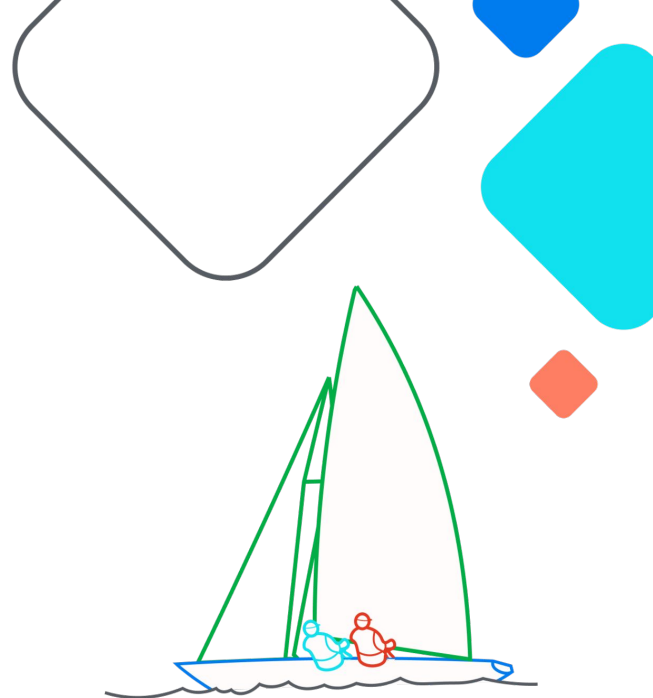


# Traps and misconceptions of running reliable workloads

Bartosz Jankiewicz



 **Airflow Summit**

Let's flow together

September 19-21, 2023,  
Toronto, Canada

# Cloud Composer: Apache Airflow in Google Cloud



Bartosz Jankiewicz  
Engineering Manager

# Lesson 1

Can I always run tasks in Apache Airflow in reliable manner?

# Sources of failures

## External

Originate from dependent services, latency, inconsistent data, network connectivity

## Internal

Intrinsic to our setup, our code running in Apache Airflow, our actions

# Lesson 2

Apache Airflow is a distributed system but you need to understand how to use the redundancy

# Some of single points of failures



**Network**



**Metadata  
database**



**Celery worker**

# What about redundancy?

## Redundancy can help only when a component is stateless or its state can be recovered

Typically executor process failure leads to task instance failure. To alleviate this, tasks should be configured with retries.

Running tasks in deferrable mode makes the tasks stateless from Airflow perspective.



# Lesson 3

How can you improve availability of metadata database?

# Database failures happen



Too many tasks



Badly written  
DAGs



Sensors



Hardware  
failures

# Some countermeasures

- Scale up
- Defer
- Limit parallelism
- Change schedule

- Increase poll interval
- Defer
- Change architecture

- Use Variables with care
- Jinja template can reduce db calls

- Highly available database

The background is a solid red color with several white diagonal lines crossing it. One line runs from the top-left towards the bottom-right. Another line runs from the top-right towards the bottom-left. A third line runs from the bottom-left towards the top-right, crossing the other two.

# Watch out!

Database scaling is not a free lunch!

**More DB  
resources**

**Bigger  
load**

**More data**

**Longer  
maintenance**

**Cost**

# More effective solutions

Schedule tasks more evenly



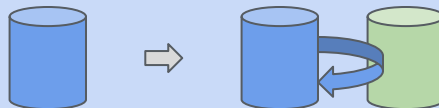
Spread load across more Airflow instances



Use Variables with care + Jinja templates



Highly available database



# Lesson 4

How can I make sure my Python code runs uninterrupted?

# What can disrupt your task?

Other scripts were running on the same host

Database being overwhelmed

Dependant service failure

Maintenance operations on VMs

Network latency



# Quick recap of Airflow executors

## Local Executors

SequentialExecutor

LocalExecutor

## Remote Executors

CeleryExecutor

KubernetesExecutor

CeleryKubernetesExecutor

# Quick summary of actions

Isolate tasks e.g. with K8s executor

Don't overload metadata database

Run tasks in deferrable mode

Plan your maintenance windows

# What if everything fails?



Have a plan 🧐



# Questions?

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

