



Apache Airflow: семь раз примерь, один раз зашедуль

Петр Ермаков
Lamoda / DataGym

Agenda

- Мотивация
- Альтернативы
- Airflow
- Best Practices

Мотивация

5	9:00am Start work 6:00pm Finish work	6
work	work	work

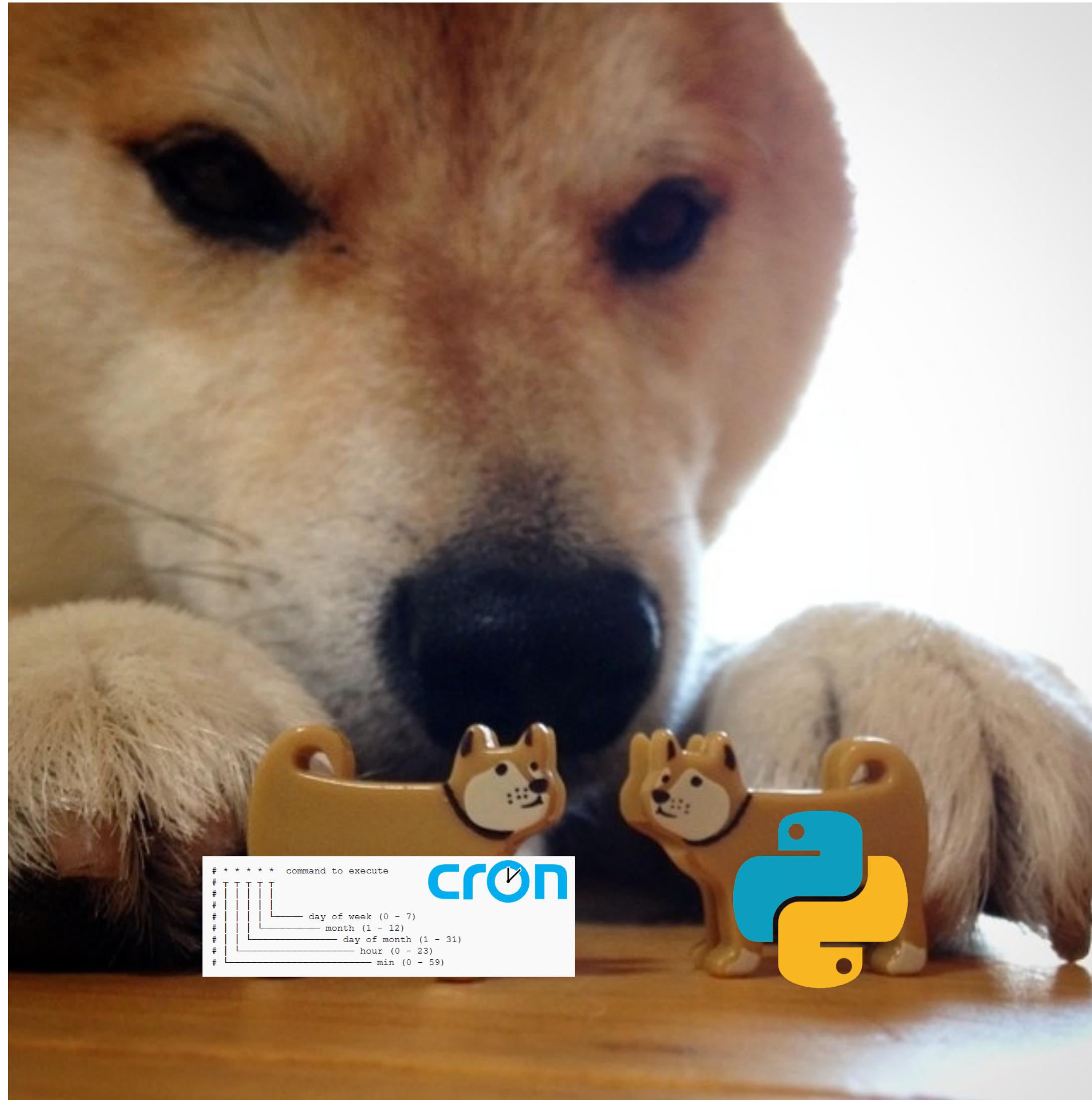
Мотивация



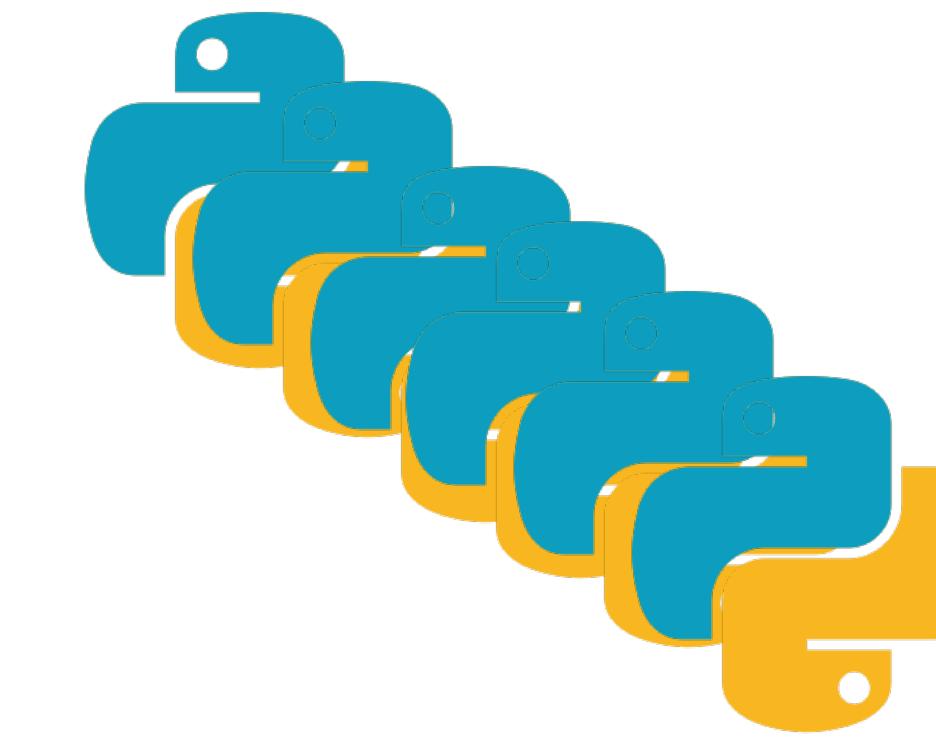
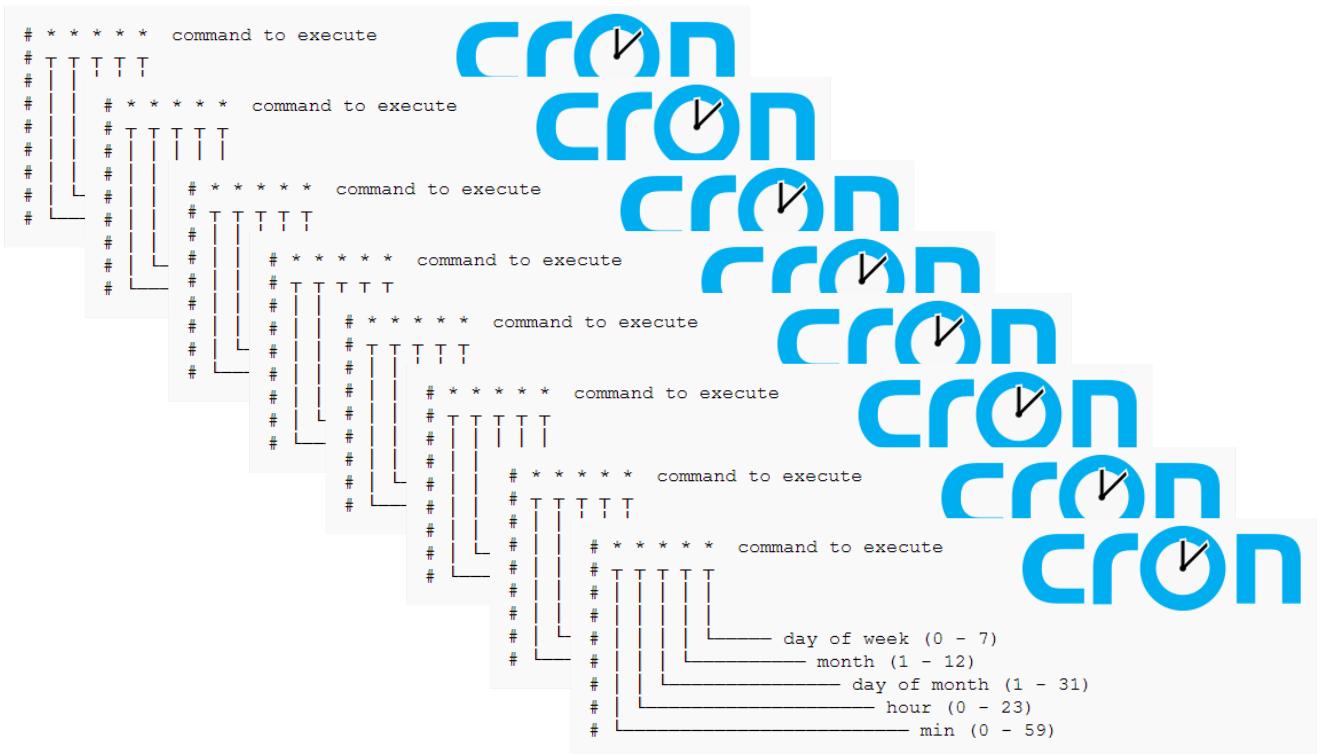
Мотивация

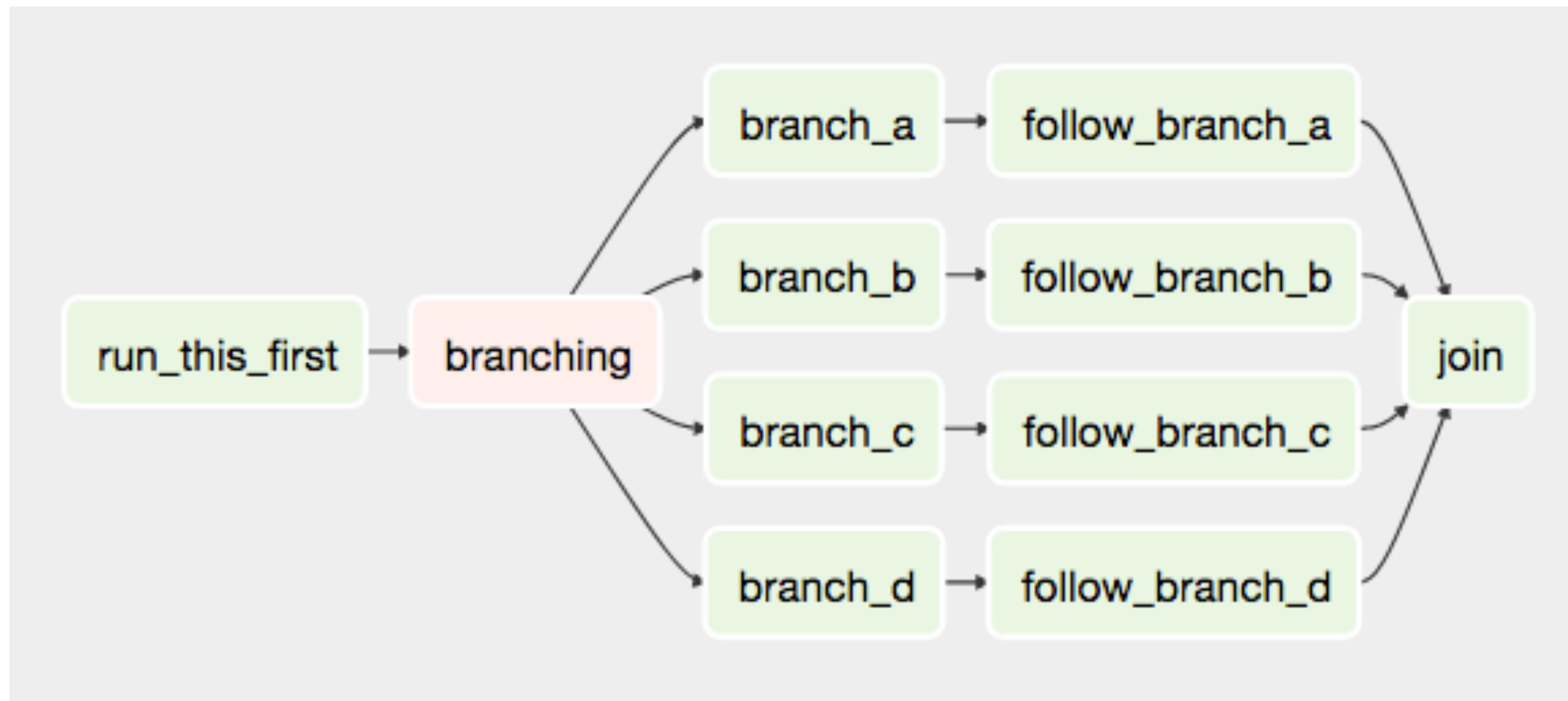
- Нужно запустить выгрузку в 8:00 из DB
- При загрузке данных в базу, сделать инференс всех новых строк?
- Если до 12:00 неактивный Slack, написать “Работаю из дома” в чат
- Найти и полайкать все новые фотографии жены, сестры, мамы, тещи
- Сделать все это, в определенном порядке и передавая состояния между задачами

Мотивация



Мотивация





Инструменты

- Oozie (<https://github.com/apache/oozie>)
- Luigi (<https://github.com/spotify/luigi>)
- AirFlow (<https://github.com/apache/airflow>)
- Azkaban (<https://github.com/azkaban/azkaban>)
- Conductor (<https://github.com/Netflix/conductor>)
- Pinball (<https://github.com/pinterest/pinball>)
- Chronos (<https://github.com/mesos/chronos>)

AirFlow

- Airflow was started in October 2014 by Maxime Beauchemin at Airbnb
- It was open source from the very first commit and officially brought under the Airbnb GitHub and announced in June 2015
- The project joined the Apache Software Foundation's Incubator program in March 2016
- The Foundation announced Apache Airflow as a Top-Level Project in January 2019.

AirFlow

- Python 2.7, 3.5, 3.6, 3.7
- Flask
- Jinja
- Relationships database
- LDAP

Составляющие

- Scheduler
- Webserver
- Worker

Webserver

Airflow DAGs Data Profiling ▾ Browse ▾ Admin ▾ Docs ▾ About ▾ 2018-09-07 22:14:10 UTC ⌂

DAGs

Search:

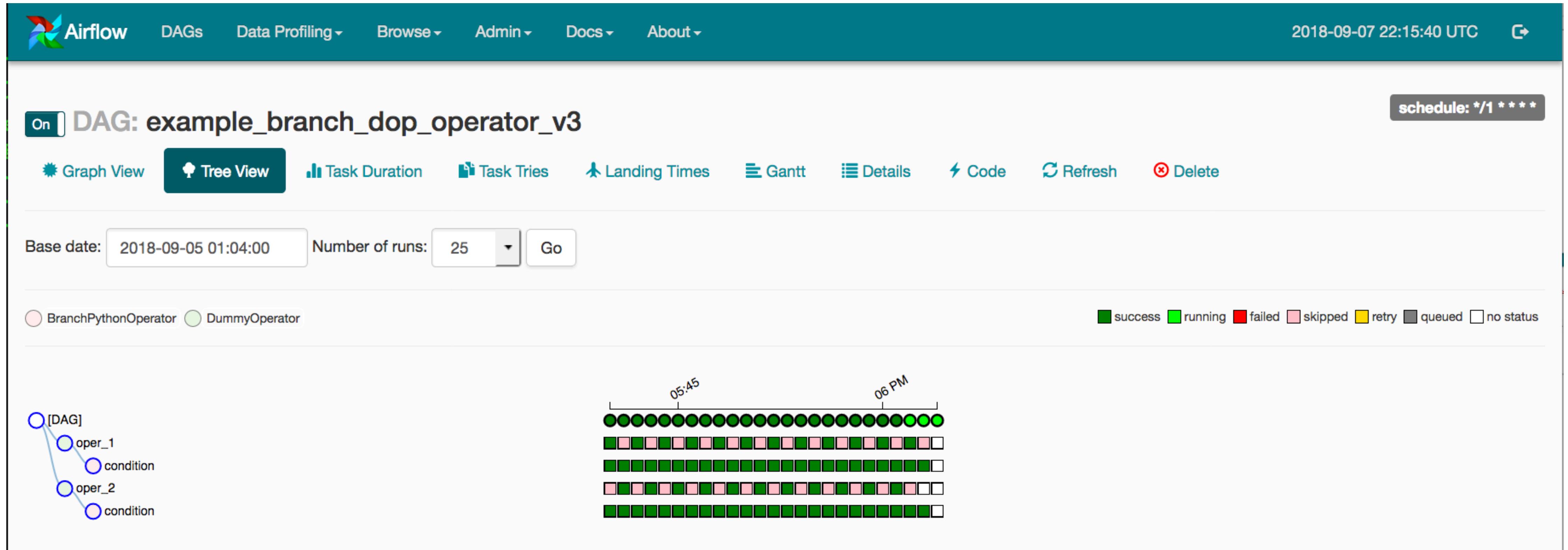
		DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
		example_bash_operator	0 0 * * *	airflow		2018-09-06 00:00		
		example_branch_dop_operator_v3	* /1 * * *	airflow		2018-09-05 00:56		
		example_branch_operator	@daily	airflow		2018-09-06 00:00		
		example_xcom	@once	airflow		2018-09-05 00:00		
		latest_only	4:00:00	Airflow		2018-09-07 16:00		

Showing 1 to 5 of 5 entries

« < 1 > »

Show Paused DAGs

Webserver



Webserver

Airflow DAGs Data Profiling ▾ Browse ▾ Admin ▾ Docs ▾ About ▾ 2018-09-07 22:29:47 UTC ⏪

On DAG: example_bash_operator schedule: 0 0 ***

Graph View Tree View Task Duration Task Tries Landing Times Gantt Details Code Refresh Delete

success Base date: 2018-09-06 00:00:01 Number of runs: 25 Run: scheduled_2018-09-06T00:00:00+00:00 Layout: Left->Right Go Search for...

BashOperator DummyOperator success running failed skipped retry queued no status

```
graph TD; runme_0[runme_0] --> run_after_loop[run_after_loop]; runme_1[runme_1] --> run_after_loop; runme_2[runme_2] --> run_after_loop; runme_2 --> also_run_this[also_run_this]; also_run_this --> run_this_last[run_this_last]; run_after_loop --> run_this_last;
```

Webserver

Screenshot of the Airflow Variables page.

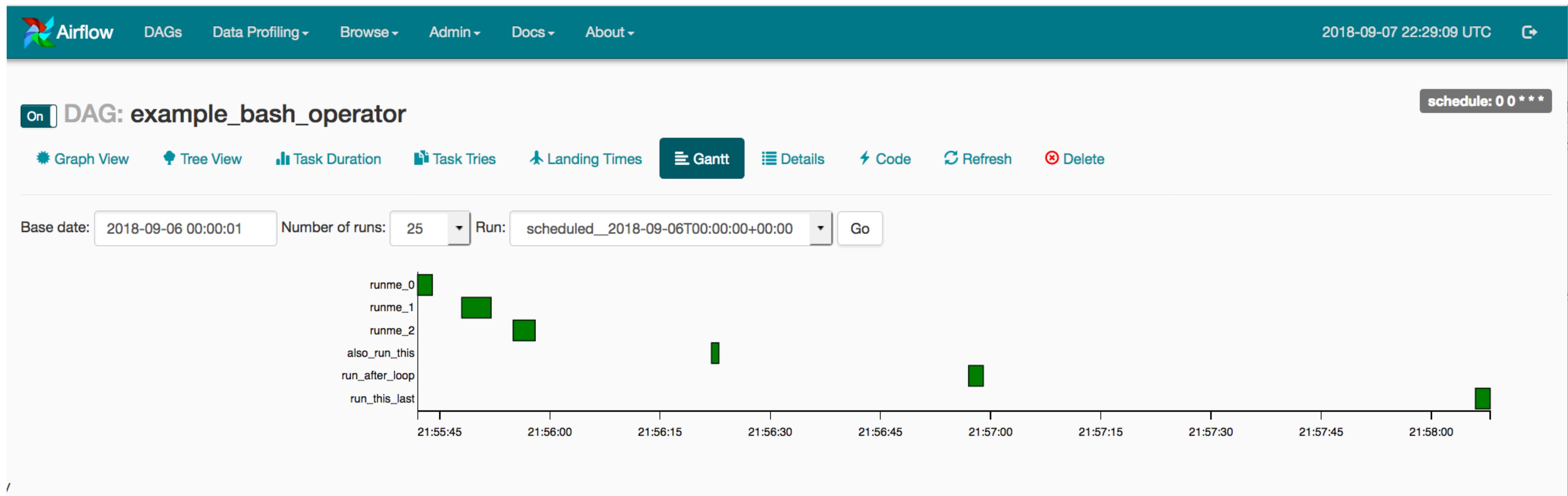
The page title is "Variables".

Header navigation includes: Airflow, DAGs, Data Profiling ▾, Browse ▾, Admin ▾, and Docs ▾.

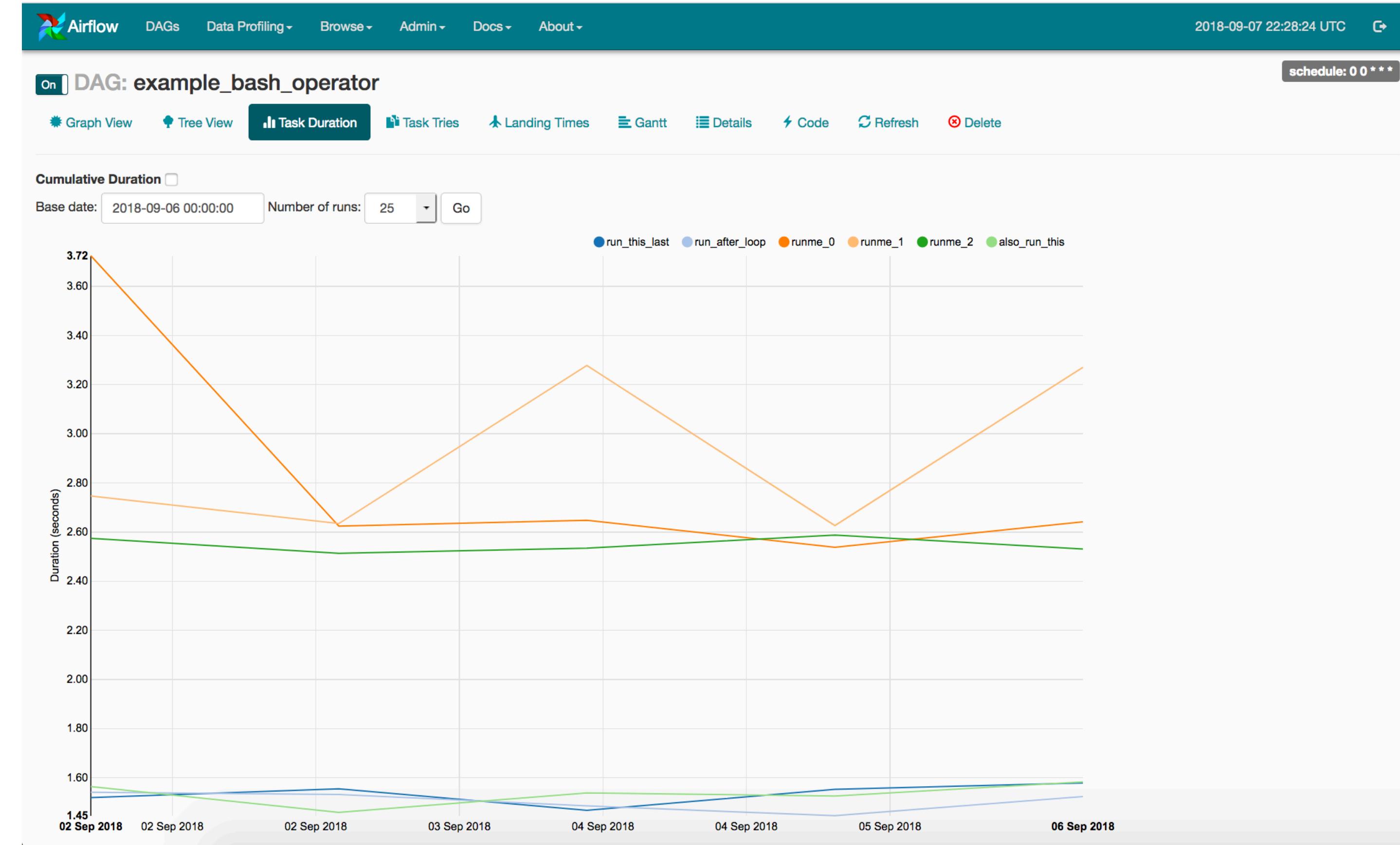
Action buttons: List (9), Create, Add Filter ▾, With selected ▾, and Search.

	Key	Val
<input type="checkbox"/>	secret_password	*****
<input type="checkbox"/>	not_so_hidden	test value
<input type="checkbox"/>	secret	*****
<input type="checkbox"/>	password	*****
<input type="checkbox"/>	passwd	*****
<input type="checkbox"/>	api_key	*****
<input type="checkbox"/>	apikey	*****
<input type="checkbox"/>	authorization	*****
<input type="checkbox"/>	access_token	*****

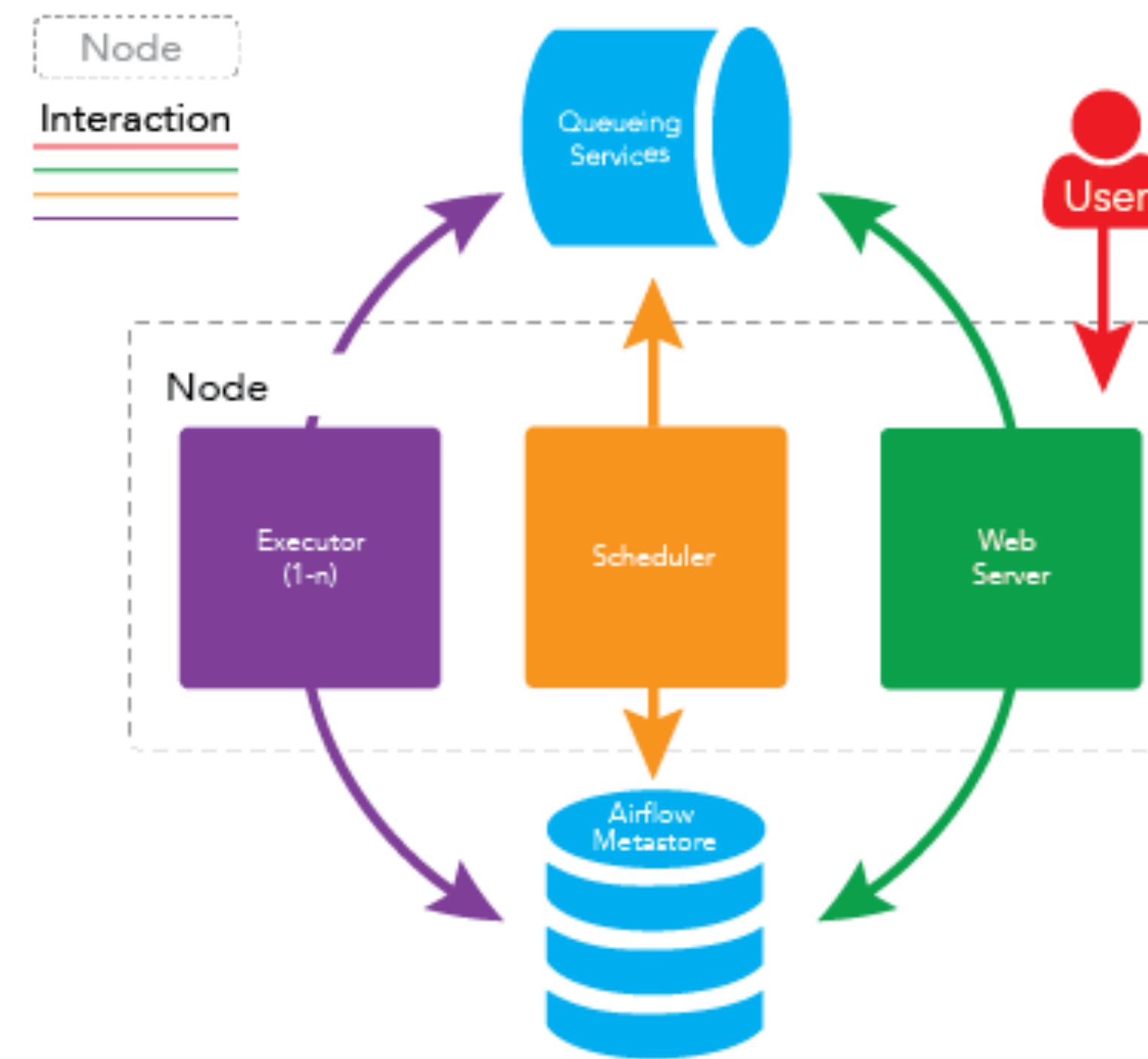
Webserver



Webserver



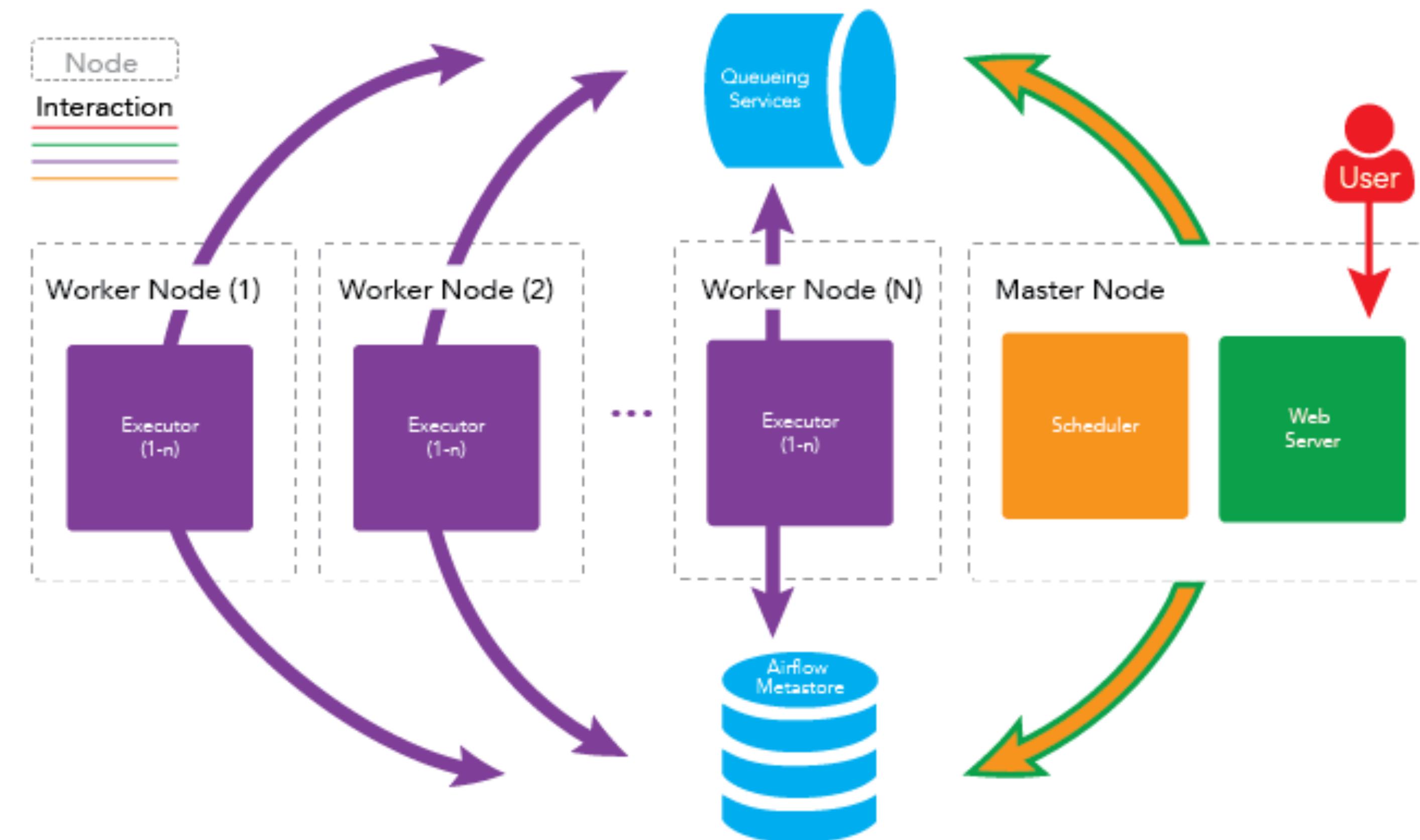
Архитектура



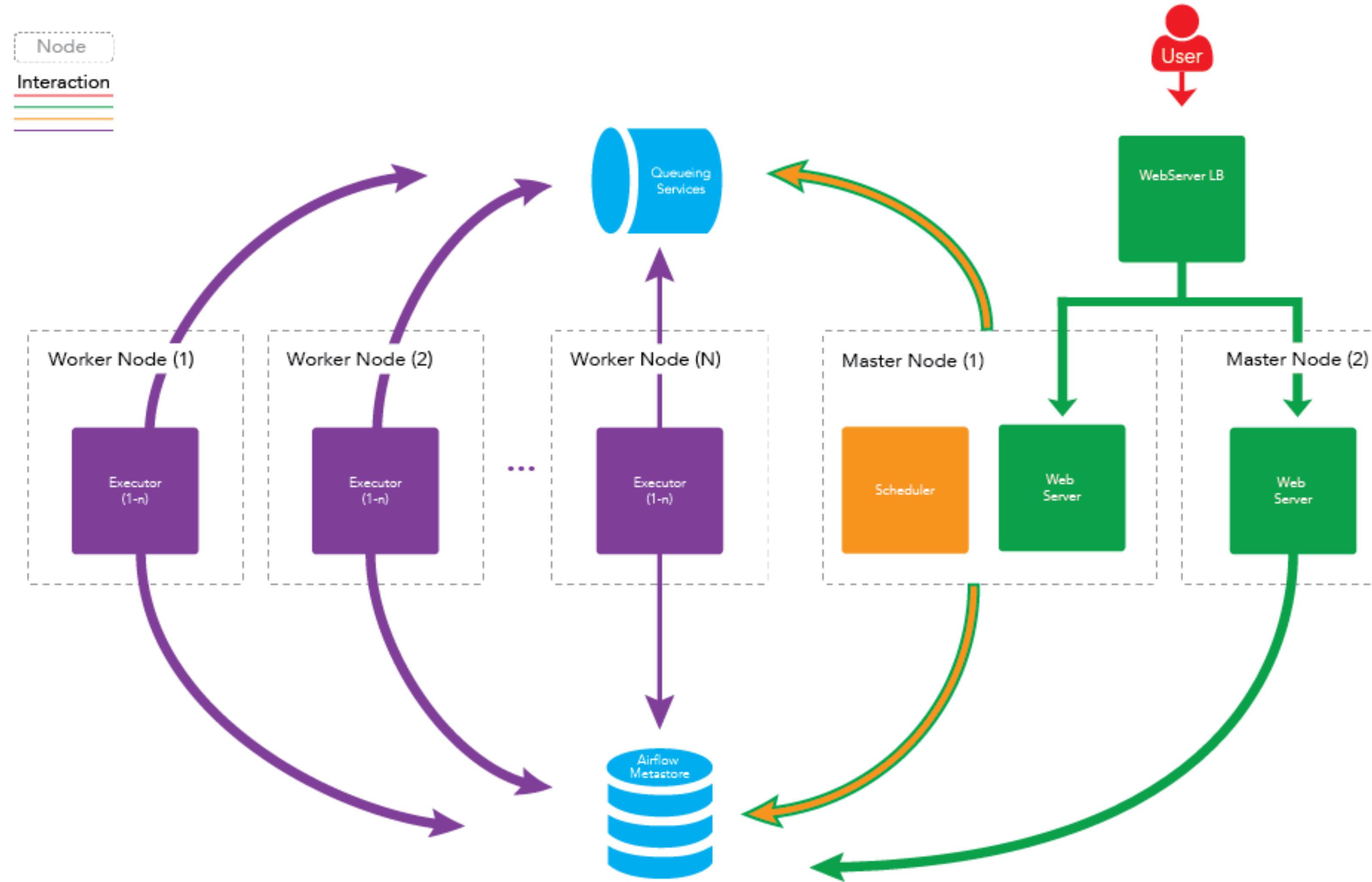
Scheduler - executors

- sequential_executor
- local_executor
- celery_executor
- dask_executor
- kubernetes_executor

Архитектура

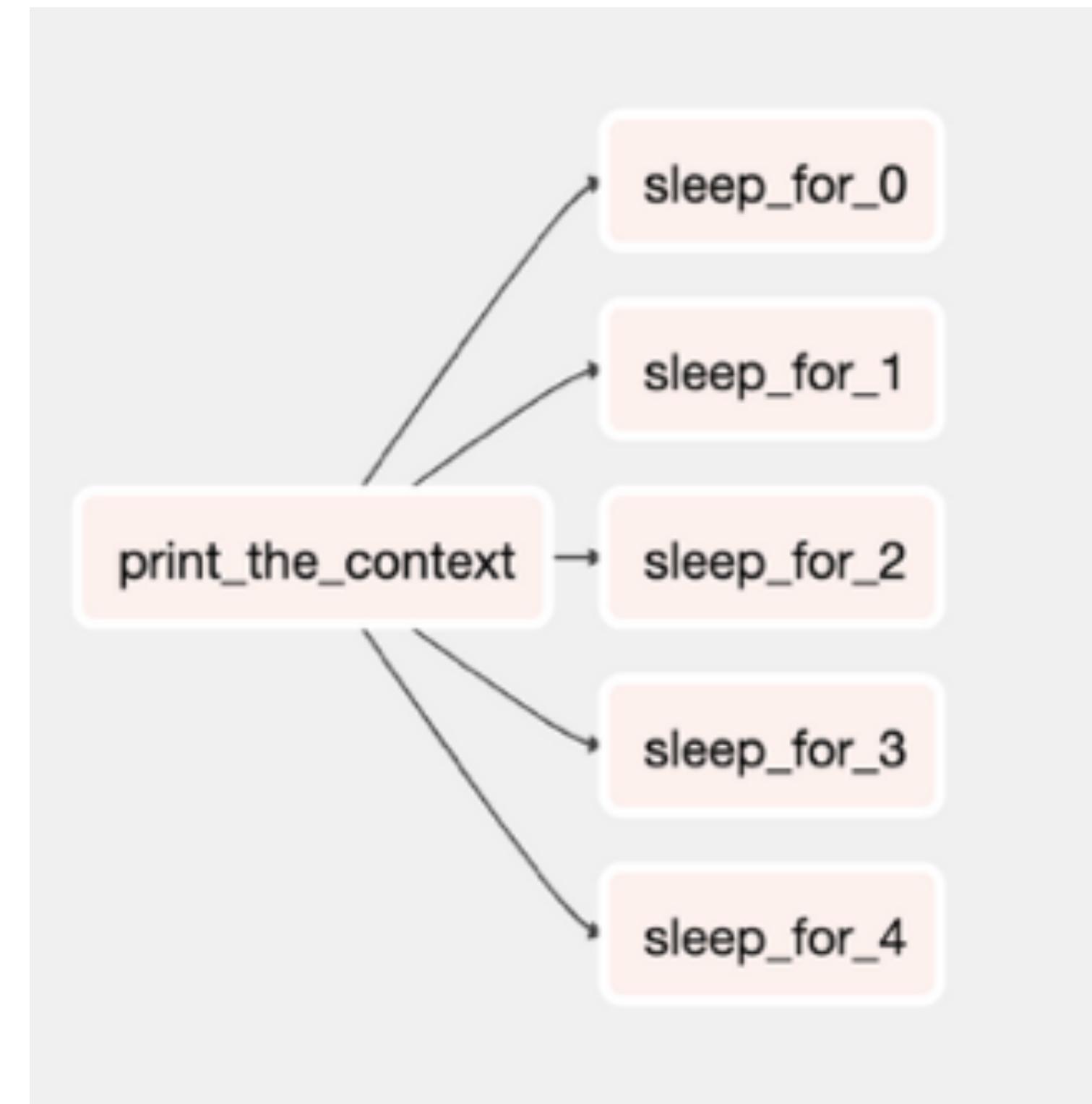


Архитектура

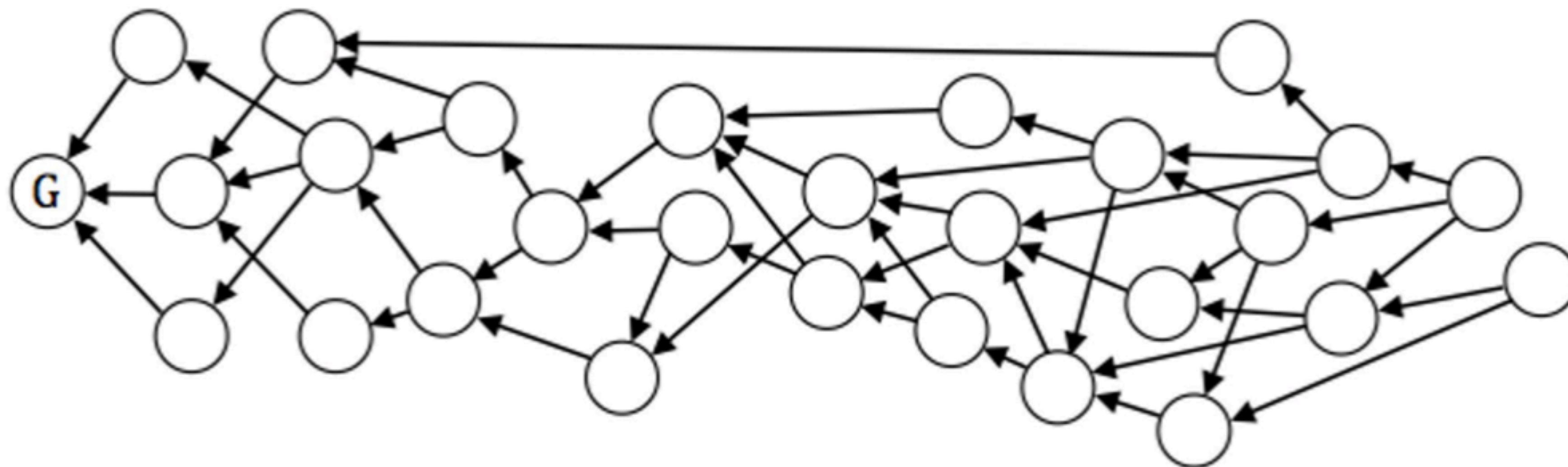


Easy

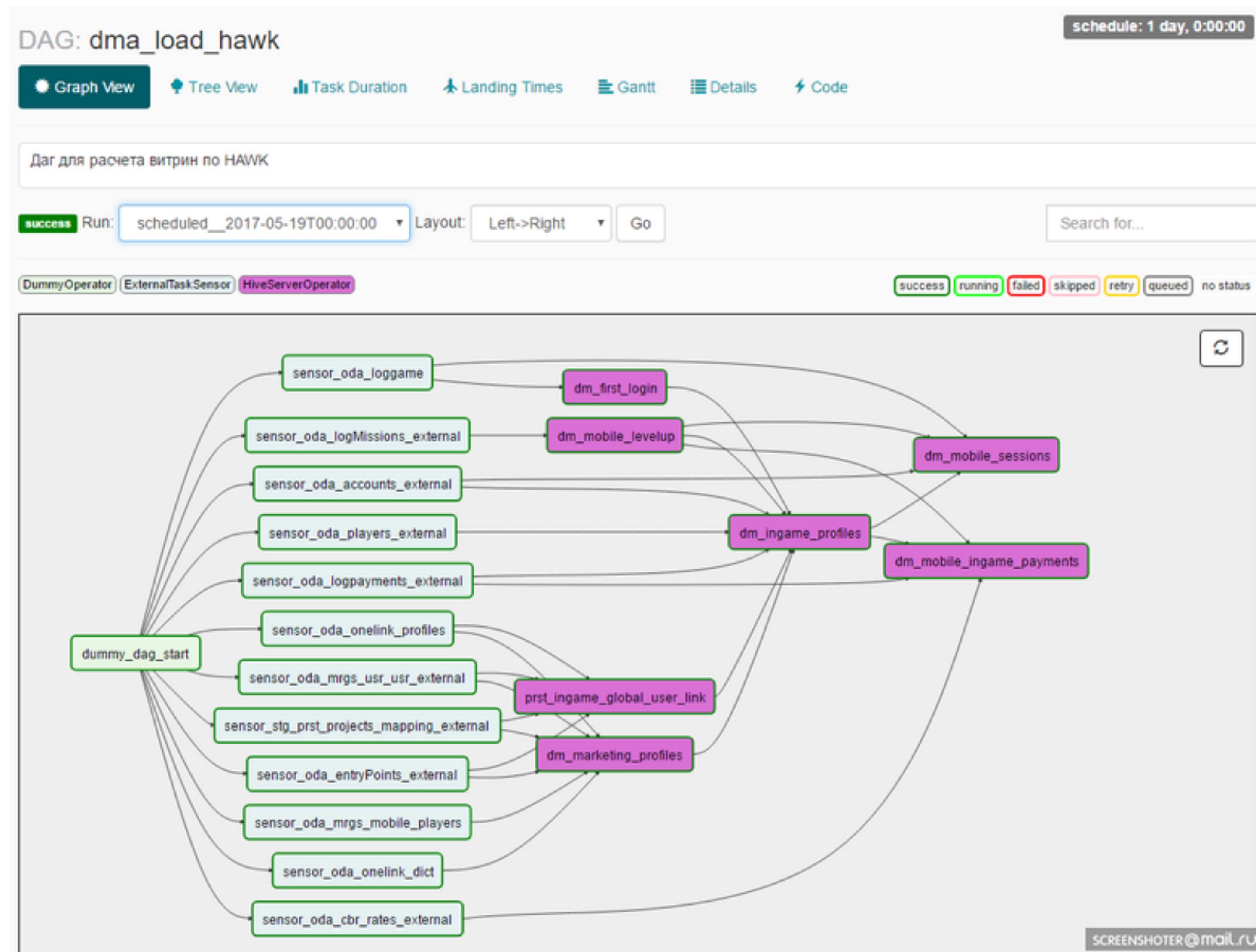
```
1 import time
2 from pprint import pprint
3
4 import airflow
5 from airflow.models import DAG
6 from airflow.operators.python_operator import PythonOperator
7
8 args = {
9     'owner': 'Airflow',
10    'start_date': airflow.utils.dates.days_ago(2),
11 }
12
13 dag = DAG(
14     dag_id='example_python_operator',
15     default_args=args,
16     schedule_interval="0 15 1 * *",
17 )
18
19 def print_context(ds, **kwargs):
20     pprint(kwargs)
21     print(ds)
22     return 'Whatever you return gets printed in the logs'
23
24 run_this = PythonOperator(
25     task_id='print_the_context',
26     python_callable=print_context,
27     dag=dag,
28 )
29
```



DAG



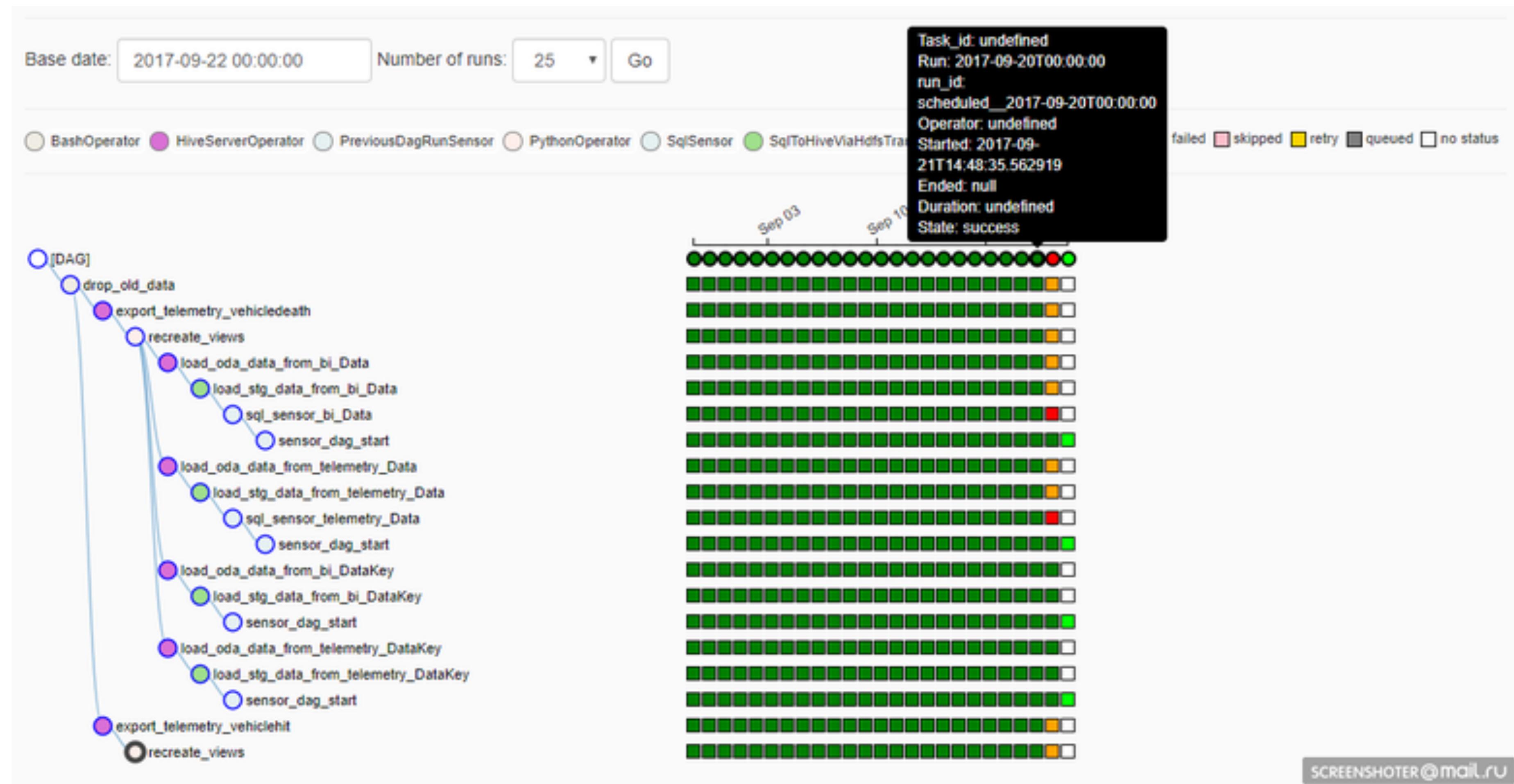
DAG



Операторы

- BashOperator – оператор для выполнения bash-команды
- PythonOperator – оператор для вызова Python-кода
- EmailOperator – оператор для отправки email'a
- HTTPOperator – оператор для работы с http-запросами
- SqlOperator – оператор для выполнения SQL-кода
- Sensor – оператор ожидания события

Execution Date



Очереди

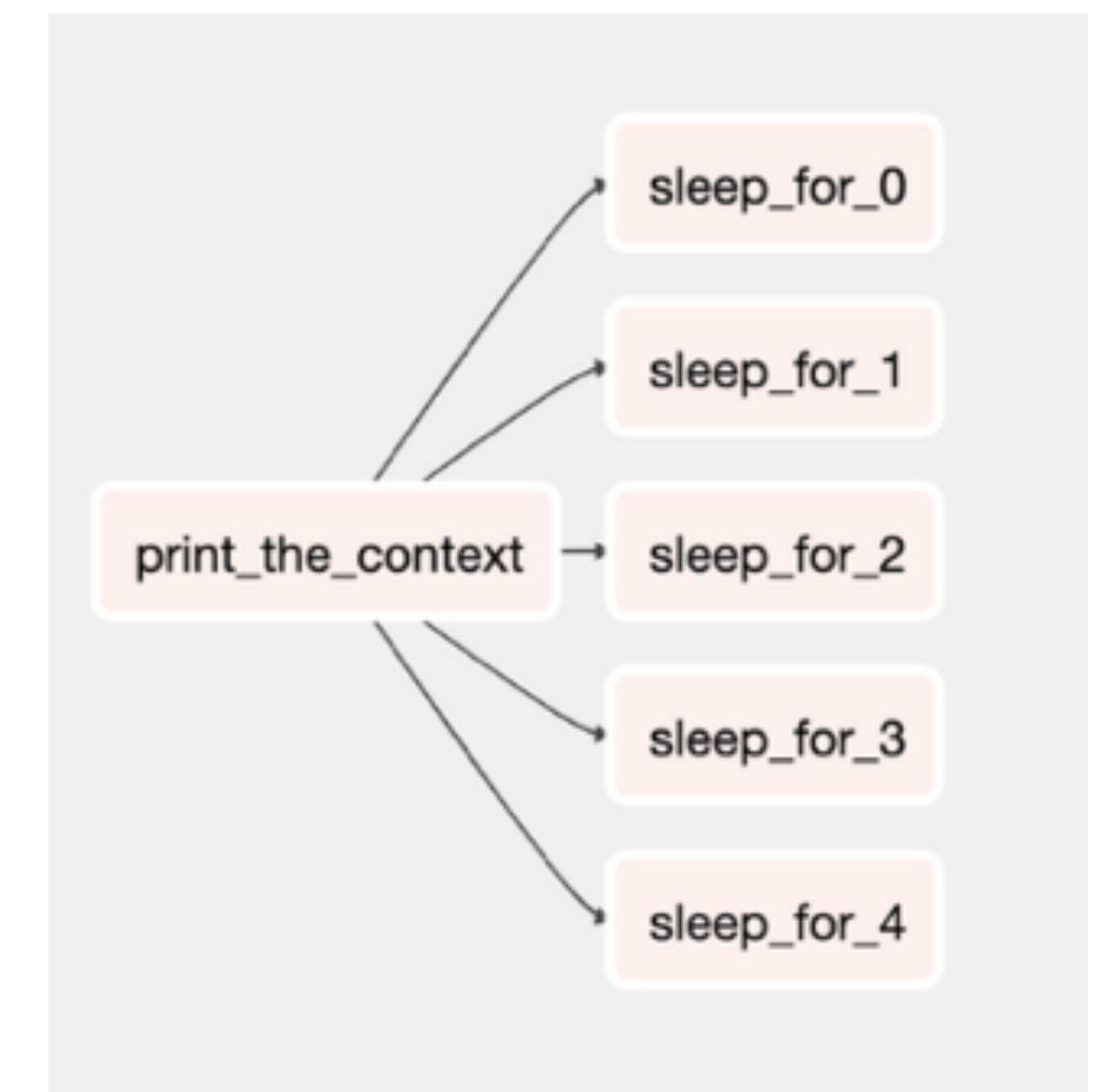


XCom

Task Details Rendered Template Log XCom

XCom

Key	Value
table_name	squads_tourney02122017_2100

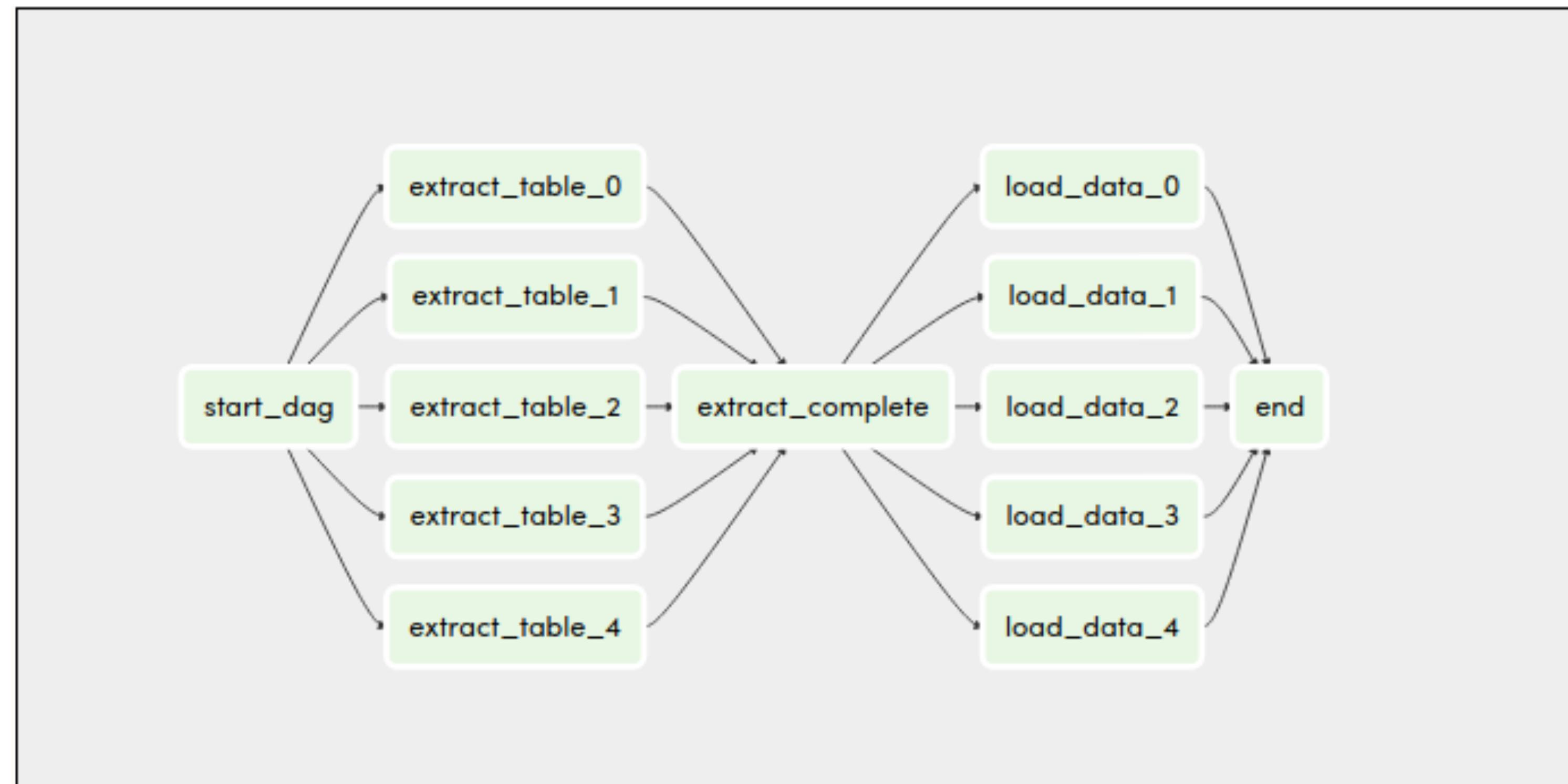


Best Practices

Идемпотентность

- $\forall x : f(f(x)) = f(x)$
- $a = a + 0 = (a + 0) + 0 = ((a + 0) + 0) + 0 =$

Гранулярность



Sensors

- ExternalTaskSensor
- HdfsSensor
- HttpSensor
- HivePartitionSensor
- SqlSensor
- TimeDeltaSensor
- TimeSensor
- WebHdfsSensor



Дополнительно

- UTC
- <https://github.com/epoch8/airflow-exporter>
- Нет долгим scheduler-ам

Генерация DAG-ов

```
def generate_dags():
    dags = []

    for hql_filename in hql_files:

        default_args = {
            'owner': 'rnd_metrics',
            'depends_on_past': False,
            'start_date': datetime.now() - timedelta(days=2),
            'retries': 1,
            'retry_delay': timedelta(minutes=5)
        }

        dag = DAG(
            dag_id=metric_name+'_dag',
            default_args=default_args,
            schedule_interval=schedule_interval,
            max_active_runs=5,
            catchup=True,
            concurrency=5
        )

        tasks = [
            PythonOperator(
                task_id=metric_name+'_operator',
                python_callable=upload_metrics,
                op_kwargs={
                    'hql_query': hql_query,
                    'metric_name': metric_name
                },
                provide_context=True,
                dag=dag
            )
        ]
        dags.append({'dag': dag, 'tasks': tasks})

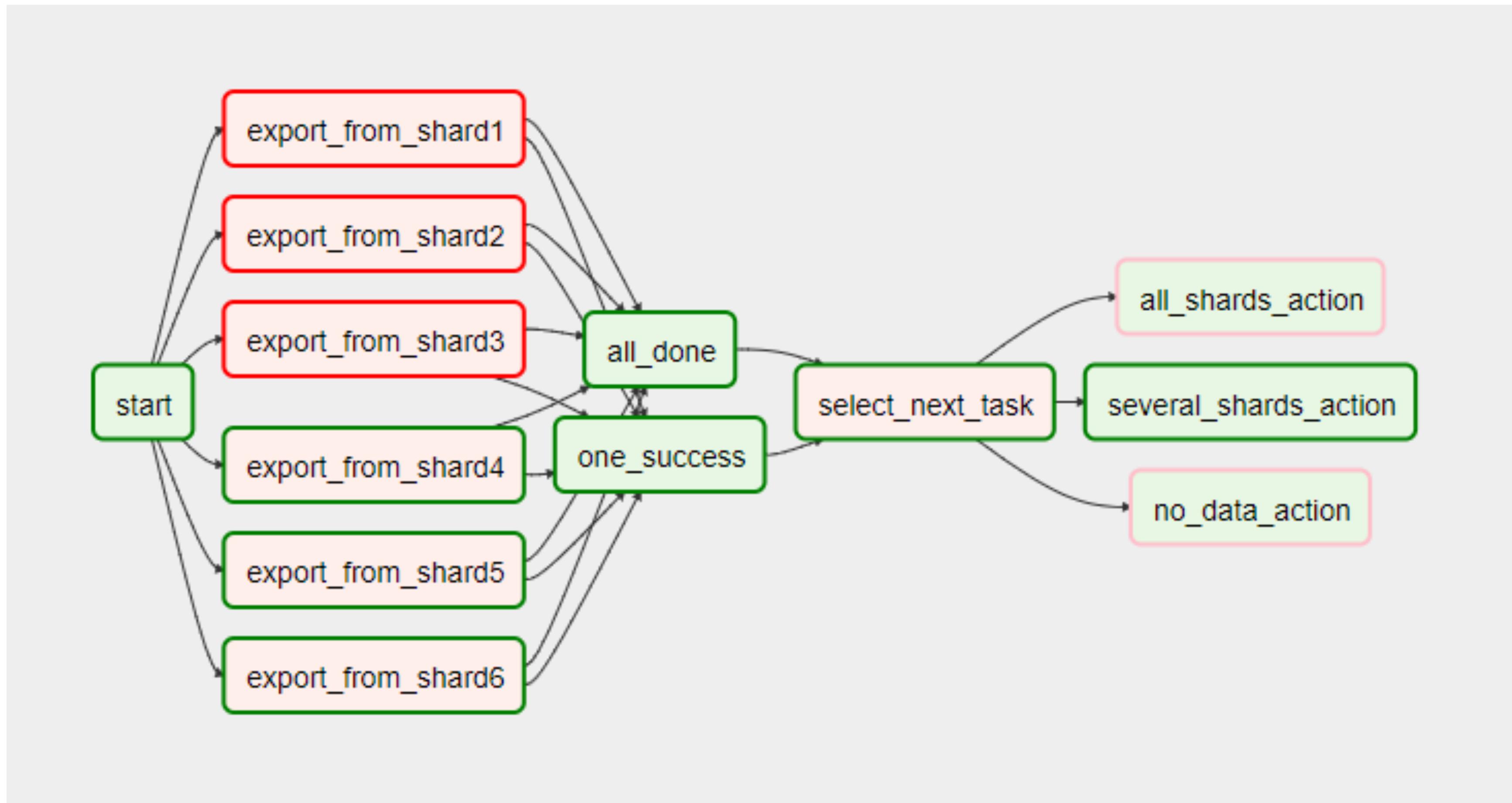
    return dags
```

```
for dag in generate_dags():
    globals()[dag['dag'].dag_id] = dag['dag']

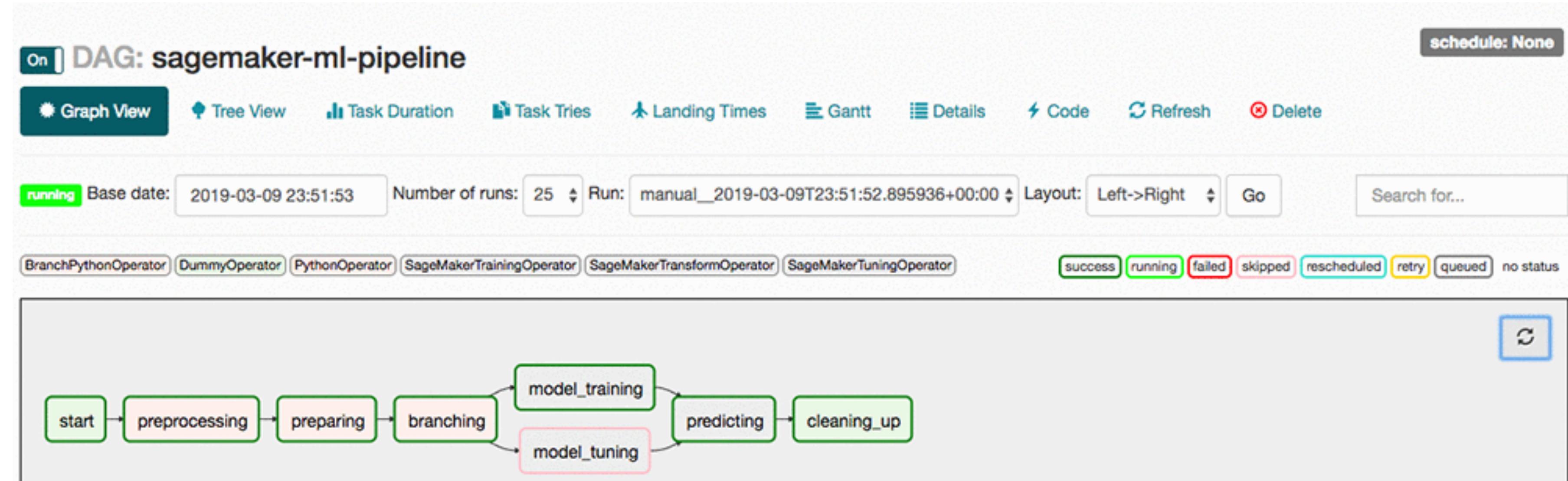
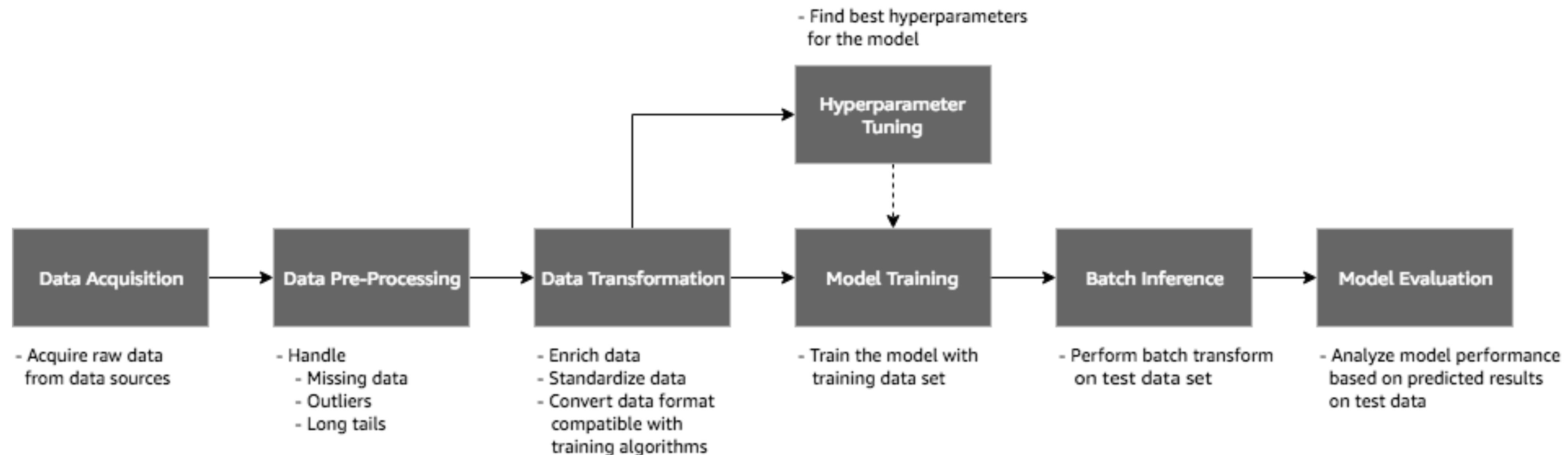
    DummyOperator(
        dag=dag['dag'],
        task_id=f'dummy_{dag["dag"].dag_id}'
    ).set_upstream(dag['tasks'])
```



Ветвление



ML



MLcomp

Projects Dags Tasks Models Reports Logs Computers Auxiliary

Filters

Project	Id	Name	Tasks	Created	Started	Last activity	Duration	Task status	Links	Image size	File size
severstal	51	kaggle_kernel	1	08.19 14:59:20	08.19 14:59:34	08.19 15:01:28	1 min 54 sec	○ ○ ○ ○ ○ 1		0 byte	0 byte
severstal	50	kaggle_kernel	1	08.19 14:59:12	08.19 14:59:14	08.19 15:00:41	1 min 27 sec	○ ○ ○ ○ ○ 1		0 byte	800 byte
severstal	49	kaggle_kernel	1	08.19 14:57:56	08.19 14:58:42	08.19 14:59:06	24 sec	○ ○ ○ ○ ○ 1		0 byte	0 byte
severstal	48	kaggle_kernel	1	08.19 14:56:12	08.19 14:57:44	08.19 15:57:37	59 min 52 sec	○ ○ ○ ○ ○ 1		0 byte	0 byte
severstal	47	kaggle_kernel	1	08.19 14:56:11	08.19 14:56:47	08.19 15:57:45	1 hour 0 min	○ ○ ○ ○ ○ 1		0 byte	0 byte

Items per page: 5 1 - 5 of 35 < >

Catalyst

MLcomp

Почитать и посмотреть

- <https://www.astronomer.io/podcast/>
- <http://site.clairvoyantsoft.com/setting-apache-airflow-cluster/>
- <https://www.astronomer.io/guides/airflow-executors-explained/>
- <https://medium.com/@maximebeauchemin/functional-data-engineering-a-modern-paradigm-for-batch-data-processing-2327ec32c42a>
- <https://habr.com/ru/company/mailru/blog/339392/>
- <https://habr.com/ru/company/mailru/blog/344398/>
- https://blog.twitter.com/engineering/en_us/topics/insights/2018/ml-workflows.html