

DATA COUNCIL SINGAPORE 2019

ARGO: **KUBERNETES NATIVE** **WORKFLOWS &** **PIPELINES**

GREG ROODT - DATA ENGINEERING LEAD

The Canva logo is centered at the bottom of the slide. It consists of the word "Canva" in a white, cursive script font, set against a dark purple circular background.

Canva

AGENDA



INTRO TO CANVA



**WHAT IS A
WORKFLOW?**



ARGO VS AIRFLOW



ARGO DEMO



Q & A

I.

INTRO TO CANVA



”

STARTUP FROM SYDNEY

“



”

**EMPOWERING THE
WORLD TO
DESIGN**

“



II.

WHAT IS A WORKFLOW?



OXFORD ENGLISH DICTIONARY

”

NOUN

The sequence of industrial, administrative, or other processes through which a piece of work passes from initiation to completion.

“

<https://www.lexico.com/en/definition/workflow>

”

**A workflow consists of an orchestrated
and repeatable pattern of activity...**

“

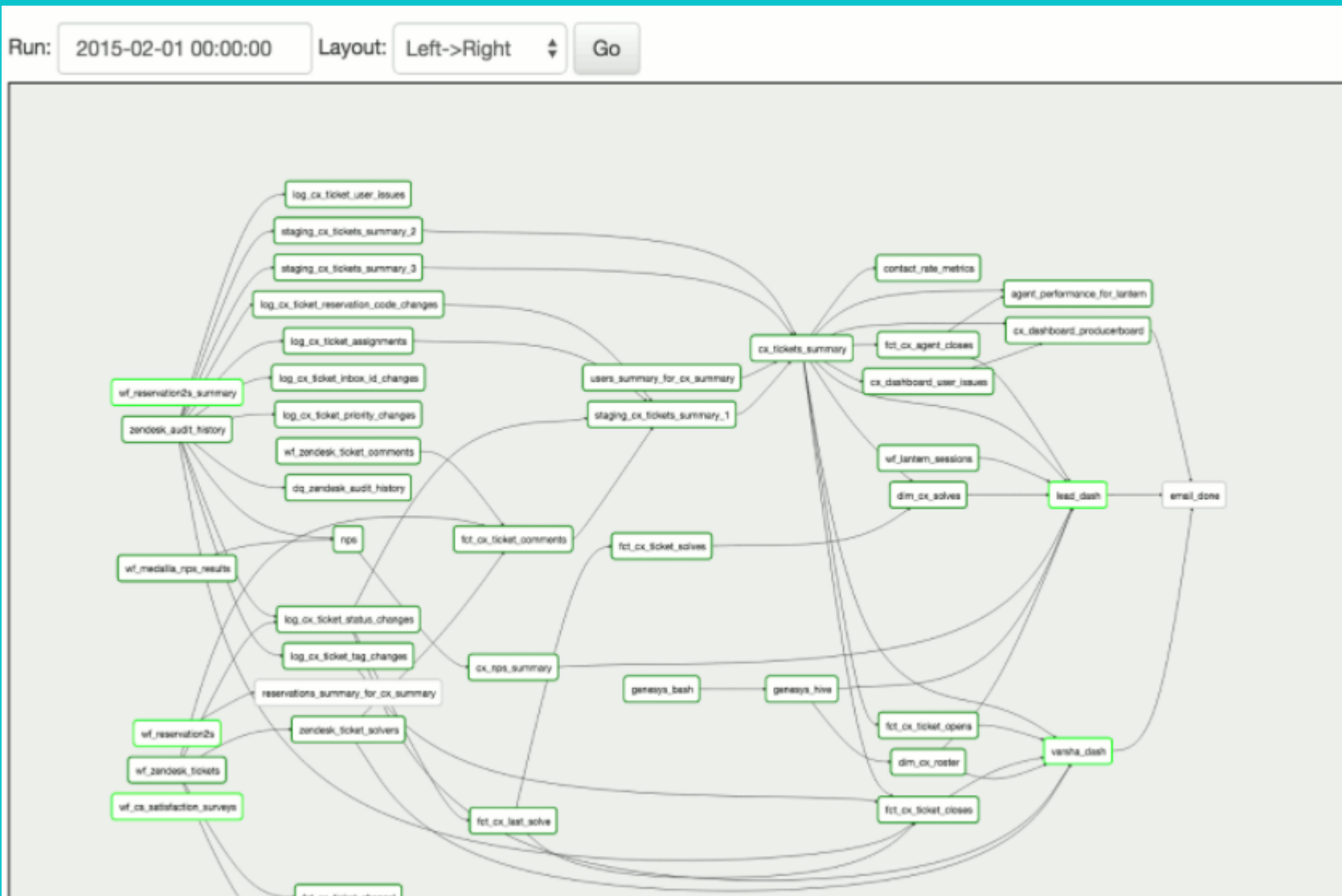
DIRECTED ACYCLIC GRAPH (DAG)

- We typically mean a DAG
- Steps in a DAG are "tasks"
- Each "task" is some discrete amount of work
- All the "tasks" together make up some larger piece of work
- Examples: ETL jobs, Spark jobs, Machine Learning jobs, CI jobs, deployment jobs etc.

WHERE TO START?

- So many different tools out there
- Oozie, Luigi, NiFi, Azkaban...
- The most popular one at the moment seems to be Apache Airflow

AIRFLOW DAG



AWS DATAPIPELINE

Data Pipeline > List Pipelines > Execution Details > Architect: new_DB_Test (df-031851030E8VOGW9LZE8) [Active]

Add Save Activate Export View/Edit tags

The screenshot displays the AWS Data Pipeline console interface. On the left, a pipeline diagram shows the following flow: 'Configuration Default' (purple) branches into 'S3PrefixNotEmpty Input Data Exists' (purple), 'Ec2Resource Micro-EC2' (grey), and 'S3PrefixNotEmpty Output Data Exists' (purple). 'S3PrefixNotEmpty Input Data Exists' leads to 'S3DataNode InputDataNode' (yellow). 'S3DataNode InputDataNode' and 'Ec2Resource Micro-EC2' both lead to 'ShellCommandActivity DB REST CALL' (blue). 'ShellCommandActivity DB REST CALL' and 'S3PrefixNotEmpty Output Data Exists' both lead to 'S3DataNode OutputDataNode' (yellow). 'S3DataNode OutputDataNode' leads to 'SnsAlarm Success SNS Alarm' (purple). On the right, the 'Activities' panel is open for 'DB REST CALL', showing configuration details: Name: DB REST CALL, Type: ShellCommandActivity, Input: InputDataNode, Runs On: Micro-EC2, and Command: curl -X POST -u. A legend at the bottom right identifies node types: DataNodes (yellow), Schedules (pink), Resources (grey), Preconditions (brown), Others (purple), and Parameters (blue).

Activities

- DB REST CALL
- Name: DB REST CALL
- Type: ShellCommandActivity
- Input: InputDataNode
- Runs On: Micro-EC2
- Command: curl -X POST -u
- Add an optional field...

Legend:

- DataNodes
- Schedules
- Resources
- Preconditions
- Others
- Parameters

ARGO DAG



III.

ARGO **VS** AIRFLOW



FIRST A STORY

- **We were a very small team at this point...**
- **First there was AWS Datapipeline**
- **And it was good**
- **Complicated tool, but it's a complicated problem and it was reliable**

TIME PASSES

- Time went on. Team grew in ambition
- Started running into limitations with Datapipeline
- Frustrations grew

SOME ISSUES

- **Limited EC2 instance types (less true now)**
- **Custom software is awkward**
- **Scheduling and retries is awkward, we built our own workarounds**
- **Not very well known tool (no community)**

ARGO **AIRFLOW**

- **We decided to compare alternatives**
- **Airflow was the obvious choice**
- **Argo was a new project that ticked a lot of boxes on paper**
- **Did a comprehensive proof-of-concept to compare**

AIRFLOW

- Originated at Airbnb
- Python based operators & DAG
- Huge community
- Open source (Host yourself)
- Commercial SaaS offerings available too



<https://airflow.apache.org/>

ARGO

- **Originated at Intuit (Applatix)**
- **K8s Native (CRD+Controller)**
- **Written in Golang**
- **DAG of Containers (Any language)**
- **Growing community**
- **Host your own on K8s**



<https://argoproj.github.io/>

PROOF OF CONCEPT

- Decided on simple, realistic example
- Setup both on AWS
- Came up with a few evaluation criteria
- Got to work!



EVALUATION CRITERIA

- Suitability for Canva (mostly Spark) workloads
- Ease of setup and operation
- Good workflow deployment
- Logging
- Live UI updates
- Timeouts & retries
- Reproducibility



THINGS WE WERE NOT LOOKING FOR

- Visual UI for editing DAGs
- Extremely Dynamic DAGs
- Python-only
- Updating a DAG after workflow submission



GOOD PARTS?

- Both tools are very capable
- Both tools could support Canva workloads
- Both tools offer live UI & logging updates
- Both tools offer timeouts and retries



ISSUES WITH AIRFLOW

- Deployment of DAGs is clumsy (poll git, rsync?)
- Mutable DAGs
- DAG deployed is not necessarily the DAG that runs
- Installation of components fairly complex
- A lot of excitement around k8s executor, so why not use k8s directly?



ISSUES WITH ARGO

- Lots of YAML
- Requires k8s
- Simpler, less powerful UI
- Early project



THE DECISION?

- Both projects very capable
- Airflow lacks solid deployment story
- Chose Argo



IV.

ARGO DEMO



**THANK
YOU**

**WE'RE HIRING!
CANVA.COM/CAREERS**

Canva