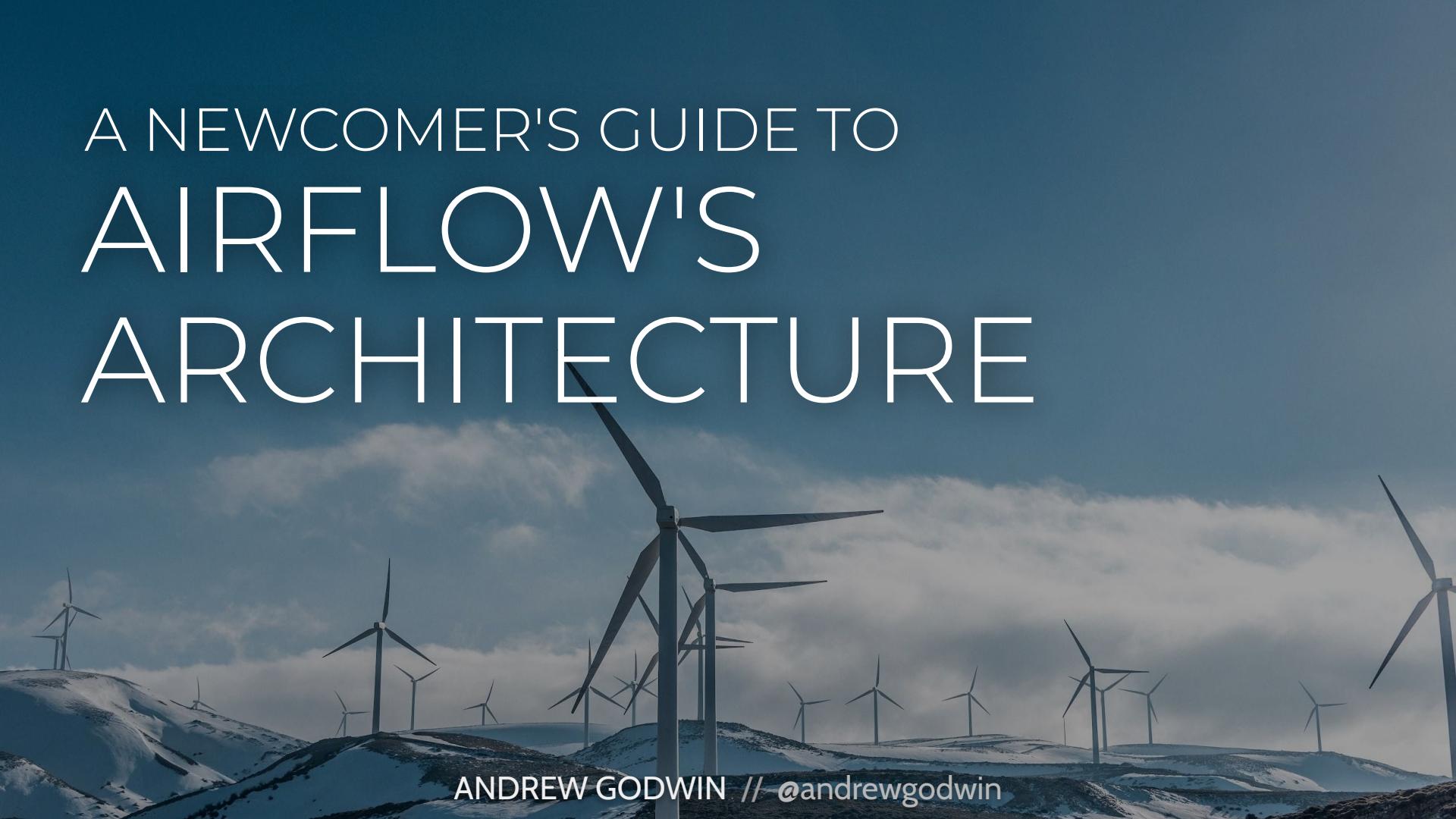


A NEWCOMER'S GUIDE TO AIRFLOW'S ARCHITECTURE



ANDREW GODWIN // [@andrewgodwin](https://twitter.com/andrewgodwin)

The background of the slide features a dark, silhouetted forest of evergreen trees against a night sky filled with numerous stars.

Hi, I'm

Andrew Godwin

- Principal Engineer at  ASTRONOMER
- Also a Django core developer, ASGI author
- Using Airflow since 2021



High-Level Concepts

What exactly is going on?

The Good and the Bad

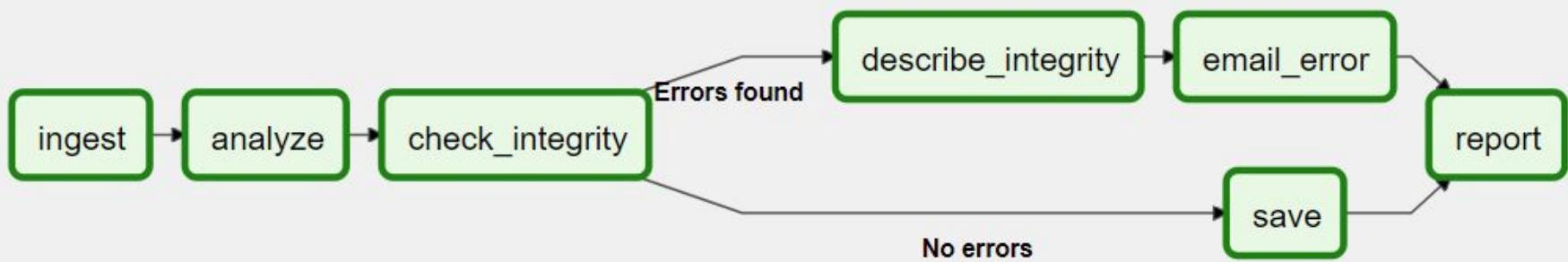
Or, How I Learned To Stop Worrying And Love The Scheduler

Problems, Fixes & The Future

Where we go from here

Airflow grew organically

It started off as an internal ETL tool



DAG ➔ DagRun

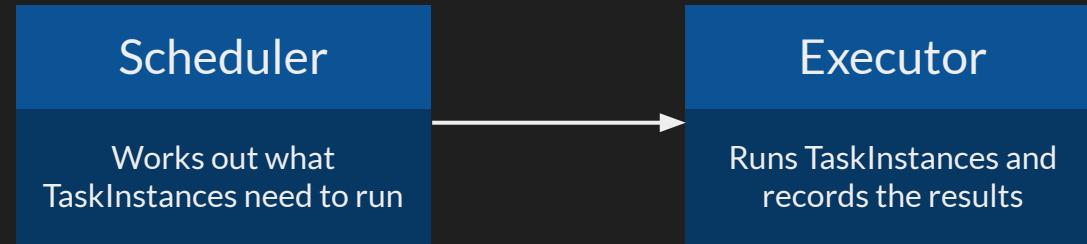
One per scheduled run, as the run starts

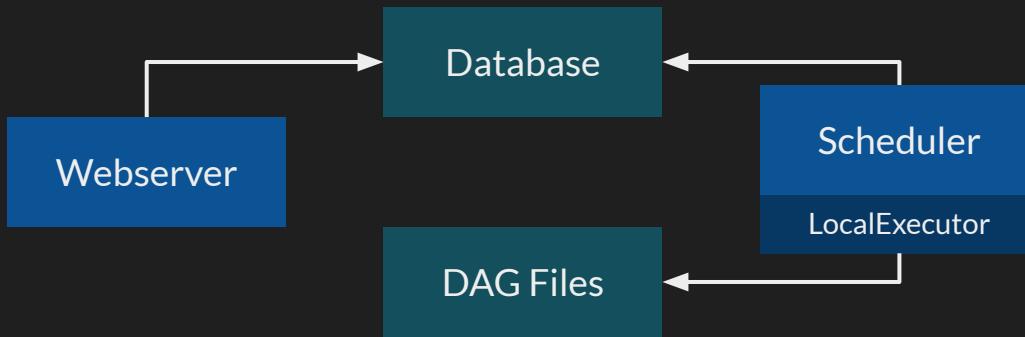
Operator ➔ Task

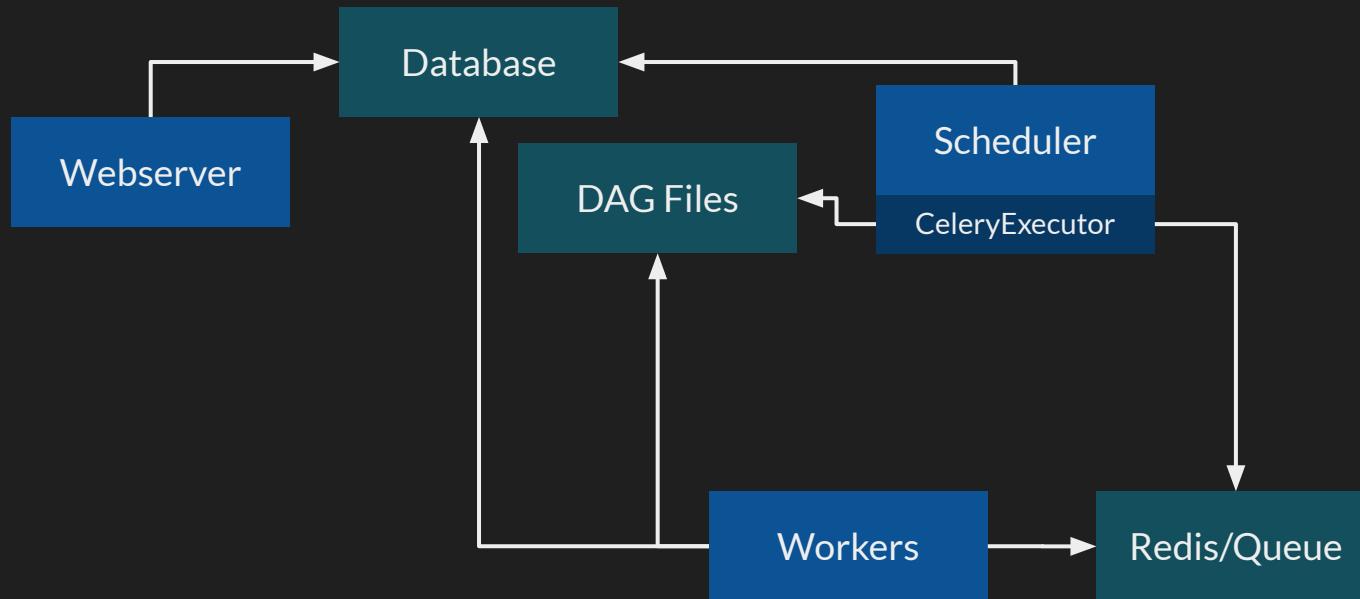
When you call an operator in a DAG

Task ➔ TaskInstance

When a Task needs to run as part of a DagRun









The Executor runs **inside** the Scheduler

Its logic, at least, and the tasks too for local ones

A wide-angle photograph of a modern, multi-story library. The building features multiple levels of white bookshelves filled with books. People are seen sitting on green couches in the central atrium. The architecture is characterized by a series of diagonal walkways and stairs. The overall atmosphere is bright and spacious.

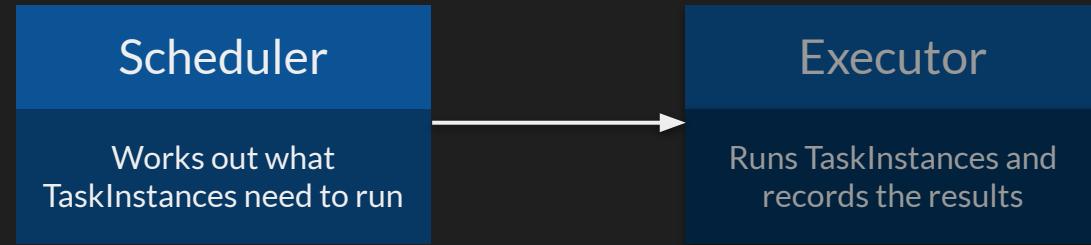
Everything talks to the database

It's the single central point of coordination

Scheduler, Workers, Webserver

All can be run in a high-availability pattern





Timing

Dependencies

Retries

Concurrency

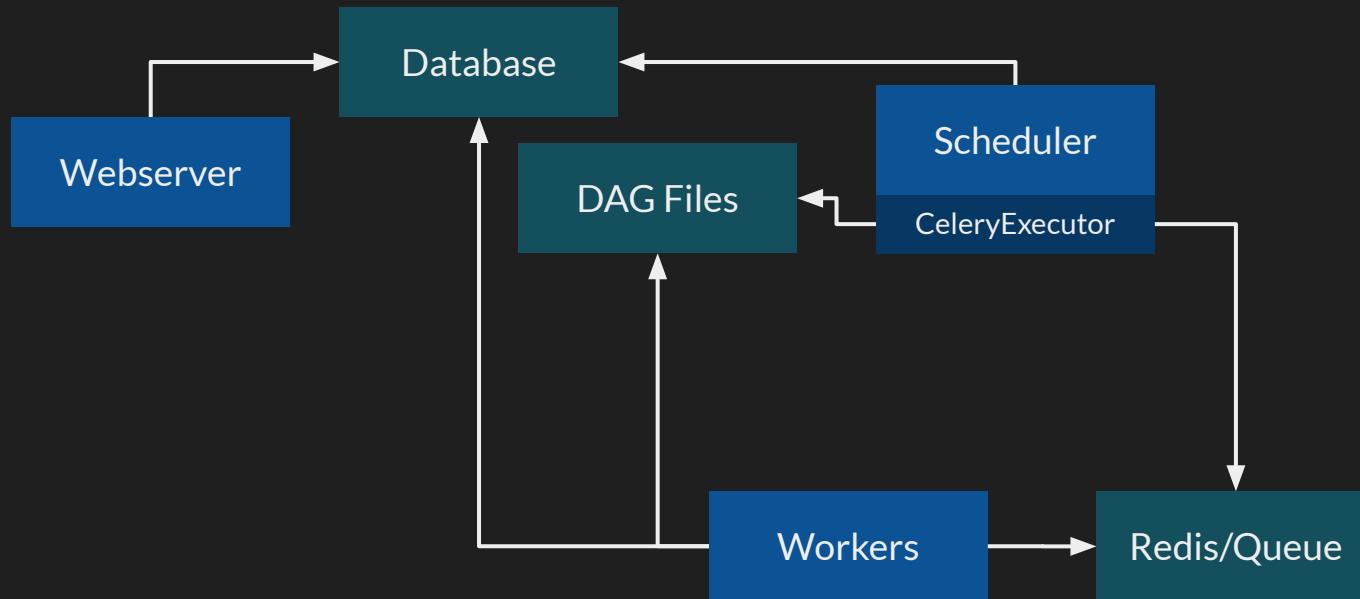
Callbacks

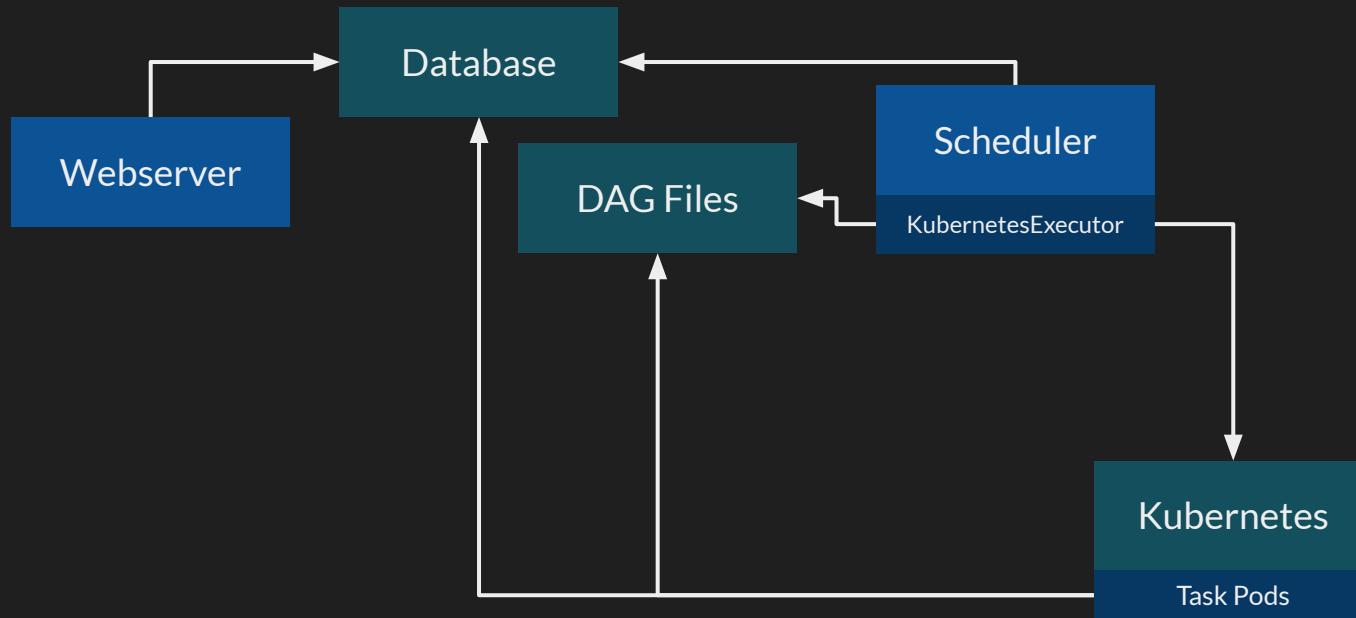
...

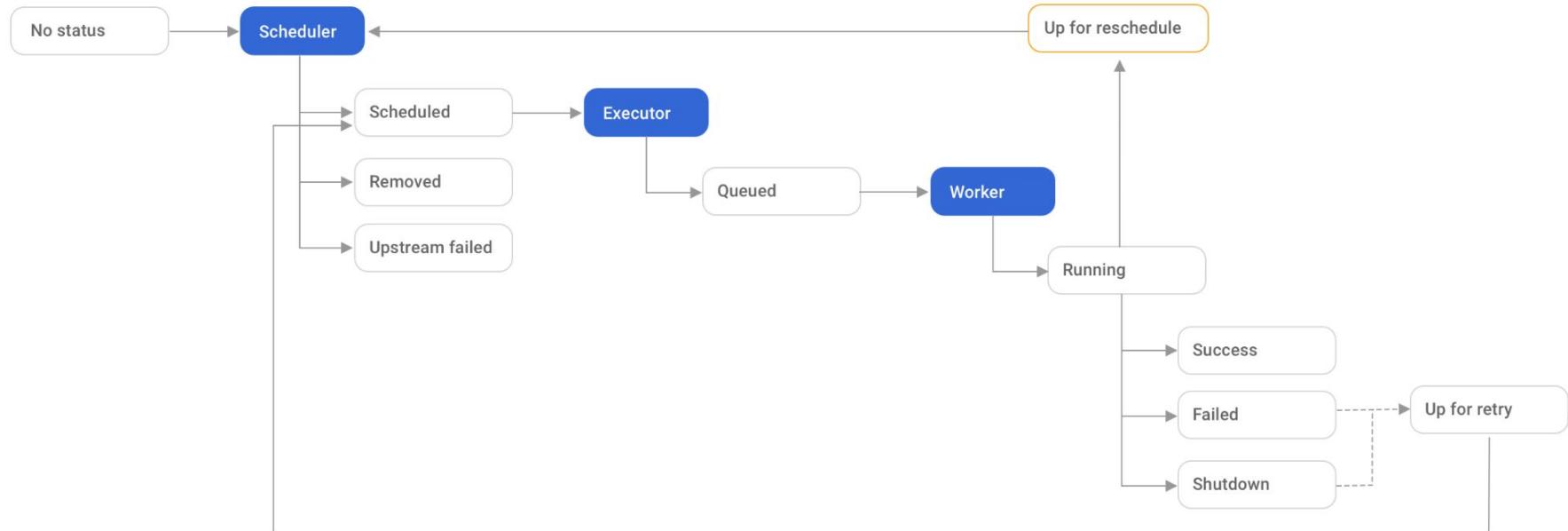


Celery or Kubernetes

Our two main options, currently







■ Component

□ Task stage

□ Task stage only for sensor

— Stage transition

--- Alternative stage transition

Tasks are the core part of the model

DAGs are more of a grouping/trigger mechanism

Very flexible runtime environments

Airflow's strength, and its weakness

Airflow doesn't know what you're running

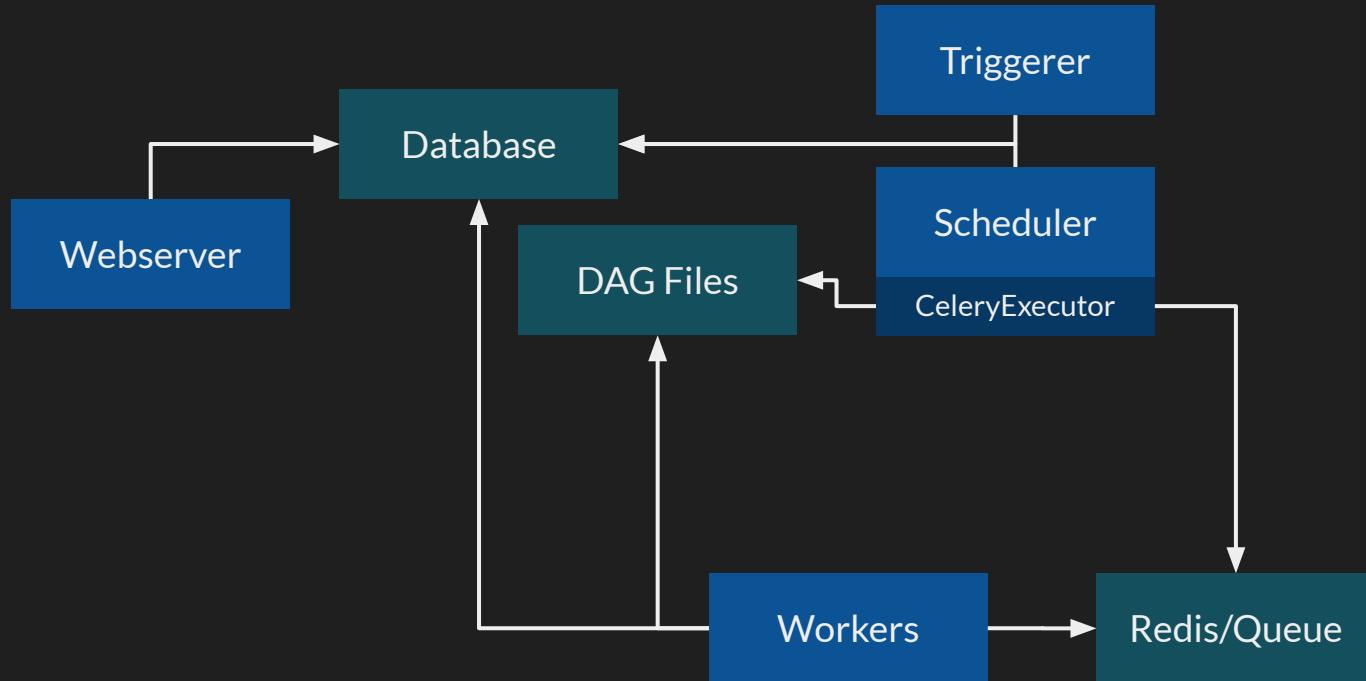
Though this is also kind of an advantage.

What can we improve?

Let's talk about The Future

More Async & Eventing

Anything that involves waiting!



Removing Database Connections

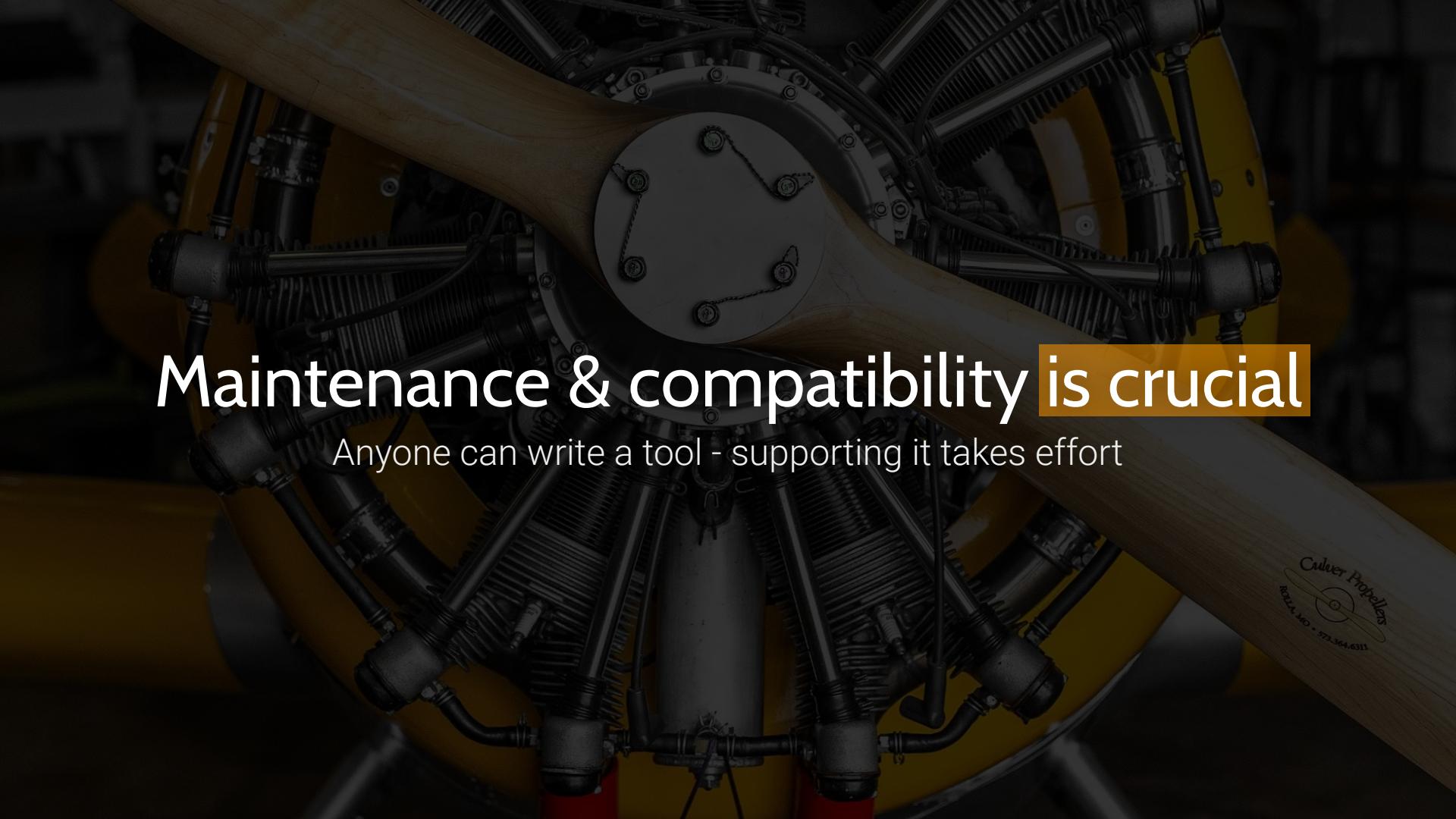
APIs scale a lot better!

I do like the database, though

There's a lot of benefit in proven technology

Software Engineering is not just coding

Any large-scale project needs documentation, architecture, and coordination



Maintenance & compatibility is crucial

Anyone can write a tool - supporting it takes effort

The background of the slide features a wide-angle photograph of a mountainous landscape. In the foreground, there's a grassy field with a paved stone path leading towards the center. The middle ground shows a valley with more green fields and a small cluster of trees. The background is filled with layers of mountains, some of which are partially obscured by low-hanging clouds or mist. The overall lighting is soft, suggesting either early morning or late afternoon.

Airflow is forged by people like you.

Coding, documentation, triage, QA, support - it all needs doing.

Thanks.

Andrew Godwin

@andrewgodwin

andrew.godwin@astronomer.io