NIKE DTC RI

# Basics of Airflow

The WHAT, WHY and HOW to use Airflows

~Pratik Saurav

# The Problem

## Companies are growing to have a complex network of processes having intricate dependencies

- Analytics ETL processes and batch processing are the critical items to be catered to
- Tons of time was spent writing jobs, monitoring and troubleshooting issues

The Classical Approach of writing batch jobs on your cluster and schedule them using a cron may not suffice.

**WHY?**

- Data lineage (the complete data cycle) might be opaque in nature
- Duplication of codes/logic
- Maintenance time grows exponentially as the ecosystem grows
- Signal to noise ratio may get out of control

# The Genesis

## Airbnb was outgrowing its workflows methodology

## 200k +Hive tables

**> 10,000**
Superset charts and dashboards

**> 6,000**
Experiments and metrics

**> 6,000**
Tableau workbooks and charts

**> 1,500**
Knowledge posts

**> 3,500 Airbnb employees**

# The Genesis

## Airbnb was outgrowing its workflows methodology

- They wanted to create a platform with dynamic pipeline generation, programmatic environment and a platform that people love working in.
- Hence they took it as an opportunity to build and share Airflow

**Brethren of airflow?**

- Facebook uses SWARM
- Yahoo uses (Used to) OOZIE
- Cron we are well aware of
- Ubisoft uses informatica based platforms

# The Genesis

**Airflow is an open source platform to programmatically author, schedule and monitor workflows.**

Maxime Beauchem (Airbnb) created the whole thing. Wanting to create a platform which much wider scopes unlike the traditional job schedulers, Airbnb tried to identify the pros, cons and features from the best job schedulers in business.

**Basic Features:**

- Basic dependency management
- Status clarity
- Scheduling
- Easy access to logs and Notification features
- Job prioritization features and Parametrized retries
- Has great capabilities to scale out
- Supports complex dependency rules like branching, joining, subworkflows
- You can define your own constructs, callbacks
- Easy state altering like : DAG surgeries, ownership and versioning, data profiles, centralized connection management, Plugins (variety of plugins with connectivity to oracle, Hadoop, query engines etc.)

# The Genesis

**Airflow is an open source platform to programmatically author, schedule and monitor workflows.**

[Maxime Beauchem](#) (Airbnb) created the whole thing. Wanting to create a platform which much wider scopes unlike the traditional job schedulers, Airbnb tried to identify the pros, cons and features from the best job schedulers in business.

**Basic Features:**

• Basic dependency management

• Status clarity and scheduling

• Easy access to logs and Notification features

• Job prioritization features and Parametrized retries

• Has great capabilities to scale out and supports complex dependency rules like branching, joining, subworkflows

• You can define your own constructs, callbacks

• Easy state altering like : DAG surgeries, ownership and versioning, data profiles, centralized connection management

# The 'how to_' guide

## An airflow has four essential blocks for execution

- **DAG**: a description of the order in which work should take place
- **Operator**: a class that acts as a template for carrying out some work
- **Task**: a parameterized instance of an operator
- **Task Instance**: a task that
    - Has been assigned to a DAG and
    - Has a state associated with a specific run of the DAG

# The 'how to_' guide

## An airflow is designed as a Directed Acyclic Graph

- When authoring a workflow, you should think how it could be divided into tasks which can be executed independently.
- You can then merge these tasks into a logical whole by combining them into a graph.



- The shape of the graph decides the overall logic of our workflow.
- Directed Acyclic Graphs (DAGs) are trees of nodes that Airflow's workers will traverse. Each node in the graph can be thought of as a steps and the group of steps make up the overall job.

# The 'how to_' guide

## An airflow is designed as a Directed Acyclic Graph

**Why are they Acyclic?**

The graphs are acyclic so you cannot jump to a previous task in the graph but you can jump to another task that hasn't yet run in the current job.

This means individual steps can be retried repeatedly if they have a failure without having to restart the whole job.

**What is a DAG?**

In Airflow, a DAG – or a Directed Acyclic Graph – is a collection of all the tasks you want to run, organized in a way that reflects their relationships and dependencies. DAGs are defined in standard Python files that are placed in Airflow's DAG_FOLDER .

# Anatomy of a DAG

## A DAG is defined by a DAG definition



```
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
from datetime import datetime, timedelta
```
Importing the libraries which are required

```
default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2015, 6, 1),
    'email': ['airflow@airflow.com'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5),
    # 'queue': 'bash_queue',
    # 'pool': 'backfill',
    # 'priority_weight': 10,
    # 'end_date': datetime(2016, 1, 1),
}
```
**Default arguments**
To define default parameters for creating tasks

```
dag = DAG('tutorial', default_args=default_args)
```
**DAG Instantiation**
This is a DAG object to nest our DAGS. You pass a (Dag ID, default args which you just created, interval (1day) here)

```
# t1, t2 and t3 are examples of tasks created by instantiating operators
t1 = BashOperator(
    task_id='print_date',
    bash_command='date',
    dag=dag)

t2 = BashOperator(
    task_id='sleep',
    bash_command='sleep 5',
    retries=3,
    dag=dag)
```
**Tasks**
To define default parameters for creating tasks
Here 't1' or 't2' gives an independent task names to differentiate
Bash_command is operator specific
You can see that we can override parameters defined in the default_args
Eg. In t2, retries is overridden by 3 which was predefined as 1

```
templated_command = """
    {% for i in range(5) %}
        echo "{{ ds }}"
        echo "{{ macros.ds_add(ds, 7)}}"
        echo "{{ params.my_param }}"
    {% endfor %}
"""

t3 = BashOperator(
    task_id='templated',
    bash_command=templated_command,
    params={'my_param': 'Parameter I passed in'},
    dag=dag)
```
**Templating**
Airflow leverages the power of Jinja Templating and provides the pipeline author with a set of built-in parameters and macros.
Airflow also provides hooks for the pipeline author to define their own parameters, macros and templates. template variable: {{ ds }}.
See that here logic is defined in %%, {{ds}}, a parameter being passed to a function etc.

```
t2.set_upstream(t1)
t3.set_upstream(t1)
```
**Dependencies**
This means that t2 will depend on t1

# Anatomy of a DAG

## A DAG is defined by a DAG definition

- A DAG definition is a python script wherein you define DAG's structure as a code.
- Airflow will execute the code in each file to dynamically build the DAG objects.
- The point to be kept in mind that it's actually just a configuration file you cannot do any data processing here at all.
- The scripts purpose is to define a DAG object. It needs to evaluate quickly (seconds, not minutes) since the scheduler will execute it periodically to reflect the changes if any

The following will come in handy while reading a DAG definition:

- [list of operators](#)
- [Standard airflow operators](#)
- [Airflow tutorial](#)

# Creating a DAG

**There are three main terminals in airflow**

- [The Airflow Terminal](#)
- [Celery Flower Terminal](#)
- [Airflow Scheduler TTY Shell](#)

# Creating a DAG

### The Airflow terminal is the basic UI terminal

[The Airflow Terminal](#)

- It's basically a GUI which has the list of all DAGs, their statuses, definitions which are very much intuitive
- Your front screen will appear pretty much like this. Where you'll see the list of DAGs, their owners and statuses
- When you click on a DAG the airflow will open up it's entire facebook profile in front of you XP (The name, status, profile picture(Acyclic Graphs), Events, Details, The codes associated) Everything!

# Creating a DAG

The flower terminal is used for performance monitoring of workers

# Creating a DAG

The scheduler tty terminal is used for deploying jobs to the workers