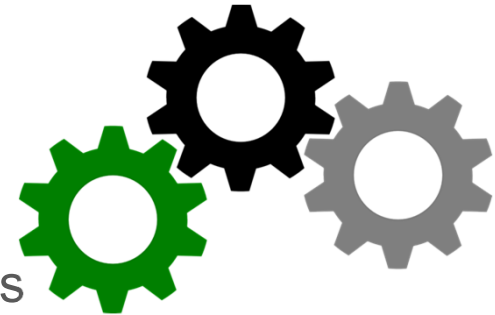# János Dani

CEU / OSA helpdesk

GitHub: /danijanos

Instagram: /dan1jan1

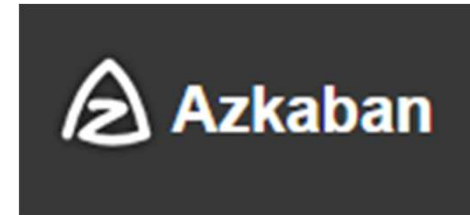# Using Airflow as an Orchestrator for microservices

Archiving at OSA

# What are microservices?

- Architectural style
- Structures an application as a collection of services
- Highly maintainable and testable
- Loosely coupled
- Independently deployable
- Communication through transactions
- Could be scaled separately and independently

# Workflow tools

- There is quite a lot...
- They are for various purposes
- Why is one better or favourable than the other?

apache / airflow

Used by  576    Watch  660    Star  14.8k    Fork  5.6k

<> Code     Pull requests 187     Actions     Projects 0     Security     Insights

Apache Airflow - A platform to programmatically author, schedule, and monitor workflows   https://airflow.apache.org/

airflow    apache    apache-airflow    python    scheduler    workflow

7,465 commits    9 branches    0 packages    124 releases    991 contributors    Apache-2.0

- ○ Email notification
  - ○ Charts
  - ○ Error checking
- ● Has a good documentation / community behind
- ● [Open source]

Apache
Airflow

# About Airflow

Started at Airbnb in October 2014

Written in Python

Becoming an Apache Incubator project in March 2016

Top-Level Software Foundation project in January 2019.

## The Apache Software Foundation Blog

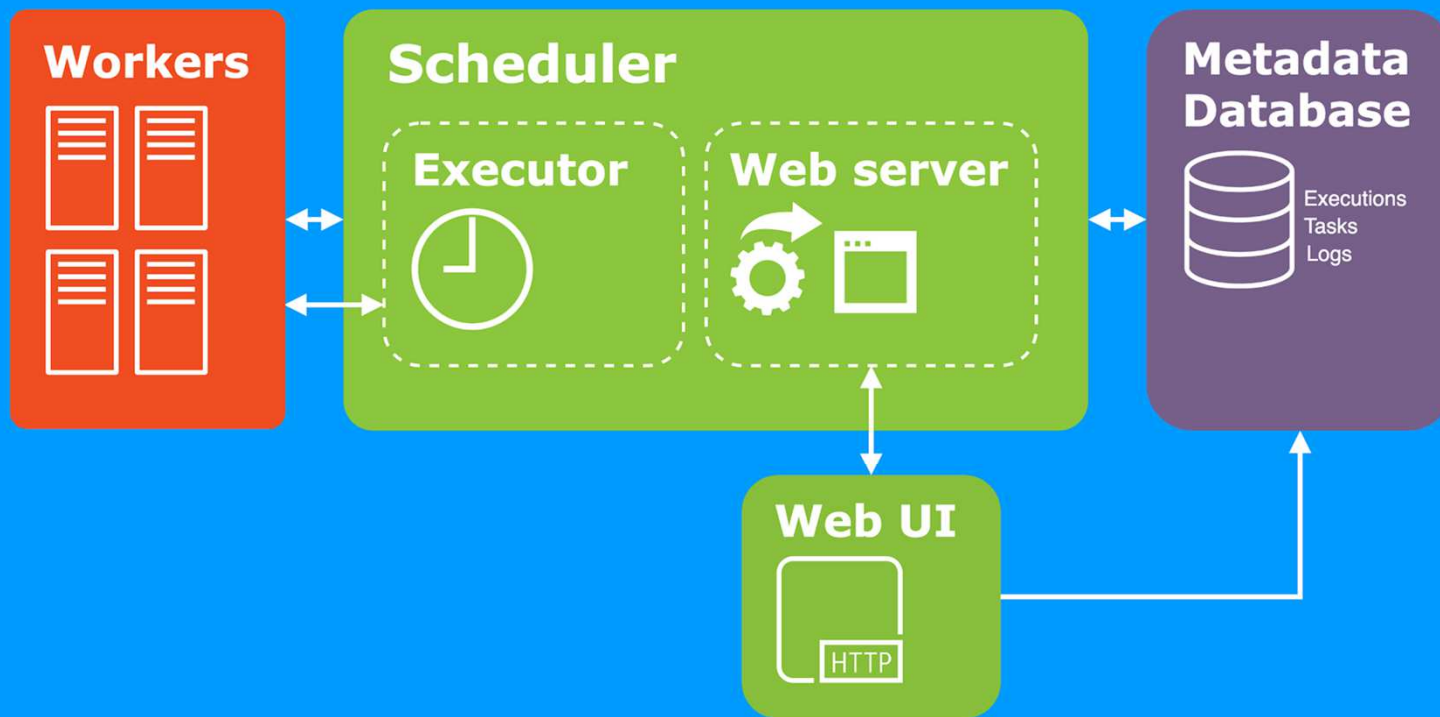« Success at Apache:... | Main | The Apache News... »

TUESDAY JANUARY 08, 2019

**The Apache Software Foundation Announces Apache® Airflow™ as a Top-Level Project**

*Open Source Big Data workflow management system in use at Adobe, Airbnb, Etsy, Google, ING, Lyft, PayPal, Reddit, Square, Twitter, and United Airlines, among others.*

- https://en.wikipedia.org/wiki/Apache_Airflow
- https://airflow.apache.org/

# The scheduler

- Executes (triggers) tasks on an array of workers
- Monitors all tasks
- Executor types:
  - Sequential executor
  - Local
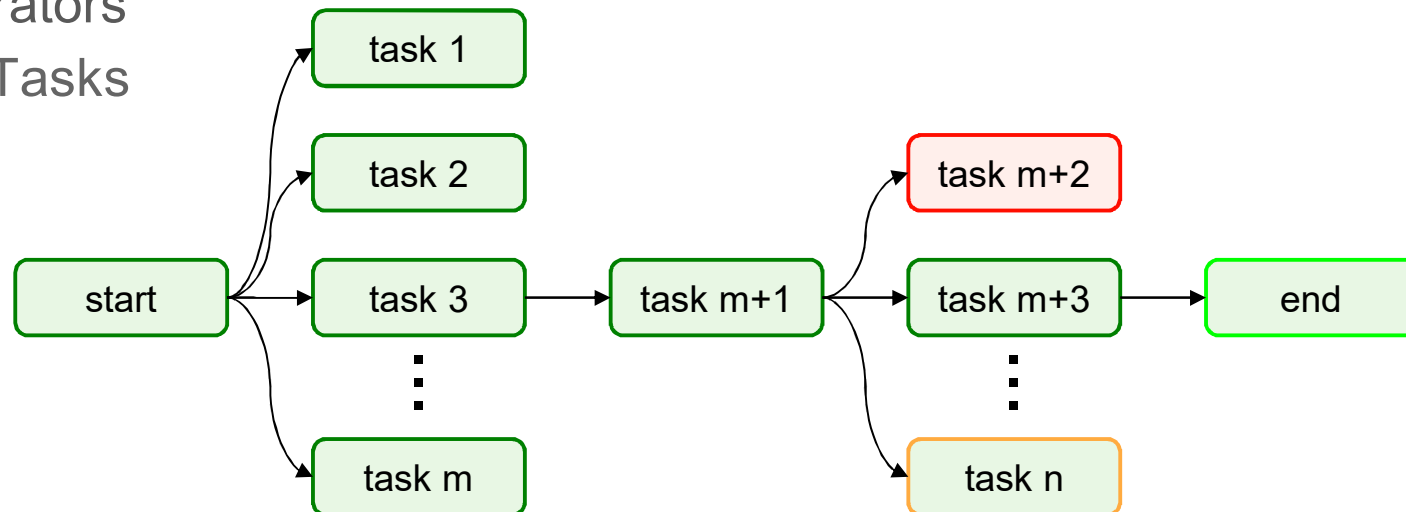  - Celery (to scale tasks on different worker nodes)

source:
- https://airflow.apache.org/docs/stable/scheduler.html

# Basic conceptual building blocks
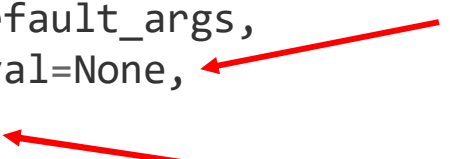
- Dags (Pipelines)
- Operators
  - Tasks

# Dag

- Holds a series of tasks connected with dependencies
- Made for avoiding cyclical dependencies between tasks
- Dags are identified by their IDs

```python
default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2018, 1, 1),
    'email': ['bonej@ceu.edu', 'danij@ceu.edu'],
    'email_on_failure': True,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5)
}
```
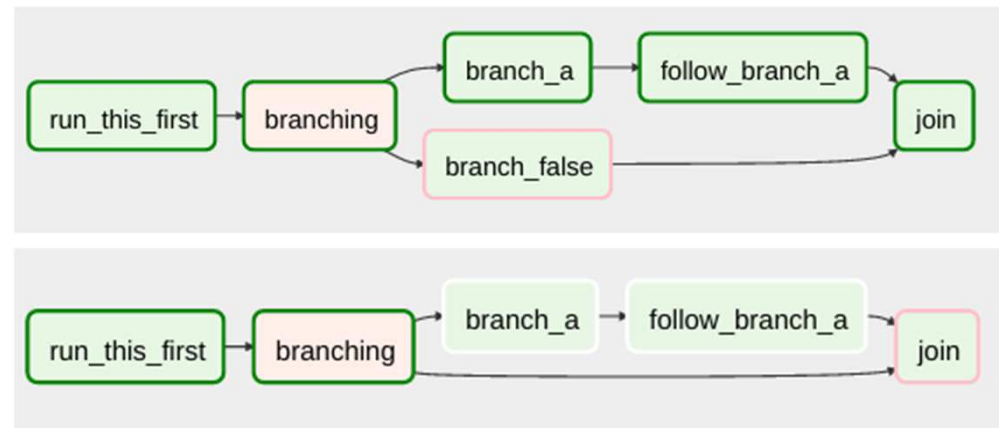
```python
# DAGS
osa_av_workflow = DAG(
    dag_id='osa-av-workflow',
    description='Main DAG for the AV preservation workflow',
    default_args=default_args,
    schedule_interval=None,
    catchup=False)
```

8/13

# Operators

```
create_master_checksums = PythonOperator(
    task_id='create_master_checksums',
    python_callable=create_checksums,
    dag=osa_av_workflow,
    op_kwargs={
        'directory': 'Preservation',
        'file_extension': MASTER_FILE_EXTENSION
    })
```
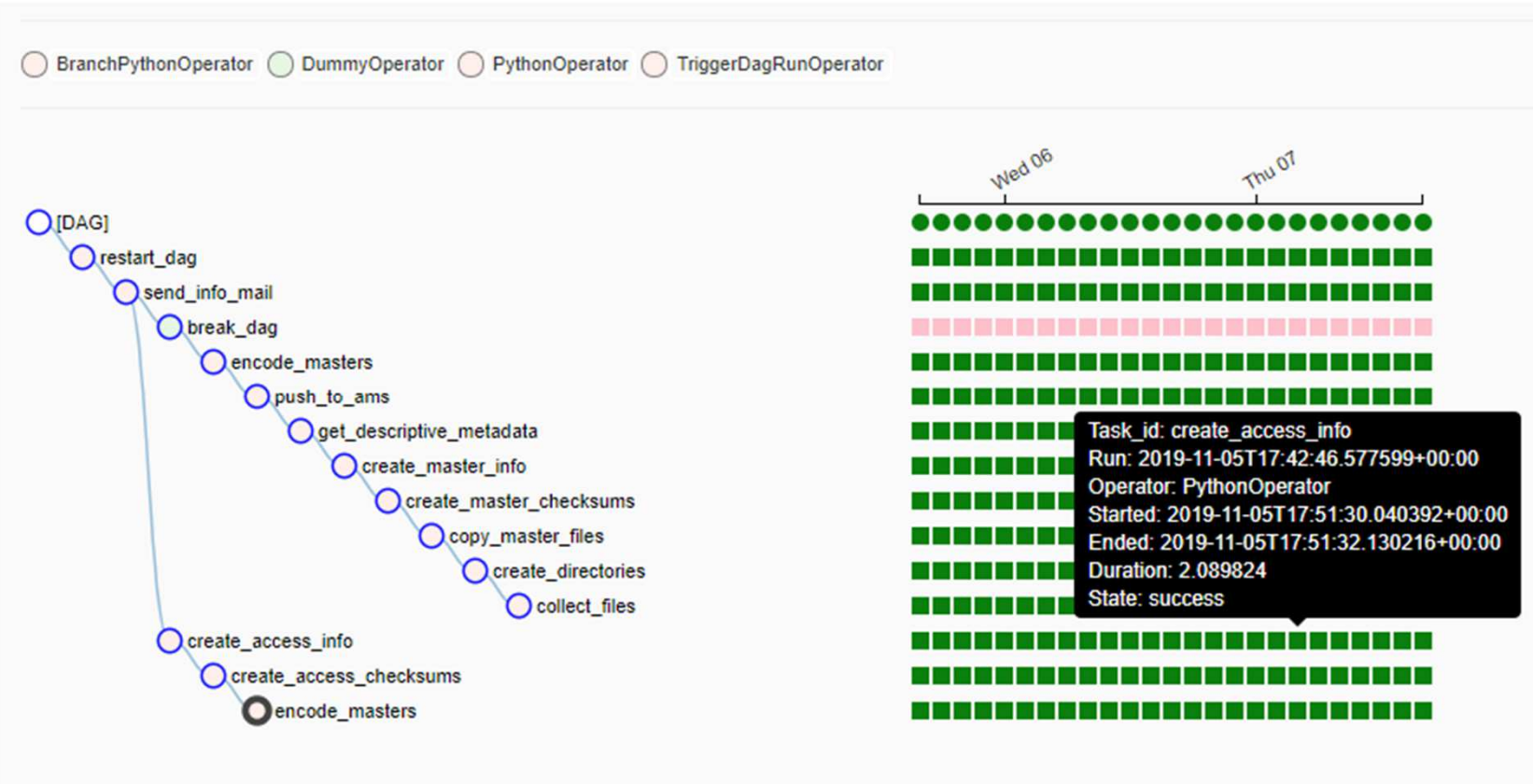
- BashOperator - executes a bash command
- PythonOperator - calls a Python function
- EmailOperator - sends an email
- SimpleHttpOperator - sends an HTTP request
- BranchOperator
- …
- Write your our own operators!

# Tasks

- Tasks are instantiated operators.
- Connected with dependencies
- Can be parallelized

```
op1.dag = dag
op1.set_downstream(op2)
# is the same as:
dag >> op1 >> op2

op1 >> op2
# is the same as:
op1.set_downstream(op2)

op2 << op1
# is the same as:
op2.set_upstream(op1)
```

# Statuses

# Airflow Web UI - Gantt chart

# Airflow Web UI - Logs

# Links

- [https://airflow.apache.org/](https://airflow.apache.org/)
- [https://hub.docker.com/r/apache/airflow](https://hub.docker.com/r/apache/airflow)
- [http://michal.karzynski.pl/blog/2017/03/19/developing-workflows-with-apache-airflow/](http://michal.karzynski.pl/blog/2017/03/19/developing-workflows-with-apache-airflow/)