



ОНЛАЙН-ОБРАЗОВАНИЕ

# Меня хорошо слышно && видно?



Напишите в чат, если есть проблемы!

Ставьте  если все хорошо

# Оркестраторы

Oozie & Airflow



После этого занятия вы будете знать

- Какие задачи решают оркестраторы
- Что такое Oozie и Airflow

После этого занятия вы сможете

- Написать свой пайплайн в Airflow

# 01

**Когда crop'a  
недостаточно**

**cron** — классический демон (компьютерная программа в системах класса UNIX), использующийся для периодического выполнения заданий в определённое время. Регулярные действия описываются инструкциями, помещёнными в файлы crontab и в специальные каталоги.

Таким образом, cron помогает

- автоматизировать рутинные задачи
- выполнять задачи периодически без ручных запусков

## Проблемы применения cron'а в больших проектах

- трудно мониторить
- непонятно, где лежит конфигурация
- нет перезапусков
- нельзя задавать связи между задачами

```
#30 2 * * * spark-submit --master yarn --queue root.daily --executor-cores=8 --executor-memory=16G --class BanReport /mnt/h
#30 4 * * * spark-submit --master yarn --queue root.daily --executor-cores=4 --executor-memory=32G --class BanGenerator /mnt
30 5 * * * spark2-submit --master yarn --queue root.daily --executor-cores=4 --executor-memory=8G --class etl.RequestOtahot
#30 6 * * * spark-submit --master yarn --queue root.daily --executor-cores=8 --executor-memory=16G --class etl.AvailScore /
#10 * * * * spark-submit --master yarn --executor-cores=4 --executor-memory=8G --class etl.DAPIToParquet /mnt/hadoop/spark-
#31 * * * * spark-submit --master yarn --executor-cores=4 --executor-memory=8G --class DiscrepancyHunter /mnt/hadoop/spark-
#33 * * * * spark-submit --master yarn --executor-cores=8 --executor-memory=16G --class DiscrepancyInStep /mnt/hadoop/spark
20 * * * * . /etc/profile; /mnt/hdfs/tools/spark-jars/exchange.py >> /root/exchange.log 2>&1
#11 * * * * spark2-submit --master yarn --queue root.hourly --executor-cores=4 --executor-memory=8G --class etl.RawDataConv
#26 * * * * spark-submit --master yarn --queue root.hourly --executor-cores=4 --executor-memory=8G --class etl.MinRates /mnt
#32 * * * * spark2-submit --master yarn --queue root.hourly --executor-cores=4 --executor-memory=8G --class etl.Discrepancy
32 * * * * spark-submit --master yarn --queue root.hourly --executor-cores=4 --executor-memory=8G --class etl.DiscrepancyHu
#40 * * * * spark-submit --master yarn --queue root.hourly --executor-cores=4 --executor-memory=8G --class etl.RAIndex /mnt
```

- Визуальное представление задач с управлением через GUI
- Единая система для всех задач во всех сервисах на всех серверах
- Обеспечение зависимостей между задачами
- Перезапуск задач со сложной логикой
- Алертинг

- Apache Oozie
- Apache Airflow
- Luigi
- Azkaban

# 02

## Обзор Apache Oozie

- Оркестратор, рассчитанный на работу с задачами кластера Hadoop
- Может выполнять задачи на основе
  - Spark
  - Hive
  - Pig
  - Sqoop
  - Shell

- Язык описания задач - XML
- Есть несколько верхнеуровневых объектов
  - **Workflow** - граф задач
  - **Coordinator** - Workflow + расписание
  - **Bundle** - набор первых и вторых

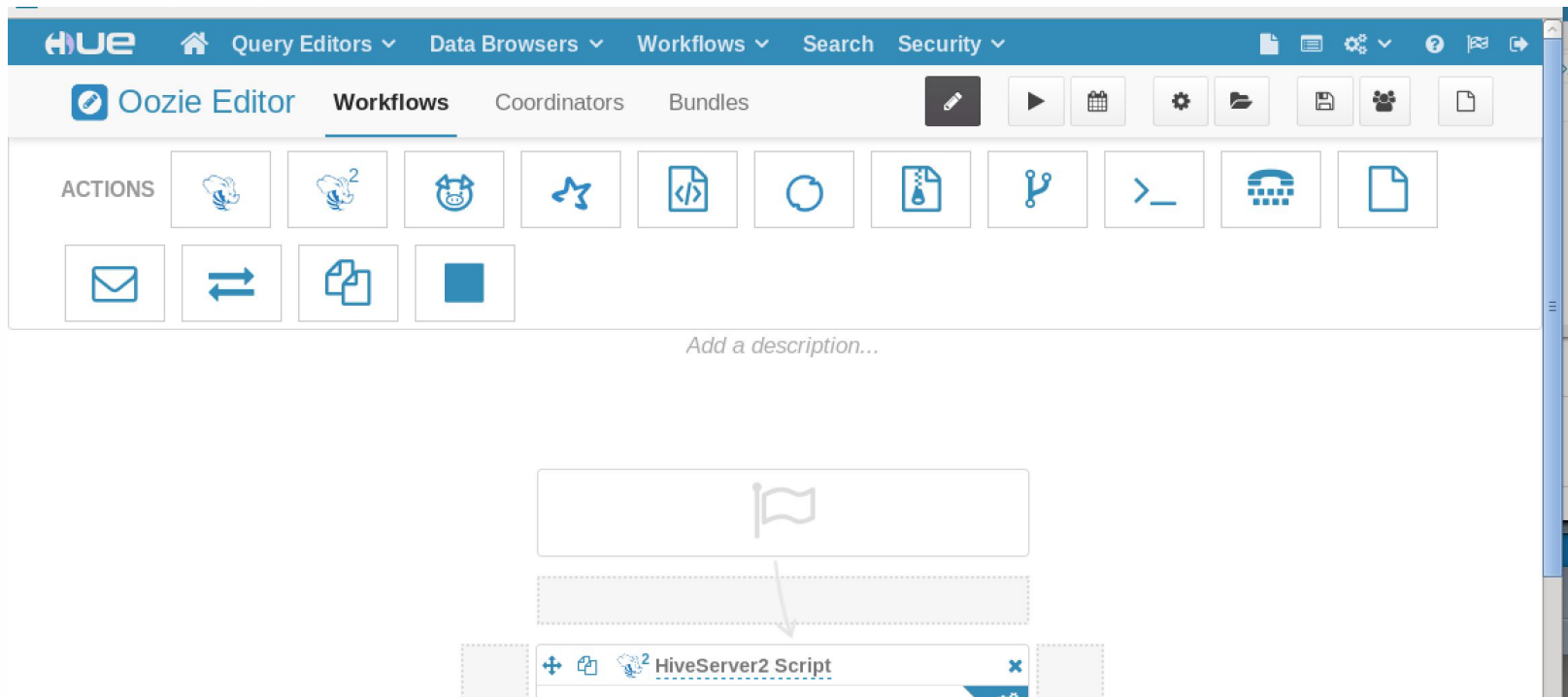
```
<workflow-app xmlns = "uri:oozie:workflow:0.4" name = "test-workflow">
  <start to = "run_hive_script" />
  <action name = "run_hive_script">
    <hive xmlns = "uri:oozie:hive-action:0.4">
      <job-tracker>xyz.com:8088</job-tracker>
      <name-node>hdfs://rootname</name-node>
      <script>hdfs_path_of_script/external.hive</script>
    </hive>
    <ok to = "end" />
    <error to = "kill_job" />
  </action>
  <kill name = "kill_job">
    <message>Job failed</message>
  </kill>
  <end name = "end" />
</workflow-app>
```

Запуск wf из командной строки, где *host\_name:8080* -  
адрес веб-морды oozie

```
oozie job --oozie http://host_name:8080/oozie -D
```

```
oozie.wf.application.path=hdfs:///path/to/workflow.xml
```

- В реальном мире чаще всего не XML, а HUE



# 03

## Обзор Apache Airflow

**Apache Airflow** — платформа для автоматического управления задачами, их расписанием и мониторингом.

Изначально был разработан в AirBNB.

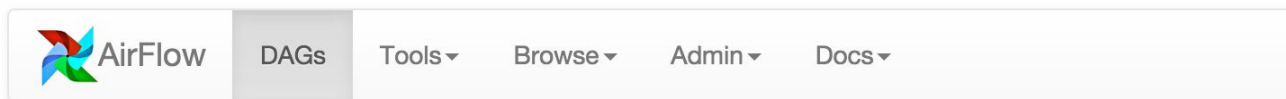
The screenshot shows the Apache Airflow web interface. At the top, there is a navigation bar with the AirFlow logo and several menu items: DAGs, Tools, Browse, Admin, and Docs. The Admin menu is open, showing options: Configuration, Connections (highlighted), Users, and Reload DAGs. Below the navigation bar, there is a table with 4 connections. The table has columns for Conn Id and type. Each row includes a checkbox, edit, and delete icons.

		Conn Id	type
<input type="checkbox"/>		local_mysql	mysql
<input type="checkbox"/>		mysql_default	mysql
<input type="checkbox"/>		presto_default	presto
<input type="checkbox"/>		hive_default	hive

- задачи описываются на Python
- поддерживает сложные процессы из нескольких задач с зависимостями, ветвлением, условиями
- удобный интерфейс
- может масштабироваться с использованием Celery
- легко расширяется под конкретные задачи

- DAG - граф вычислений
- Operator - тип задачи
- Task - задача, граф состоит из задач
- DAG run - запуск графа для конкретного периода
- Task instance - запуск задачи для конкретного периода

В Airflow задачи собираются в длинные цепочки, называемые DAG.

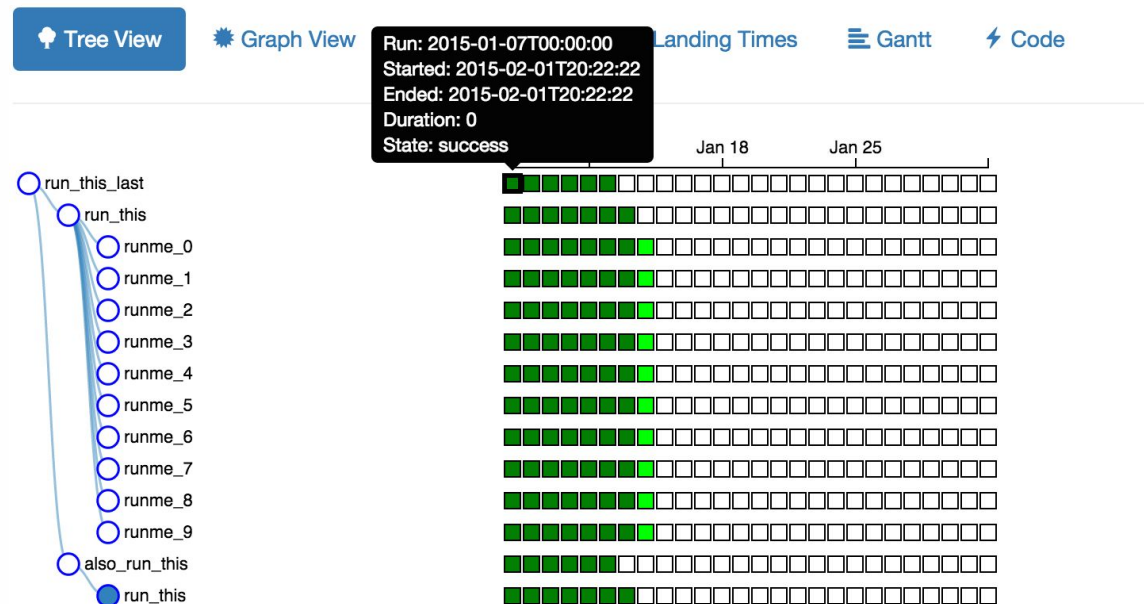


## DAGs

DAG	Filepath	Owner	Task by State	Links
<a href="#">example1</a>	<a href="#">example_dags/example1.py</a>	airflow	<span>80</span> <span>1</span> <span>0</span>	
<a href="#">example2</a>	<a href="#">example_dags/example2.py</a>	airflow	<span>128</span> <span>10</span> <span>0</span>	
<a href="#">example3</a>	<a href="#">example_dags/example3.py</a>	airflow	<span>138</span> <span>5</span> <span>0</span>	

DAG (Directed Acyclic Graph) – это набор задач, которые вы хотите запускать, организованный таким образом, чтобы отражать их связи и зависимости.

## DAG: example2



```
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
from datetime import datetime, timedelta

default_args = {
    # ...
}

dag = DAG('tutorial', default_args=default_args,
          schedule_interval=timedelta(days=1))

t1 = BashOperator(task_id='print_date', bash_command='date', dag=dag)

t2 = BashOperator(
    task_id='sleep',
    bash_command='sleep 5',
    retries=3,
    dag=dag)

t1 >> t2
```

```
default_args = {  
    'start_date': datetime(2016, 1, 1),  
    'owner': 'Airflow'  
}  
  
dag = DAG('my_dag', default_args=default_args)  
op = DummyOperator(task_id='dummy', dag=dag)  
print(op.owner) # Airflow
```

- По расписанию
- Через веб-интерфейс
- Через API
- По сенсору

- BashOperator
- PythonOperator
- SimpleHttpOperator
- PostgresOperator
- HiveOperator
- DockerOperator
- SSHOperator
- SlackAPIOperator

Сенсоры позволяют выполнять задачи по наступлению определенного события:

- Появление файла на FTP
- Выполнение SQL-запроса
- Успешный ответ по HTTP

и т. д.

 [external\\_task\\_sensor.py](#)

 [hdfs\\_sensor.py](#)

 [hive\\_partition\\_sensor.py](#)

 [http\\_sensor.py](#)

 [metastore\\_partition\\_sensor.py](#)

 [named\\_hive\\_partition\\_sensor.py](#)

 [s3\\_key\\_sensor.py](#)

 [s3\\_prefix\\_sensor.py](#)

 [sql\\_sensor.py](#)

 [time\\_delta\\_sensor.py](#)

 [time\\_sensor.py](#)

 [web\\_hdfs\\_sensor.py](#)

Одна из главных полезных фишек - зависимость между задачами внутри DAG

```
dag = DAG('my_dag', start_date=datetime(2016, 1, 1))
```

```
op1 = DummyOperator(task_id='op1', dag=dag)
```

```
op2 = DummyOperator(task_id='op2', dag=dag)
```

```
op1.set_upstream(op2)
```

ИЛИ

```
op1 >> op2
```

Hooks и operators выглядят похожими объектами, но

- hook - отвечает за техническую часть (запрос к базе, выполнение команды, установка соединения)
- operator - за бизнес-логику (использование результатов запроса, составление команды)

Операторы используют хуки под капотом. Один оператор может использовать несколько хуков.

[Пример](#)

Пулы позволяют ограничить число исполняющихся задач одного типа. Где это может быть полезно:

- Использование общего ресурса
- Препятствие “накоплению” задач
- Управление приоритетами

Этот механизм позволяет безопасно хранить и переиспользовать общие настройки для подключения к сервисам (Main > Admin > Connections)

```
task1 = SSHOperator(  
    task_id="ssh1",  
    command="echo -n \"airflow works!\"",  
    ssh_conn_id="ssh_connection1",  
    dag=dag)
```

На уровне сервиса есть возможность задавать “глобальные” переменные (Admin -> Variables)

```
from airflow.models import Variable
```

```
foo = Variable.get("foo")
```

```
bar = Variable.get("bar", deserialize_json=True)
```

```
baz = Variable.get("baz", default_var=None)
```

Для общения между задачами в рамках одного DAG run

есть механизм

```
# inside a PythonOperator called 'pushing_task'
```

```
def push_function():
```

```
    return value
```

```
# inside another PythonOperator where provide_context=True
```

```
def pull_function(**context):
```

```
    value =
```

```
context['task_instance'].xcom_pull(task_ids='pushing_task')
```

Для заполнения “истории” есть специальная команда

**backfill**, которая вызывается консольной командой

```
airflow backfill -s 2019-08-01 -e 2019-08-10 example_dag_id
```

Её можно использовать для того, чтобы очистить и запустить заново выполненные задачи

**04**

**Подведем итоги**

После этого занятия вы будете знать

- Какие задачи решают оркестраторы
- Что такое Oozie и Airflow

После этого занятия вы сможете

- Написать свой пайплайн в Airflow

Заполните, пожалуйста, опрос в личном кабинете!

- [Документация Airflow](#)
- [Docker-образ для быстрого запуска](#)
- [Сравнение оркестраторов](#)



**Егор Матешук**

[egor@mateshuk.com](mailto:egor@mateshuk.com)

**Спасибо  
за внимание!**

