# Airflow

Arun

May 27, 2019

# Airflow

# What is Airflow?

# What is Airflow?

- It's a workflow manager

# What is Airflow?

- It's a workflow manager
- We can also call it a fancy Scheduler

# What is Airflow?

- It's a workflow manager
- We can also call it a fancy Scheduler
- It can be a easy backend job Scheduler in a distributed system
- Ref, this blog

# Other tools

# Other tools

- Oozie - XML code

# Other tools

- Oozie - XML code
- Pinball (pinterest) - no community interest

# Other tools

- Oozie - XML code
- Pinball (pinterest) - no community interest
- Luigi (Spotify) - Uses HDFS, no alerts or monitoring

# Other tools

- Oozie - XML code
- Pinball (pinterest) - no community interest
- Luigi (Spotify) - Uses HDFS, no alerts or monitoring
- Azkaban (LinkedIn) - Uses custom config file for workflow setups

# Airflow

# Airflow

- Python code base + for workflow definitions

# Airflow

- Python code base + for workflow definitions
- Trigger rules

# Airflow

- Python code base $+$ for workflow definitions
- Trigger rules
- Xcoms

# Airflow

- Python code base $+$ for workflow definitions
- Trigger rules
- Xcoms
- Cool UI & CLI

# Airflow

- Python code base + for workflow definitions
- Trigger rules
- Xcoms
- Cool UI & CLI
- Queues & Pools

# Airflow

- Python code base + for workflow definitions
- Trigger rules
- Xcoms
- Cool UI & CLI
- Queues & Pools
- Zombie cleanup

# Airflow
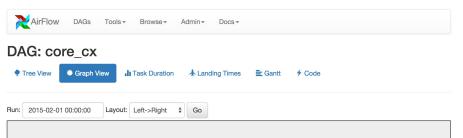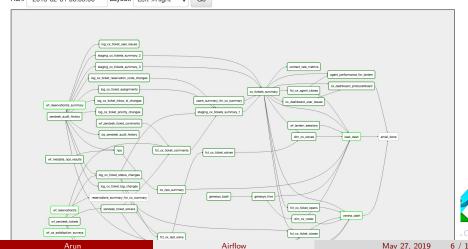
- Python code base $+$ for workflow definitions
- Trigger rules
- Xcoms
- Cool UI & CLI
- Queues & Pools
- Zombie cleanup
- Large community

# DAG: core_cx

♠ Tree View    ✸ Graph View    ▮▮ Task Duration    ✈ Landing Times    ≡ Gantt    ⚡ Code

Run: 2015-02-01 00:00:00    Layout: Left->Right ⬍    Go

# Airflow - Executors

- Sequential
- Local (parallel)
- Celery
- Mesos
- Kubernetes (still in alpha)

# Our cluster

- Single node (8 core, 20 GB node)
- Local Executor (with 24 parallel schedulable tasks)
- 3000+ tasks triggred a day
- Code auto deployed from master to prod

# Tasks

- BQ job
- Run spark jobs daily for ML
- All batched jobs in Google Dataflow
- Scheduling backend jobs
- All errors to SLACK
- Critical DAG failures will invoke Victorops issue