# Moving Data with Apache Airflow at Zendesk

An overview of Airflow and how we use it at Zendesk for data engineering

Pitt Fagan
21 September 2017

Pitt Fagan

zendesk

# Today's agenda

1. Introductions

2. Overview of Airflow (components, architecture, etc.)

3. Demo of the GUI

4. Using Airflow at Zendesk

5. Questions

# A bit about me …

1. Education
   - B.A. Geology from Rice University
   - M.S. Soil Science from UW-Madison
   - M.S. Statistics from UW-Madison

2. Work history
   - Scientist at US Geological Survey
   - V.P. Data Engineering at Earthling Interactive
   - Senior Data Analyst at Zendesk

3. BigDataMadison Meetup
   Organizer, 2012 - present

4. BigDataWisconsin Conference
   Co-Organizer, 2016 - present

Provenance

- Development began in 2014
- Open-sourced by AirBnB in 2015
- Accepted into Apache Software Foundation (ASF) Incubator in 2016
- Current version 1.8.1
- 100% Python!

# Apache Airflow

- URL: https://github.com/apache/incubator-airflow
- 300+ committers
- 4000+ commits
- 6000+ stars

# Apache Airflow

- Task-based data orchestration platform
- Create, schedule and monitor workflows
- Workflows consist of DAGs of tasks
- DAGs control WHEN something happens, and under what conditions
- Tasks control WHAT happens

**DAGs** - Directed Acyclic Graphs
- a collection of all the tasks you want to run for a job, organized in a way that reflects their ordering, relationships and dependencies.

DAGs aren't concerned with what its constituent tasks do;
It's job is to make sure that whatever they do happens:
- at the right time,
- in the right order,
- with the right handling of any unexpected issues

# Apache Airflow

**Operators** - a single task in the workflow

Many types:
- SqlOperator
- PythonOperator
- BashOperator
- HTTPOperator
- AWS and GCS operators
- plenty more, including community supplied and custom

General Architecture – Major Components

- Auth - authentication protocols (password, Kerberos, LDAP, Google)
- Hooks – interface/connect to external platforms and databases
- Sensors - wait for files, change of state, db rows, etc.
- Executors - control scheduling/how the DAGs are run (local, sequential, Celery, Mesos are defaults available)

## General Architecture – Controlling the DAG

- Linear DAGs are the simplest variant
- Branching logic - using the BranchPythonOperator
- Sub-dags - useful for repeated patterns (where the subdag is returned from a function
- SLAs (Service Level Agreement) - time-based checks on task success

Running Airflow

- Command line interface - full suite of commands: DAG operation, scheduling
- GUI - Lots of functionality here (live demo coming up)
  - Scheduling
  - Dig into code
  - Administration
  - Activity history

LIVE DEMO!

zendesk

# Apache Airflow

- Populate BigQuery (GCS data warehouse) with information from multiple third party applications used throughout company
  - NetSuite
  - HEAP
  - Zuora
  - Salesforce
  - Bizible
  - MySQL
  - Hadoop

# Apache Airflow

- Use a custom executor for Kubernetes*, to scale up in a cloud environment
- Container spawner
- Uses Kubernetes and Docker** instances to spin up each task (sequentially or in parallel) in separate containers

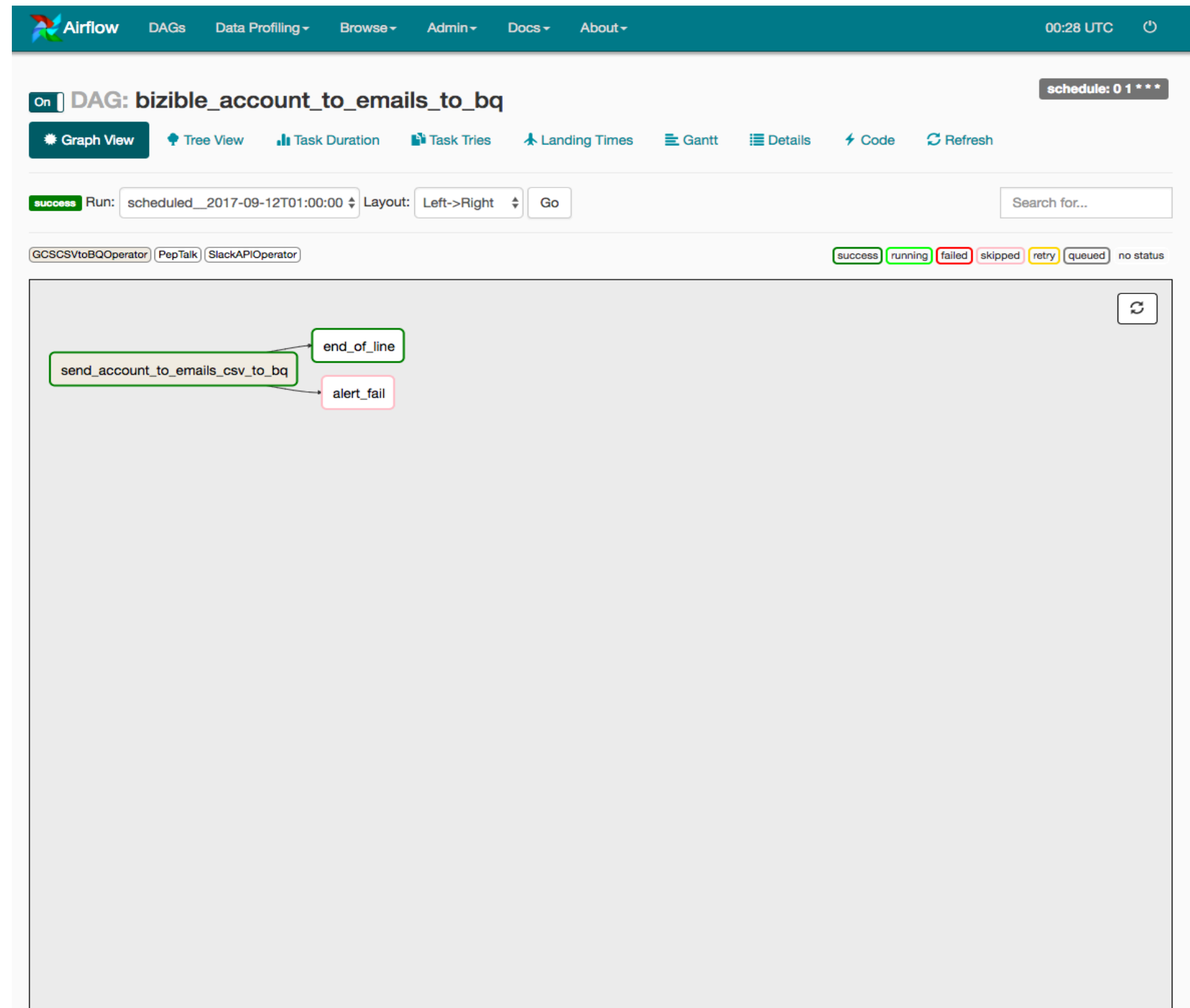* A system for management/deployment of containers in a cluster
** Container software

# Apache Airflow

- Use hooks to extract datafile to GCS buckets, employing a consistent nested folder construct
- Pick up the file and load the data to a stage table
  - Sensor detects new file(s)
  - Container spawner fires up new container to process the first task in the DAG (KubeExecutor).
  - Once finished spins down and launches new one (if necessary)
- Perform quality checks on the data, notify via Slack
- Move data to the production table using an SCD type 2 pattern
- Share data out between the separated datasets using views, to avoid exposing underlying source tables

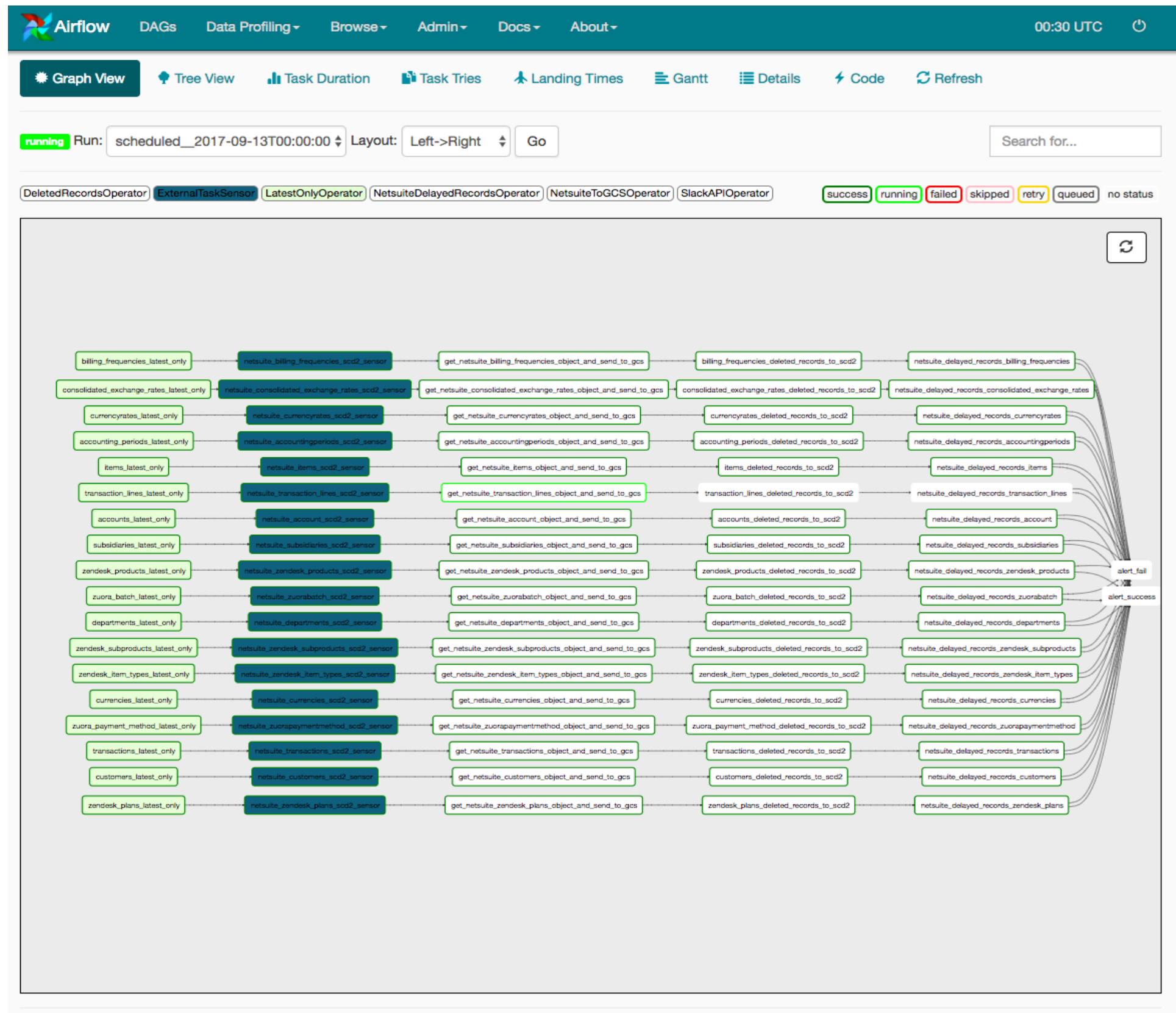# Apache Airflow

How we use it at Zendesk



Graph View simple DAG
• Failure vs Success

# Apache Airflow

## How we use it at Zendesk



Graph View – complex DAG
- A series of linear DAGs tied together and dependent on one another for success or failure

# Apache Airflow

- Documentation - http://pythonhosted.org/airflow

- Repo - https://github.com/apache/incubator-airflow

- Community: https://gitter.im/apache/incubator-airflow

## Questions

KEEP CALM AND ASK ANY QUESTION

Pitt Fagan

- http://www.pittfagan.com
- https://www.linkedin.com/in/pittfagan/
- https://www.meetup.com/BigDataMadison/

# Thanks for coming!