

# Why Apache Airflow Sucks

---

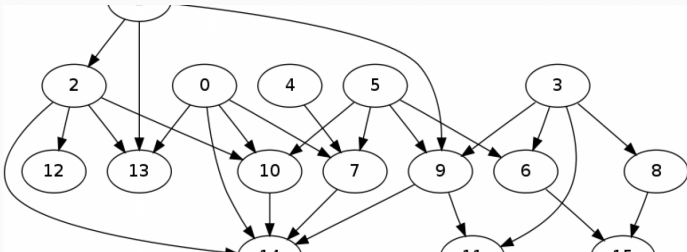
Слава Могилевский

5 октября 2017 г.

Provectus

# Введение: для чего нужен Airflow

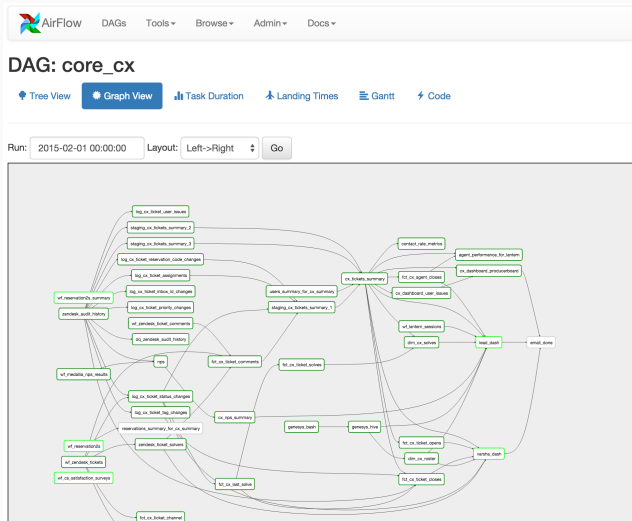
- запустить задачи (hive-запрос, spark, запрос в mist?) по расписанию
- у задачи есть зависимости от других (DAG)
- мониторинг, логи, и т.д



## Airflow don't suck sometimes

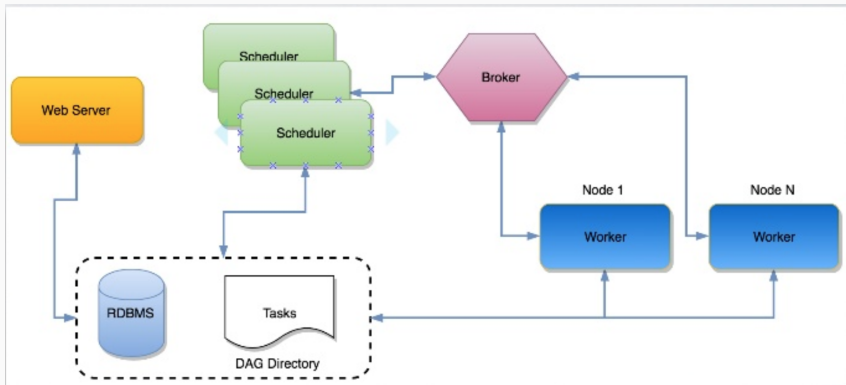
- Динамика: dynamic pipeline generation
- Гибкость: parameterizing scripts by Jinja templating engine
- Масштабируемость: ready to scale to infinity

# Airflow don't suck sometimes



Удобный UI

# Архитектура Airflow



Обновление DAG - это боль:

- Нужно каждый раз рестартить: UI, scheduler
- Обновленный DagBag как-то нужно расшарить между UI, scheduler и Worker'ами
  - NFS - привет, race-conditions
  - Делаем копии DagBag везде, куда нужно - проблемы с синхронизации
- Удаление тасок - часто их еще нужно удалять ручками с БД

Для распределенных worker'ов используется только Celery

А Celery - это тоже боль

Невозможно понять, что происходит с воркерами:

- какая задача выполняется
- почему в airflow висит задача, а в celery все воркеры свободны!?

Логи - это боль

- Логи хранятся у worker'а - они пропадают вместе с ним
- Есть опция хранить логи в S3 - не всегда это нужно
- Хочется иметь живую консоль, как у Jenkins



Системные зависимости в worker'е - это боль

Для каждого нового типа тасок (hive query, spark query) необходимо:

- пересобрать образ для воркера для новой зависимости
- передеплоить новый образ на все воркеры
- ничего при этом не ломать

Workflow as a Code vs. Workflow as a Infrastructure

Все должно быть взаимозаменяемым (плагины)

Задача = собираем зависимости + go for it!

<https://github.com/meirwah/awesome-workflow-engines> - список движков для построения и запуска workflow ( Jenkins'а в нем нет ;) )

- Airflow
- Luigi
- Pinball
- ...

Внимание!!!



Спасибо за внимание