

- Workflows Management Systems
- Architecture
- Building blocks
- More features
- User Interface
- Security
- CLI
- Demo

WTH is a Workflow Management System ?

A workflow Management system is:

Is a data-centric software (framework) for :

- **Setting up**
- **Performing**
- **Monitoring**

of a defined **sequence of processes** and **tasks**

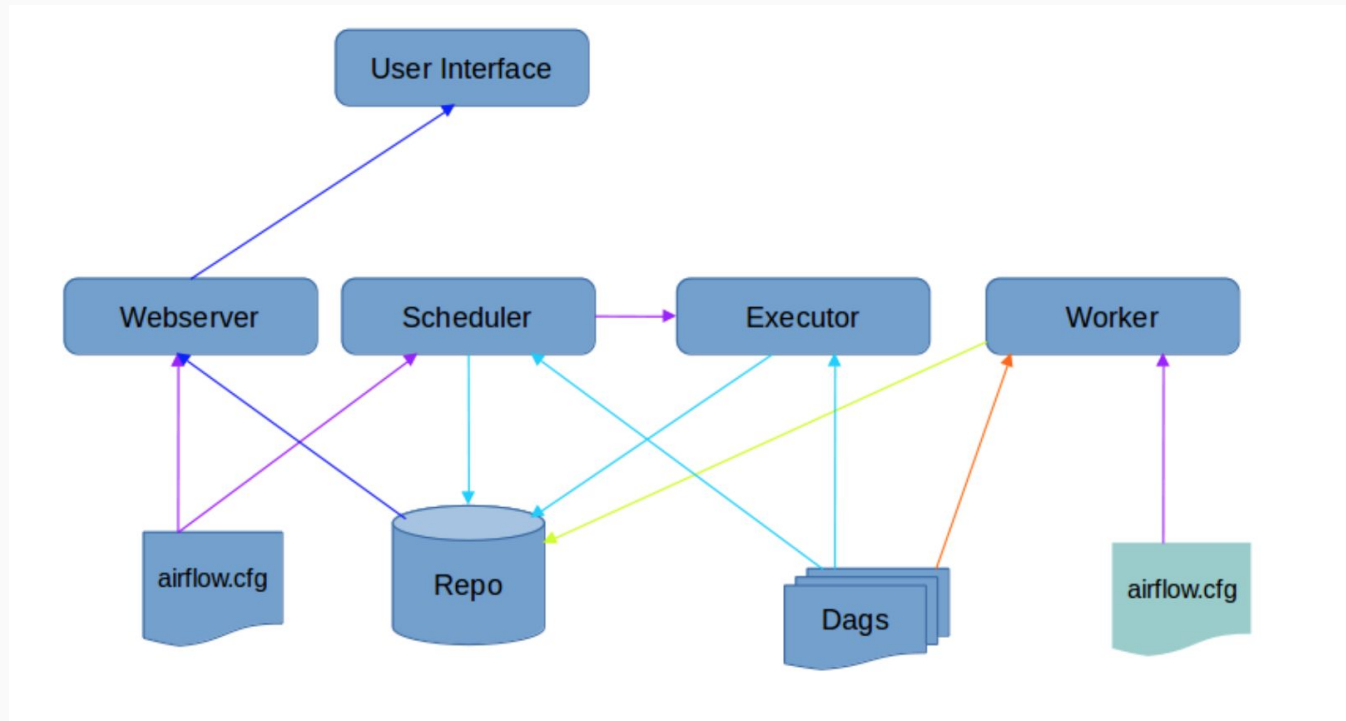
Popular Workflow Management Systems



Airflow Architecture

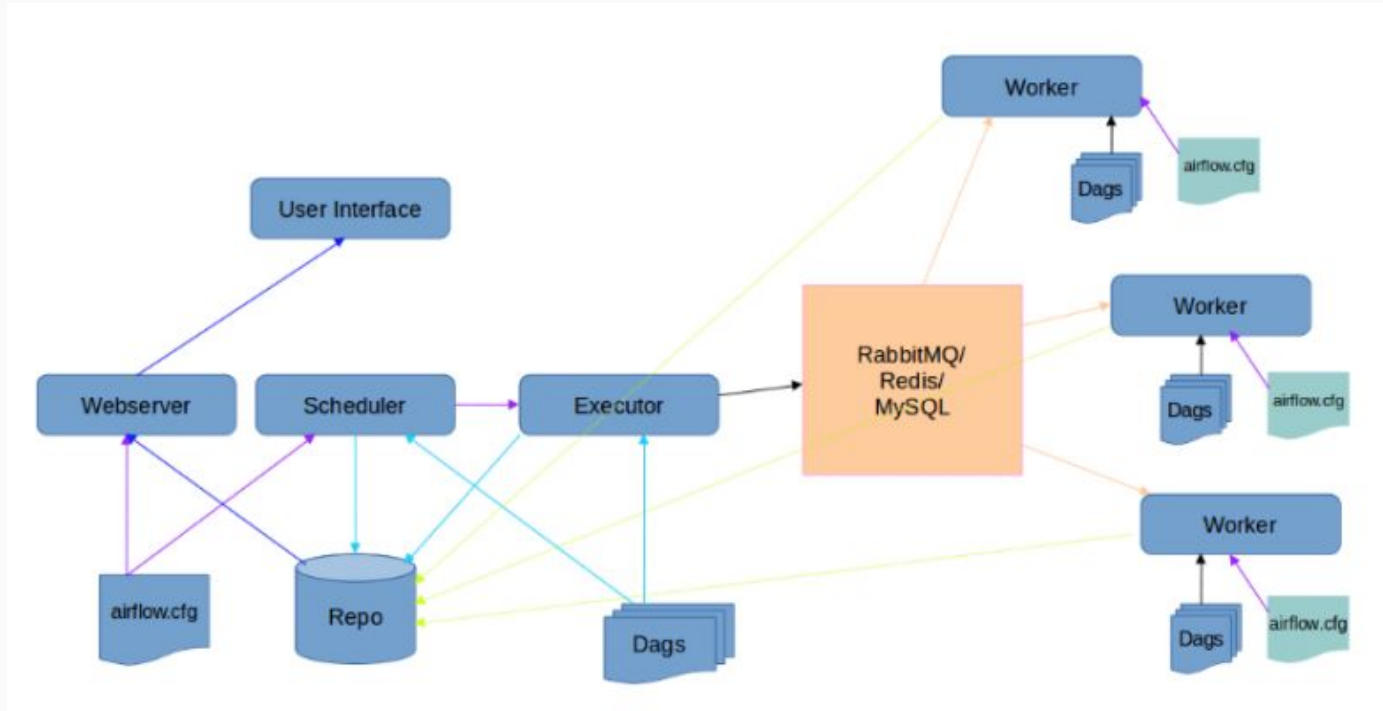
Airflow architecture

SequentialExecutor / LocalExecutor



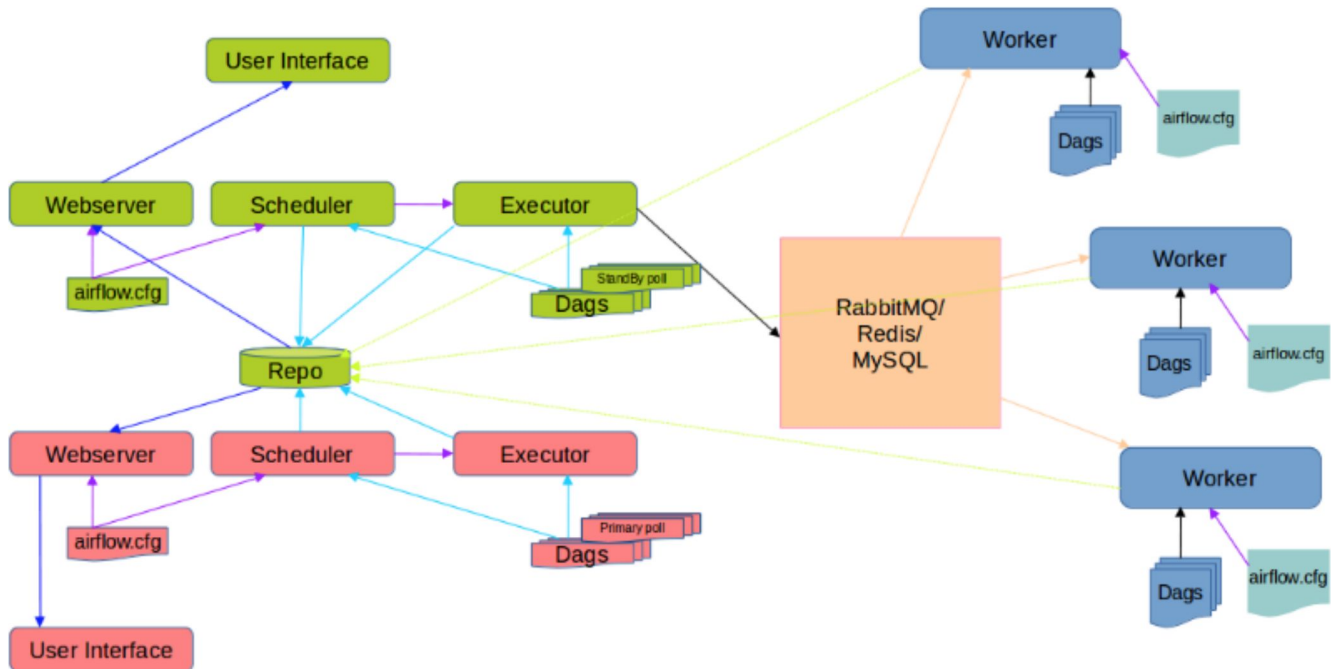
Airflow architecture

CeleryExecutor



Airflow architecture

HA + CeleryExecutor



Airflow architecture

- MesosExecutor : already available in contrib package
- KubernetesExecutor ??

Building blocks

Dags :

- Directed Acyclic Graph
- Is a collection of all the tasks you want to run
- DAGs describe *how* to run a workflow

Dags :

```
1  from airflow.models import DAG
2
3  default_args = {
4      'owner': 'airflow',
5      'start_date': datetime(2018, 3, 3),
6      'schedule_interval' : '0 12 * * *',
7      'depends_on_past': False,
8      'retry_delay': timedelta(minutes=5),
9      'retries': 1,
10     'email': [''],
11     'email_on_failure': True,
12     'email_on_retry': True
13 }
14
15 facebook_insights_dag = DAG(
16     dag_id='facebook_insights_dag',
17     default_args=default_args
18 )
19
20
```

Operators :

- Describes a **single task** in a workflow.
- Determine what actually gets done
- Operators generally run independently (atomic)
- The DAG make sure that operators run in the correct certain **order**
- They may run on completely **different machines**

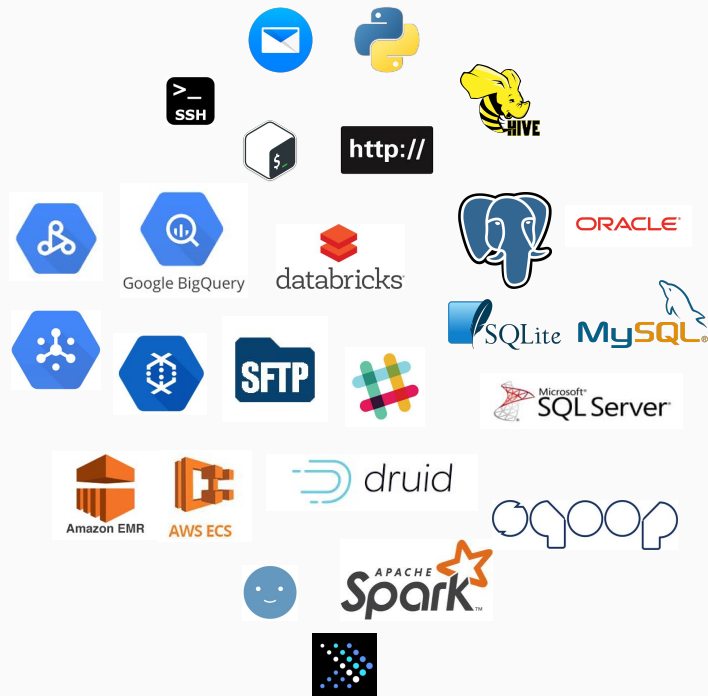
Operators : There are 3 main types of operators:

- **Operators** that performs an **action**, or tell another system to perform an action
- **Transfer** operators move data from one system to another
- **Sensors** are a certain type of operator that will **keep running** until a certain **criterion** is met.
 - Examples include a specific file landing in HDFS or S3.
 - A partition appearing in Hive.
 - A specific time of the day.

Building blocks

Operators :

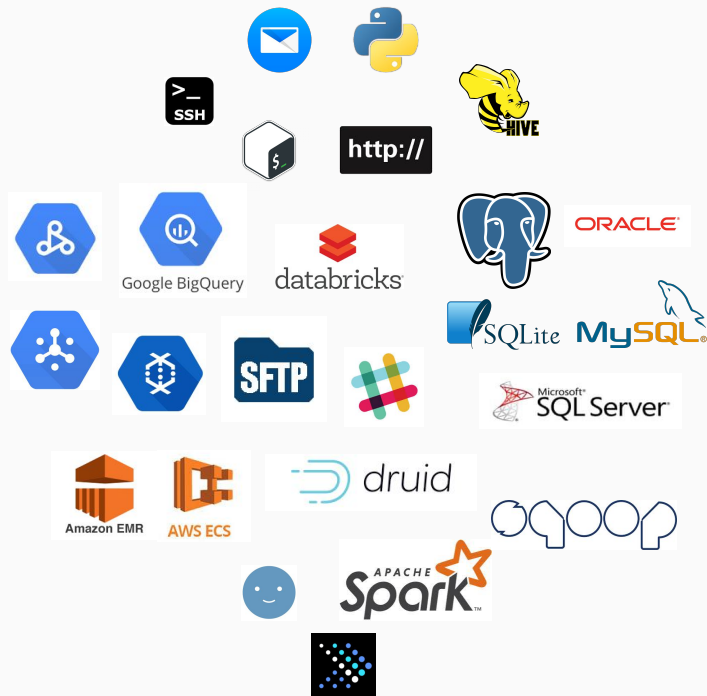
- Operators :
 - BashOperator
 - PythonOperator
 - EmailOperator
 - HTTPOperator
 - MySqlOperator
 - SqliteOperator
 - PostgresOperator
 - MsSqlOperator
 - OracleOperator
 - JdbcOperator
 - DockerOperator
 - HiveOperator
 - SlackOperator



Building blocks

Operators :

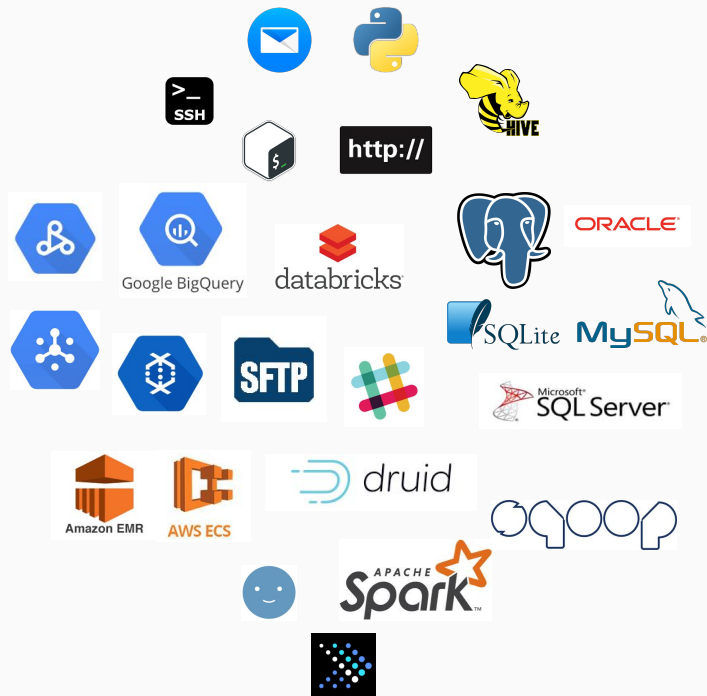
- Transfers :
 - S3FileTransferOperator
 - PrestoToMySQLOperator
 - MySQLToHiveTransfer
 - S3ToHiveTransfer
 - BigQueryToCloudStorageOperator
 - GenericTransfer
 - HiveToDruidTransfer
 - HiveToMySQLTransfer



Building blocks

Operators :

- Sensors :
 - ExternalTaskSensor
 - HdfsSensor
 - HttpSensor
 - MetastorePartitionSensor
 - HivePartitionSensor
 - S3KeySensor
 - S3PrefixSensor
 - SqlSensor
 - TimeDeltaSensor
 - TimeSensor
 - WebHdfsSensor



Operators :

```
1 def fetch_ads_infos(**kwargs):
2     ...
3 def enrich_campaigns_insights(**kwargs):
4     ...
5
6 ads_infos_task = PythonOperator(
7     task_id='ads_infos_task',
8     python_callable=fetch_ads_infos,
9     provide_context=True,
10    dag=facebook_insights_dag
11 )
12
13 enrich_campaigns_lifetime_insights_task = PythonOperator(
14     task_id='enrich_campaigns_lifetime_insights_task',
15     op_kwargs={
16         'granularity': 'lifetime'
17     },
18     python_callable=enrich_campaigns_insights,
19     provide_context=True,
20     dag=facebook_insights_dag
21 )
22
23 ads_infos_task >> enrich_campaigns_lifetime_insights_task
24 ads_infos_task.set_downstream(enrich_campaigns_lifetime_insights_task)
25
26 enrich_campaigns_lifetime_insights_task << ads_infos_task
27 enrich_campaigns_lifetime_insights_task.set_upstream(ads_infos_task)
28
29
30
```

Tasks : a parameterized instance of an operator

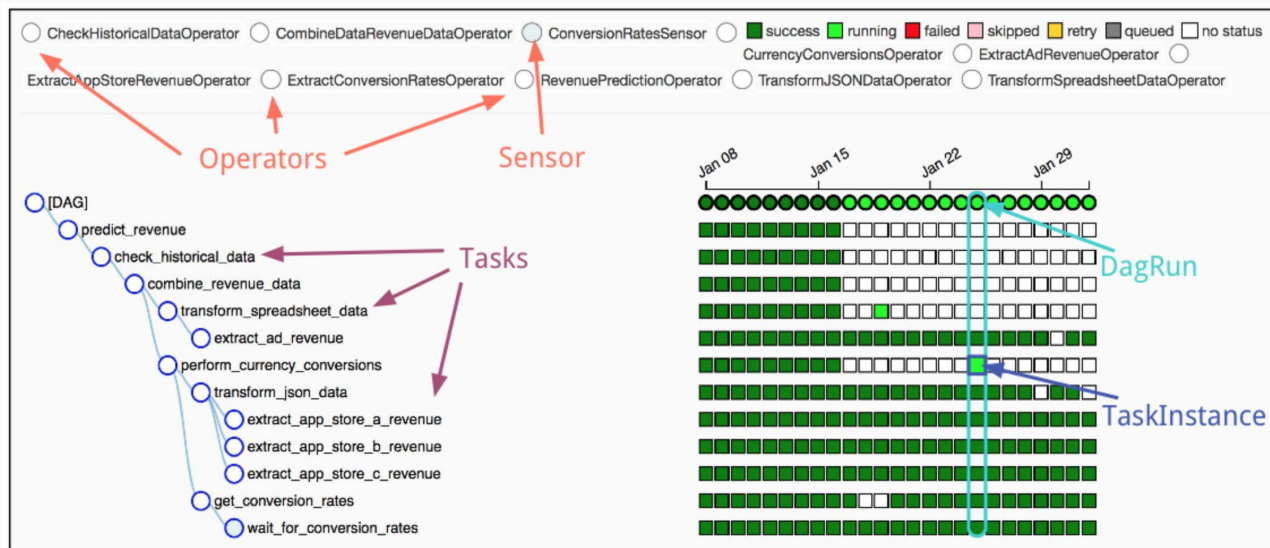
Task Instance : Dag + Task + point in time

- Specific run of a Task
- A task assigned to a DAG
- Has State associated with a specific run of the DAG
- States : it could be
 - **running**
 - **success,**
 - **failed**
 - **skipped**
 - **up for retry**
 - ...

Workflows :

- **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
 - Has a state associated with a specific run of the DAG
- By combining **DAGs** and Operators to create **TaskInstances**, you can build **complex workflows**.

Building blocks



More features

More features

- Features :
 - Hooks
 - Connections
 - Variables
 - XComs
 - SLA
 - Pools
 - Queues
 - Trigger Rules
 - ~~- Branchings~~
 - ~~- SubDags~~

Hooks :

- Interface to external platforms and databases :
 - Hive
 - S3
 - MySQL
 - PostgreSQL
 - HDFS
 - Hive
 - Pig
 - ...
- Act as building block for Operators
- Use **Connection** to retrieve authentication informations
- Keep authentication infos out of pipelines.

Connections :

Connection informations to external systems are stored in the airflow metadata Database and managed in the UI

More features

Airflow								
DAGs		Data Profiling		Browse		Admin		Docs
About								15:58 UTC
<input type="checkbox"/>		druid_ingest_default				Pools	druid-overlord	8081
<input type="checkbox"/>		emr_default				Configuration		
<input type="checkbox"/>		fs_default				Users		
<input type="checkbox"/>		google_cloud_default				Connections		
<input type="checkbox"/>		hive_cli_default		hive_cli		Variables		
<input type="checkbox"/>		hiveserver2_default		hiveserver2	localhost	XComs		
<input type="checkbox"/>		http_default		http	https://www.google.com/			
<input type="checkbox"/>		local_mysql		mysql	localhost			
<input type="checkbox"/>		metastore_default		hive_metastore	localhost			
<input type="checkbox"/>		mssql_default		mssql	localhost			
<input type="checkbox"/>		mysql_default		mysql	localhost			
<input type="checkbox"/>		postgres_default		postgres	localhost			
<input type="checkbox"/>		presto_default		presto	localhost			
<input type="checkbox"/>		redis_default		redis	localhost			
<input type="checkbox"/>		s3_pr_connection		s3				
<input type="checkbox"/>		spark_default		spark	yarn			
<input type="checkbox"/>		sqlite_default		sqlite	/tmp/sqlite_default.db			
<input type="checkbox"/>		sqoop_default		sqoop	rmdbs			
<input type="checkbox"/>		ssh_default		ssh	localhost			
<input type="checkbox"/>		vertica_default		vertica	localhost			
<input type="checkbox"/>		wasb_default		wasb				
<input type="checkbox"/>		webhdfs_default		hdfs	localhost			

More features

Exemple de Hook + connection :

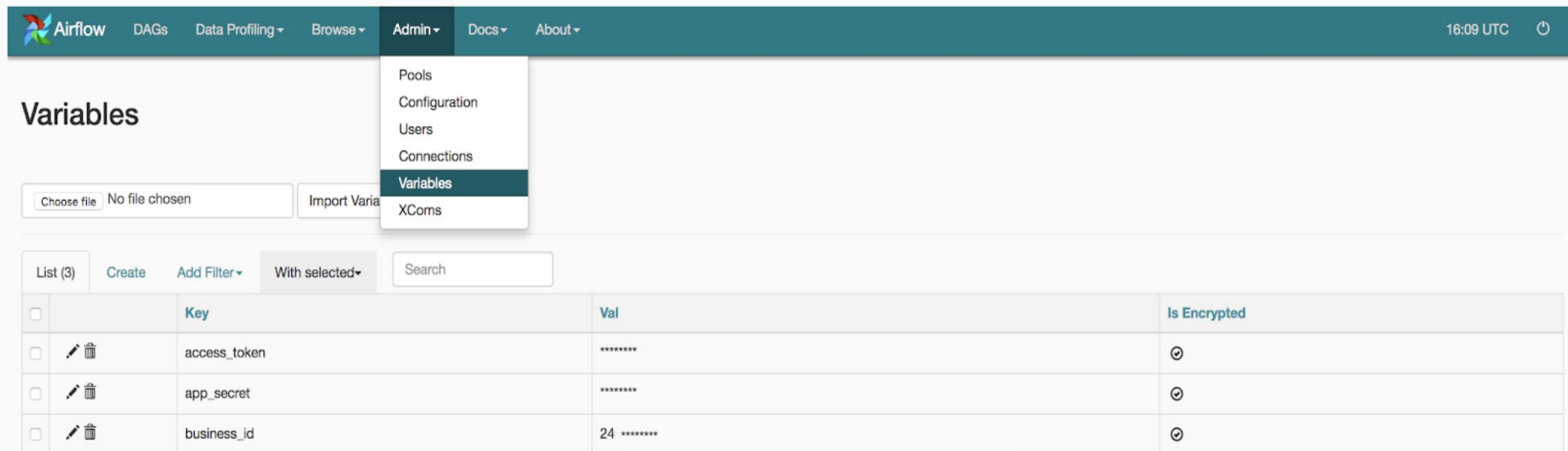
```
143
144 def enrich_campaigns_insights(granularity, **kwargs):
145     s3hook = S3Hook(aws_conn_id=AWS_AF_CONNECTION_NAME)
146     extraction_date = kwargs['execution_date'].strftime('%Y-%m-%d')
147
148     campaigns_infos_path = 'facebook-ads-and-campaigns/referential/campaigns_infos.json'
149     campaigns_insights_path = f'facebook-ads-and-campaigns/data/campaigns_{granularity}_{extraction_date}.json'
150     output_path = f'facebook-ads-and-campaigns/data/enriched_campaigns_{granularity}_{extraction_date}.json'
151
152     campaigns_infos_select_cols = ['account_id', 'id', 'start_time', 'stop_time', 'updated_time', 'status', 'created_time', 'effective_status']
153
154     campaigns_infos_content = s3hook.read_key(campaigns_infos_path, S3_BUCKET)
155     campaigns_insights_content = s3hook.read_key(campaigns_insights_path, S3_BUCKET)
156     campaigns_referential_df = pd.read_json(campaigns_infos_content, lines=True)
157     campaigns_insights_df = pd.read_json(campaigns_insights_content, lines=True)
158
159     enriched_campaigns_df = campaigns_insights_df.merge(campaigns_referential_df[campaigns_infos_select_cols],
160                                                         how='inner',
161                                                         left_on=['account_id', 'campaign_id'],
162                                                         right_on=['account_id', 'id'])
163
164     jl_data = enriched_campaigns_df.to_json(orient='records', lines=True)
165     s3hook.load_string(
166         jl_data,
167         output_path,
168         S3_BUCKET,
169         replace=True,
170         encrypt=False)
171
```

Variables :

- A generic way to store and retrieve arbitrary content or settings as a simple key value store within Airflow.
- Variables can be listed, created, updated and deleted from the UI (Admin -> Variables), code or CLI.
- While your pipeline code definition and most of your constants and variables should be defined in code and stored in source control, it can be useful to have some variables or configuration items accessible and modifiable through the UI.

```
1 from airflow.models import Variable
2 foo = Variable.get("foo")
3 bar = Variable.get("bar", deserialize_json=True)
```

More features









The screenshot displays the Airflow Admin web interface. The top navigation bar is dark teal with the Airflow logo and links for DAGs, Data Profiling, Browse, Admin, Docs, and About. The Admin dropdown menu is open, showing options like Pools, Configuration, Users, Connections, Variables (highlighted), and XComs. The main content area is titled 'Variables' and includes a file upload section with a 'Choose file' button and an 'Import Variables' button. Below this is a table of variables with columns for selection, key, value, and encryption status. The table contains three entries: 'access_token', 'app_secret', and 'business_id'. The 'business_id' value is '24' followed by redacted characters. The interface also shows a 'List (3)' button, a 'Create' button, an 'Add Filter' dropdown, and a 'With selected' dropdown.

Variables

Choose file No file chosen Import Variables

List (3) Create Add Filter With selected Search

<input type="checkbox"/>		Key	Val	Is Encrypted
<input type="checkbox"/>	 	access_token	*****	<input checked="" type="checkbox"/>
<input type="checkbox"/>	 	app_secret	*****	<input checked="" type="checkbox"/>
<input type="checkbox"/>	 	business_id	24 *****	<input checked="" type="checkbox"/>

XCom or Cross-communication:

- Let tasks **exchange** messages allowing **shared state**.
- Defined by a **key**, **value**, and **timestamp**.
- Also **track** attributes like the task/DAG that created the XCom and when it should become visible.
- **Any object** that can be **pickled** can be used as an XCom value.

XComs can be :

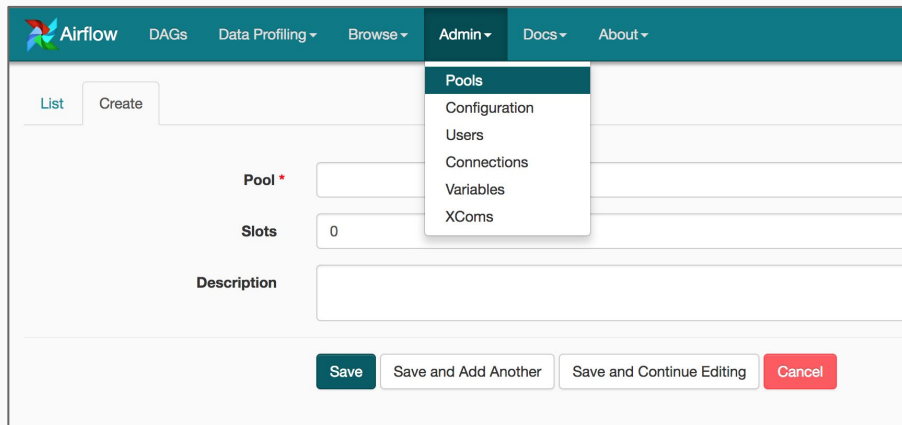
- Pushed (sent) :
 - Calling `xcom_push()`
 - If a task return a value (from its operator `execute()` method) or from a PythonOperator's `python_callable`
- Pulled (received) : calling `xcom_pull()`

SLA :

- Service Level Agreements, or time by which a task or DAG should have succeeded,
- Can be set at a task level as a timedelta.
- An alert email is sent detailing the list of tasks that missed their SLA.

Pools :

- Some systems can get overwhelmed when too many processes hit them at the same time.
- **Limit the execution parallelism** on arbitrary sets of tasks.



The screenshot displays the Apache Airflow Admin interface. The top navigation bar includes links for DAGs, Data Profiling, Browse, Admin, Docs, and About. The Admin menu is open, showing options for Pools, Configuration, Users, Connections, Variables, and XComs. The Pools page is in 'Create' mode, featuring a form with the following fields:

- Pool ***: A text input field.
- Slots**: A numeric input field with the value '0'.
- Description**: A text area.

At the bottom of the form are four buttons: 'Save' (dark green), 'Save and Add Another' (light green), 'Save and Continue Editing' (light green), and 'Cancel' (red).

Pools :

```
1 aggregate_db_message_job = BashOperator(  
2     task_id='aggregate_db_message_job',  
3     execution_timeout=timedelta(hours=3),  
4     pool='ep_data_pipeline_db_msg_agg',  
5     bash_command=aggregate_db_message_job_cmd,  
6     dag=dag)  
7 aggregate_db_message_job.set_upstream(wait_for_empty_queue_)
```

Queues : (only on CeleryExecutors) :

- Every Task can be assigned a specific queue name
- By default, both worker and tasks are assigned with the **default_queue** queue
- Workers can be assigned **multiple queues**
- Very useful feature when specialized workers are needed (GPU, Spark...)

Trigger Rules:


Though the normal workflow behavior is to trigger tasks when all their directly upstream tasks have succeeded, Airflow allows for more complex dependency settings.

All operators have a `trigger_rule` argument which defines the rule by which the generated task get triggered. The default value for `trigger_rule` is `all_success` and can be defined as “trigger this task when all directly upstream tasks have succeeded”. All other rules described here are based on direct parent tasks and are values that can be passed to any operator while creating tasks:

- **all_success:** (default) all parents have succeeded
- **all_failed:** all parents are in a failed or upstream_failed state
- **all_done:** all parents are done with their execution
- **one_failed:** fires as soon as at least one parent has failed, it does not wait for all parents to be done • **one_success:** fires as soon as at least one parent succeeds, it does not wait for all parents to be done • **dummy:** dependencies are just for show, trigger at will.

User Interface

Dags view :

 AirFlow

DAGs






















Tools ▾

Browse ▾

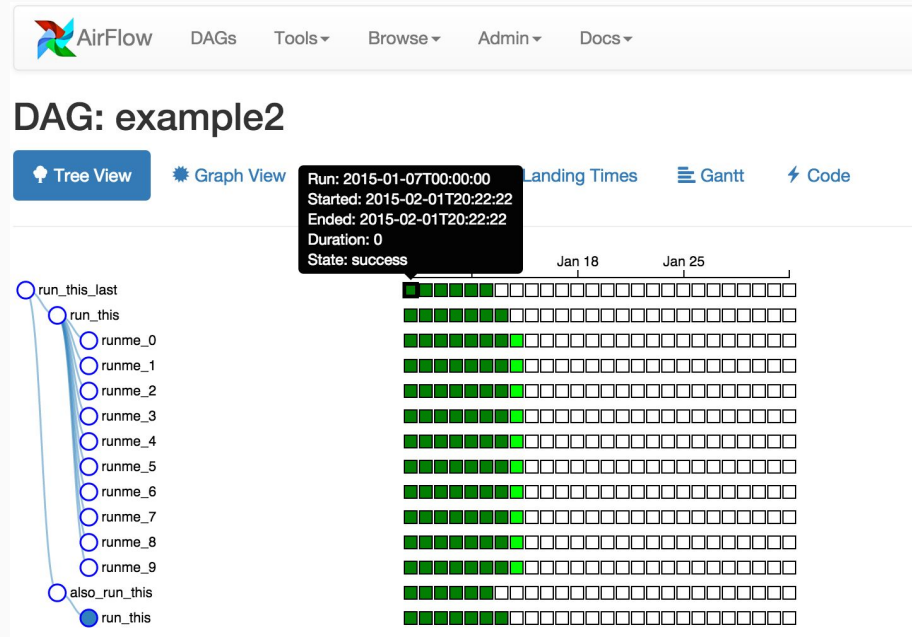
Admin ▾

Docs ▾

DAGs

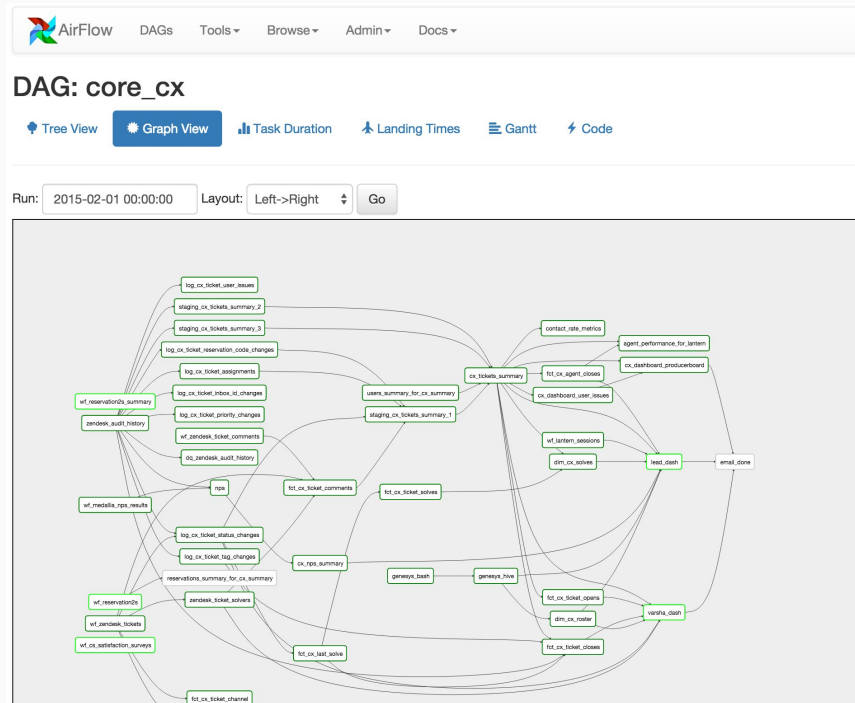
DAG	Filepath	Owner	Task by State	Links
example1	example_dags/example1.py	airflow	<div><div>80</div><div>1</div><div>0</div></div>	      
example2	example_dags/example2.py	airflow	<div><div>128</div><div>10</div><div>0</div></div>	      
example3	example_dags/example3.py	airflow	<div><div>138</div><div>5</div><div>0</div></div>	      

Tree view :

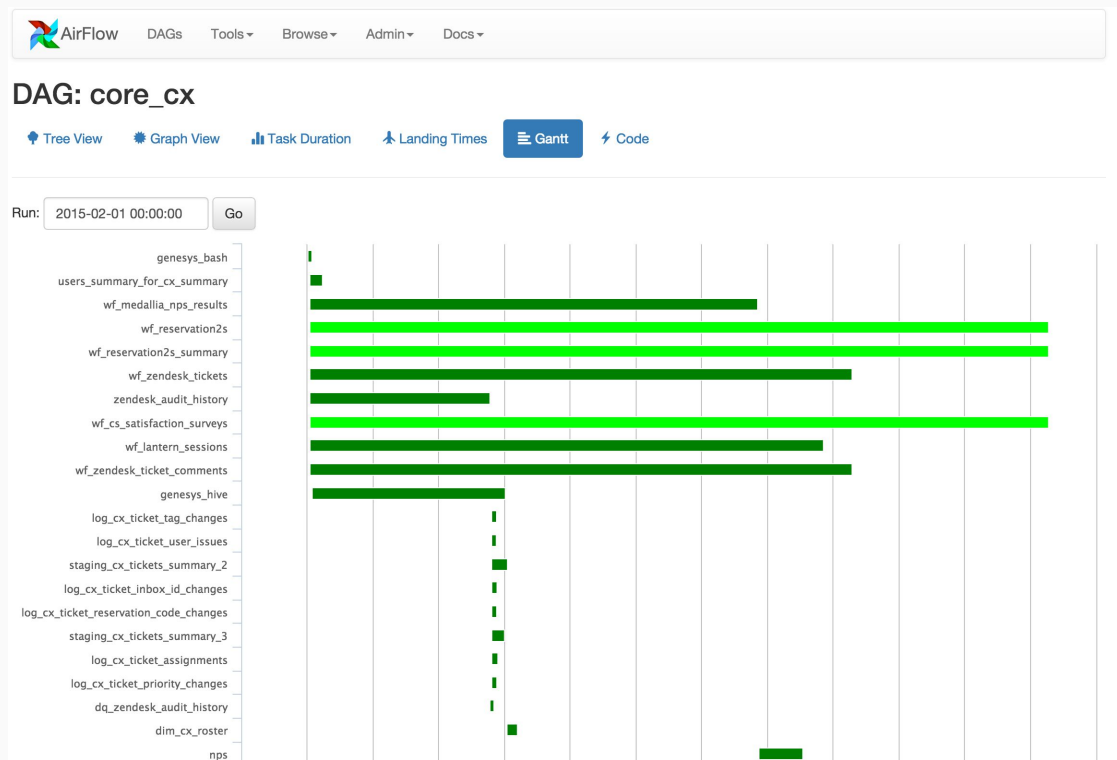


User Interface

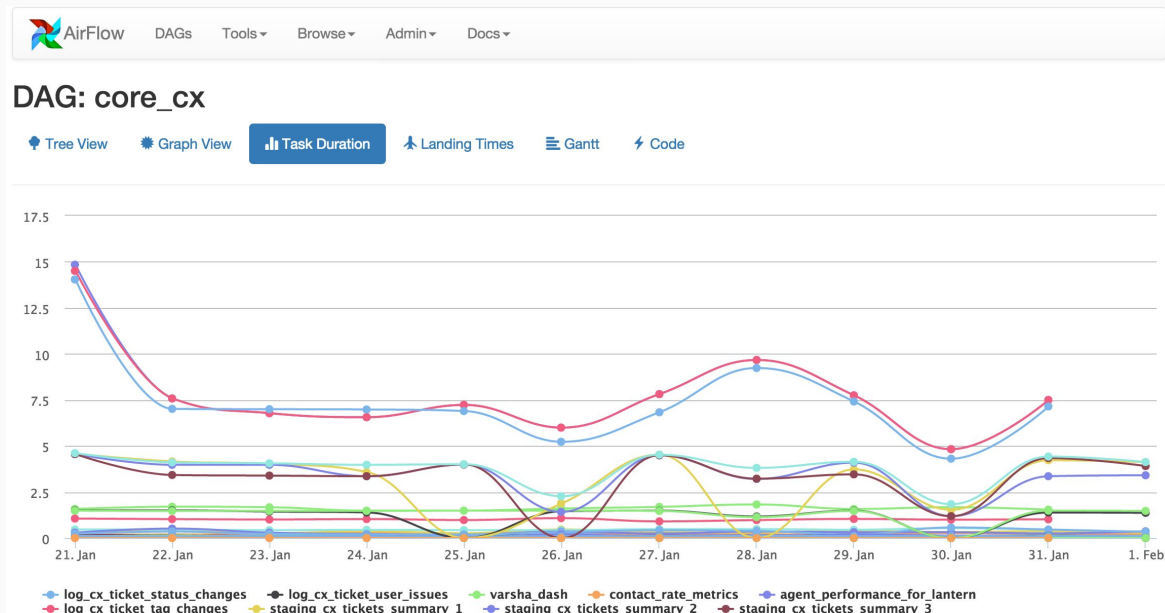
Graph view :





Gantt view :



Task duration :



Data Profiling : SQL Queries

 DAGs **Data Profiling** ▾ Browse ▾ Admin ▾ Docs ▾ About ▾ 05:04 UTC 

Ad Hoc Query

airflow_db ▾

Run!

.csv

1

```
select * from task_instance limit 10;
```

Show

100 ▾

entries

Search:

task_id	dag_id	execution_date	start_date	end_date	duration	state	try_number	hostname	unixname	job_id	pool	queue	priority_weight	operator
convert_to_csv_task	facebook_insights_dag	2018-03-20	2018-03-21 16:38:52.541027	2018-03-21 16:39:09.613790	17.072763	success	1	OKACHAs-MacBook-Pro.local	IOKACHA	452		default	1	PythonOperator
campaign_lifetime_insights_task	facebook_insights_dag	2018-03-17	2018-03-21 16:22:00.403601	2018-03-21 16:26:53.214537	292.810936	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	357		default	3	PythonOperator
campaign_lifetime_insights_task	facebook_insights_dag	2018-03-16	2018-03-21 16:22:00.964350	2018-03-21 16:27:11.234319	310.269969	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	359		default	3	PythonOperator
campaign_lifetime_insights_task	facebook_insights_dag	2018-03-19	2018-03-21 16:22:00.702286	2018-03-21 16:27:11.709757	311.007471	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	358		default	3	PythonOperator
campaigns_infos_task	facebook_insights_dag	2018-03-20	2018-03-21 16:20:51.355984	2018-03-21 16:21:31.989539	40.633555	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	330		default	4	PythonOperator
campaigns_infos_task	facebook_insights_dag	2018-03-18	2018-03-21 16:20:51.400999	2018-03-21 16:21:34.013295	42.612296	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	333		default	4	PythonOperator
campaign_lifetime_insights_task	facebook_insights_dag	2018-03-20	2018-03-21 16:22:01.067421	2018-03-21 16:27:13.030706	311.963285	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	360		default	3	PythonOperator
enrich_ads_daily_insights_task	facebook_insights_dag	2018-03-15	2018-03-21 16:22:53.110321	2018-03-21 16:23:13.671750	20.561429	success	1	OKACHAs-MacBook-Pro.local	IOKACHA	382		default	2	PythonOperator
campaign_lifetime_insights_task	facebook_insights_dag	2018-03-14	2018-03-21 16:23:25.384377	2018-03-21 16:28:03.737195	278.352818	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	401		default	3	PythonOperator
campaign_lifetime_insights_task	facebook_insights_dag	2018-03-15	2018-03-21 16:23:52.229709	2018-03-21 16:28:07.791497	255.561788	success	5	OKACHAs-MacBook-Pro.local	IOKACHA	419		default	3	PythonOperator


Showing 1 to 10 of 10 entries

Previous

1

Next

Data Profiling : Charts

 DAGs


Data Profiling ▾

Browse ▾

Admin ▾

Docs ▾

About ▾

05:09 UTC 

List>Create>Edit

Label

label

Can include {{ templated_fields }} and {{ macros }}

Owner

The chart's owner, mostly used for reference and filtering in the list view.

Conn Id ▾

airflow_db

Source database to run the query against

Chart Type

Bar Chart

The type of chart to be displayed

Show Datatable

☐

Whether to display an interactive data table under the chart.

X Is Date

☒

Whether the X axis should be casted as a date field. Expect most intelligible date formats to get casted properly.

Y Log Scale

☐

Whether to use a log scale for the Y axis.

Display the SQL Statement

☒

Whether to display the SQL statement as a collapsible section in the chart page.

Chart Height

600

Height of the chart, in pixels.

SQL Layout

SELECT series, x, y FROM ...

Defines the layout of the SQL that the application should expect. Depending on the tables you are sourcing from, it may make more sense to pivot / unpivot the metrics.

SQL

Can include {{ templated_fields }} and {{ macros }}.

```
1 select dag_id, task_id, sum(try_number) from task_instance group by dag_id , task_id ;
```

Default Parameters

{}

A dictionary of {"key": "values,"} that define what the templated fields (parameters) values should be by default. To be valid, it needs to "eval" as a Python dict. The key values will show up in the url's querystring and can be altered there.

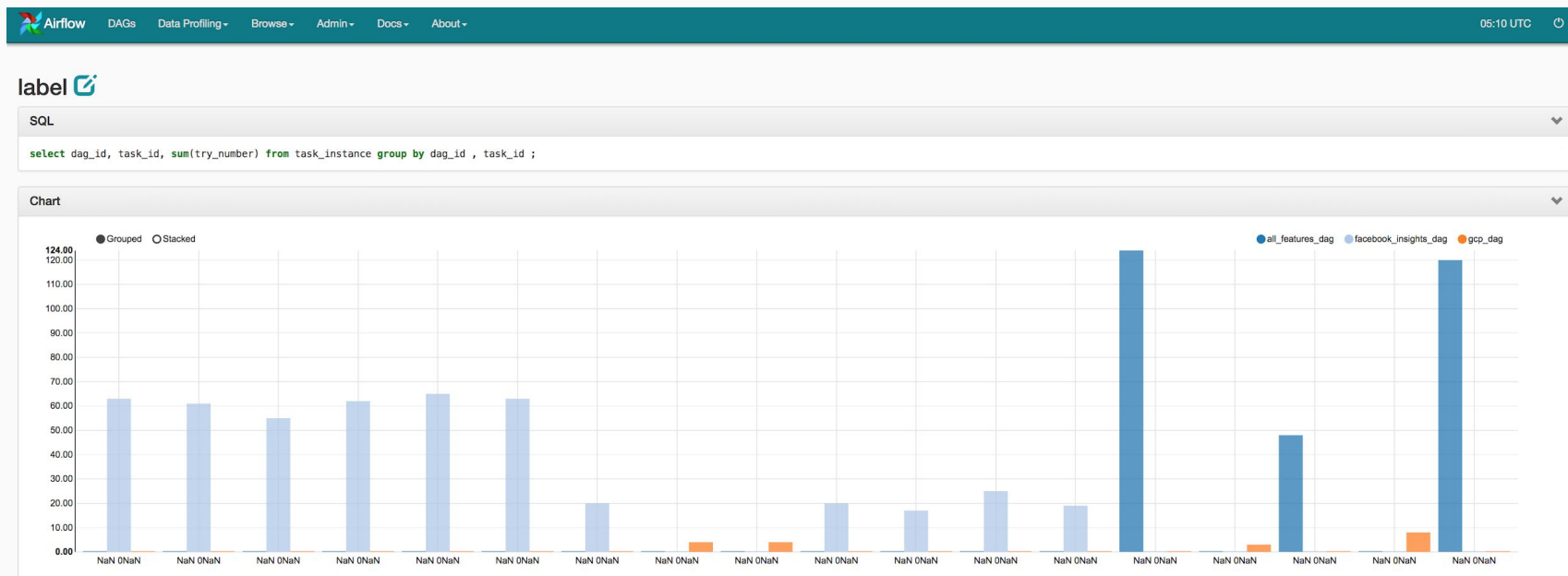
Save

Save and Add Another

Save and Continue Editing

Cancel

Data Profiling : Charts



CLI

```
airflow variables [-h] [-s KEY VAL] [-g KEY] [-j] [-d VAL] [-i FILEPATH] [-e FILEPATH] [-x KEY]
```

```
airflow connections [-h] [-l] [-a] [-d] [--conn_id CONN_ID]
                    [--conn_uri CONN_URI] [--conn_extra CONN_EXTRA]
                    [--conn_type CONN_TYPE] [--conn_host CONN_HOST]
                    [--conn_login CONN_LOGIN] [--conn_password CONN_PASSWORD]
                    [--conn_schema CONN_SCHEMA] [--conn_port CONN_PORT]
```

```
airflow pause [-h] [-sd SUBDIR] dag_id
```

```
airflow test [-h] [-sd SUBDIR] [-dr] [-tp TASK_PARAMS] dag_id task_id execution_date
```

```
airflow backfill dag_id task_id -s START_DATE -e END_DATE
```

```
airflow clear DAG_ID
```

```
airflow resetdb [-h] [-y]
```

Security

Security

By default : all access are open

Support ;

- Web authentication with :
 - Password
 - LDAP
 - Custom auth
 - Kerberos
 - OAuth
 - Github Enterprise Authentication
 - Google Authentication
- Impersonation (run as other \$USER)
- Secure access via SSL

Demo

Demo

1. **Facebook Ads insights** data pipeline.
2. Run a pyspark script on a ephemeral **dataproc** cluster only when s3 data input is available
3. Useless workflow : Hook + Connection + Operators + Sensors + XCom +(SLA):
 - List s3 files (hooks)
 - Share state with the next task (xcom)
 - Write content to s3 (hooks)
 - Resume the workflow when an S3 DONE.FLAG file is ready (sensor)

Resources

<https://airflow.apache.org>

<http://www.clairvoyantsoft.com/assets/whitepapers/GuideToApacheAirflow.pdf>

<https://speakerdeck.com/artwr/apache-airflow-at-airbnb-introduction-and-lessons-learned>

<https://www.slideshare.net/sumitmaheshwari007/apache-airflow>

Thanks