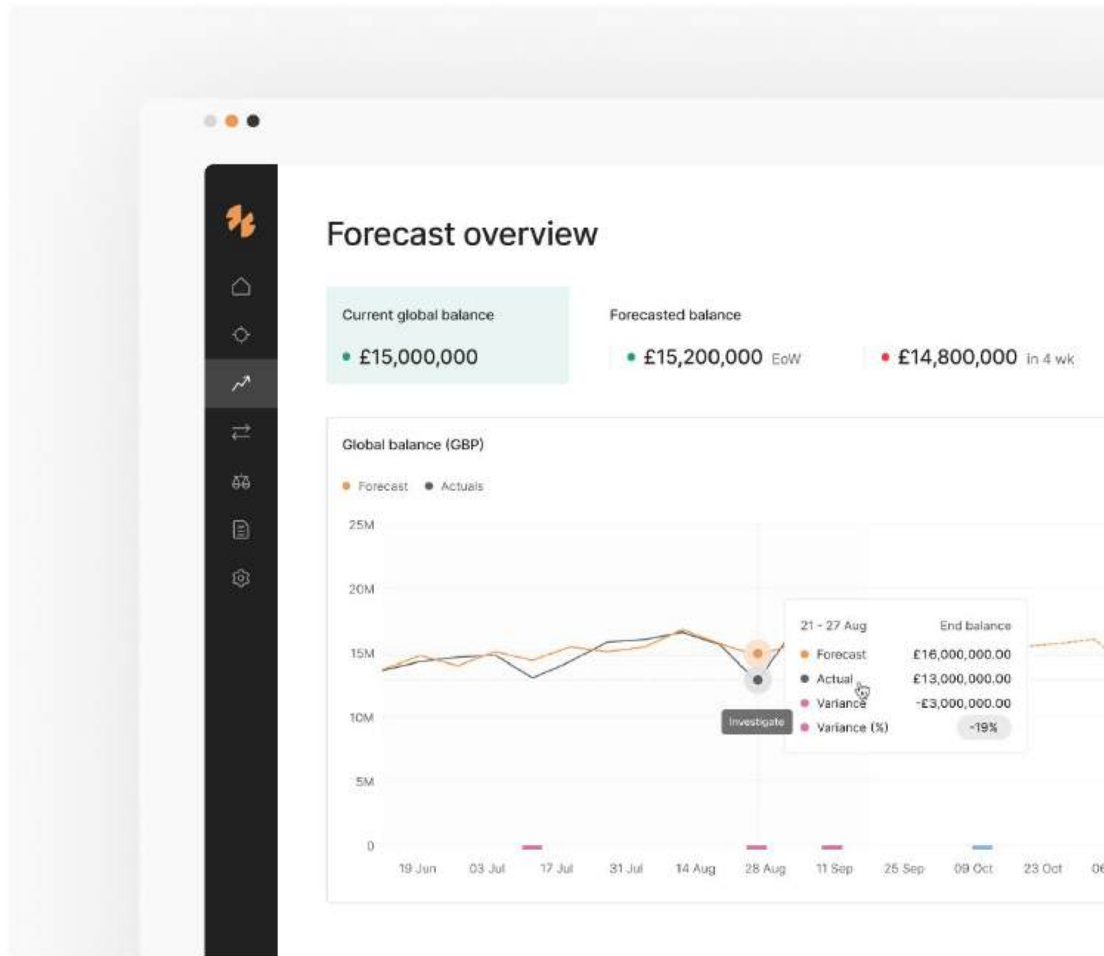
- Summary

Start early, don't wait

Use what you are familiar with

Design with value-add products in mind

I am Joining
 Palm



We are Hiring

Back-end Engineers

ML Scientists



usepalm.com

