
Airflow Documentation

Release

Maxime Beauchemin

May 18, 2017

1	Principles	3
2	Beyond the Horizon	5
3	Content	7
3.1	Project	7
3.1.1	History	7
3.1.2	Committers	7
3.1.3	Resources & links	7
3.1.4	Roadmap	8
3.2	License	8
3.3	Quick Start	12
3.3.1	What's Next?	12
3.4	Installation	12
3.4.1	Getting Airflow	12
3.4.2	Extra Packages	13
3.5	Tutorial	15
3.5.1	Example Pipeline definition	15
3.5.2	It's a DAG definition file	16
3.5.3	Importing Modules	16
3.5.4	Default Arguments	16
3.5.5	Instantiate a DAG	17
3.5.6	Tasks	17
3.5.7	Templating with Jinja	17
3.5.8	Setting up Dependencies	18
3.5.9	Recap	19
3.5.10	Testing	20
3.5.10.1	Running the Script	20
3.5.10.2	Command Line Metadata Validation	20
3.5.10.3	Testing	20
3.5.10.4	Backfill	21
3.5.11	What's Next?	21
3.6	Configuration	21
3.6.1	Setting Configuration Options	21
3.6.2	Setting up a Backend	22
3.6.3	Connections	22
3.6.4	Scaling Out with Celery	23

3.6.5	Scaling Out with Dask	24
3.6.6	Logs	24
3.6.7	Scaling Out on Mesos (community contributed)	25
3.6.8	Integration with systemd	25
3.6.9	Integration with upstart	26
3.6.10	Test Mode	26
3.7	UI / Screenshots	26
3.7.1	DAGs View	26
3.7.2	Tree View	27
3.7.3	Graph View	28
3.7.4	Variable View	28
3.7.5	Gantt Chart	29
3.7.6	Task Duration	30
3.7.7	Code View	31
3.7.8	Task Instance Context Menu	32
3.8	Concepts	33
3.8.1	Core Ideas	33
3.8.1.1	DAGs	33
3.8.1.2	Operators	34
3.8.1.3	Tasks	36
3.8.1.4	Task Instances	36
3.8.1.5	Workflows	36
3.8.2	Additional Functionality	37
3.8.2.1	Hooks	37
3.8.2.2	Pools	37
3.8.2.3	Connections	37
3.8.2.4	Queues	38
3.8.2.5	XComs	38
3.8.2.6	Variables	39
3.8.2.7	Branching	39
3.8.2.8	SubDAGs	40
3.8.2.9	SLAs	42
3.8.2.10	Trigger Rules	42
3.8.2.11	Latest Run Only	42
3.8.2.12	Zombies & Undeads	43
3.8.2.13	Cluster Policy	44
3.8.2.14	Documentation & Notes	44
3.8.2.15	Jinja Templating	45
3.8.3	Packaged dags	45
3.9	Data Profiling	46
3.9.1	Adhoc Queries	46
3.9.2	Charts	47
3.9.2.1	Chart Screenshot	48
3.9.2.2	Chart Form Screenshot	49
3.10	Command Line Interface	49
3.10.1	Positional Arguments	49
3.10.2	Sub-commands:	50
3.10.2.1	resetdb	50
3.10.2.2	render	50
3.10.2.3	variables	50
3.10.2.4	connections	51
3.10.2.5	pause	51
3.10.2.6	task_failed_deps	52
3.10.2.7	version	52

3.10.2.8	trigger_dag	52
3.10.2.9	initdb	53
3.10.2.10	test	53
3.10.2.11	unpause	53
3.10.2.12	dag_state	54
3.10.2.13	run	54
3.10.2.14	list_tasks	55
3.10.2.15	backfill	56
3.10.2.16	list_dags	57
3.10.2.17	kerberos	57
3.10.2.18	worker	57
3.10.2.19	webserver	58
3.10.2.20	flower	59
3.10.2.21	scheduler	60
3.10.2.22	task_state	60
3.10.2.23	pool	61
3.10.2.24	serve_logs	61
3.10.2.25	clear	61
3.10.2.26	upgradedb	62
3.11	Scheduling & Triggers	62
3.11.1	DAG Runs	62
3.11.2	Backfill and Catchup	63
3.11.3	External Triggers	64
3.11.4	To Keep in Mind	64
3.12	Plugins	64
3.12.1	What for?	64
3.12.2	Why build on top of Airflow?	65
3.12.3	Interface	65
3.12.4	Example	65
3.13	Security	67
3.13.1	Web Authentication	67
3.13.1.1	Password	67
3.13.1.2	LDAP	67
3.13.1.3	Roll your own	68
3.13.2	Multi-tenancy	68
3.13.3	Kerberos	69
3.13.3.1	Limitations	69
3.13.3.2	Enabling kerberos	69
3.13.3.3	Using kerberos authentication	70
3.13.4	OAuth Authentication	70
3.13.4.1	GitHub Enterprise (GHE) Authentication	70
3.13.4.2	Google Authentication	71
3.13.5	SSL	71
3.13.6	Impersonation	72
3.13.6.1	Default Impersonation	72
3.14	Experimental Rest API	72
3.14.1	Endpoints	72
3.14.2	CLI	72
3.14.3	Authentication	73
3.15	Integration	73
3.15.1	Azure: Microsoft Azure	73
3.15.1.1	Azure Blob Storage	73
3.15.2	AWS: Amazon Webservies	73
3.15.3	Databricks	74

3.15.3.1	DatabricksSubmitRunOperator	74
3.15.4	GCP: Google Cloud Platform	76
3.15.4.1	BigQuery	76
3.15.4.2	Cloud DataFlow	81
3.15.4.3	Cloud DataProc	83
3.15.4.4	Cloud Datastore	86
3.15.4.5	Cloud Storage	87
3.16	FAQ	89
3.16.1	Why isn't my task getting scheduled?	89
3.16.2	How do I trigger tasks based on another task's failure?	90
3.16.3	Why are connection passwords still not encrypted in the metadata db after I installed air-flow[crypto]?	90
3.16.4	What's the deal with <code>start_date</code> ?	90
3.16.5	How can I create DAGs dynamically?	91
3.16.6	What are all the <code>airflow run</code> commands in my process list?	91
3.17	API Reference	91
3.17.1	Operators	91
3.17.1.1	BaseOperator	92
3.17.1.2	BaseSensorOperator	94
3.17.1.3	Operator API	94
3.17.1.4	Community-contributed Operators	101
3.17.2	Macros	107
3.17.2.1	Default Variables	107
3.17.2.2	Macros	109
3.17.3	Models	110
3.17.4	Hooks	120
3.17.4.1	Community contributed hooks	122
3.17.5	Executors	126
3.17.5.1	Community-contributed executors	127



Important: Disclaimer: Apache Airflow is an effort undergoing incubation at The Apache Software Foundation (ASF), sponsored by the Apache Incubator. Incubation is required of all newly accepted projects until a further review indicates that the infrastructure, communications, and decision making process have stabilized in a manner consistent with other successful ASF projects. While incubation status is not necessarily a reflection of the completeness or stability of the code, it does indicate that the project has yet to be fully endorsed by the ASF.

Airflow is a platform to programmatically author, schedule and monitor workflows.

Use airflow to author workflows as directed acyclic graphs (DAGs) of tasks. The airflow scheduler executes your tasks on an array of workers while following the specified dependencies. Rich command line utilities make performing complex surgeries on DAGs a snap. The rich user interface makes it easy to visualize pipelines running in production, monitor progress, and troubleshoot issues when needed.

When workflows are defined as code, they become more maintainable, versionable, testable, and collaborative.

CHAPTER 1

Principles

- **Dynamic:** Airflow pipelines are configuration as code (Python), allowing for dynamic pipeline generation. This allows for writing code that instantiates pipelines dynamically.
- **Extensible:** Easily define your own operators, executors and extend the library so that it fits the level of abstraction that suits your environment.
- **Elegant:** Airflow pipelines are lean and explicit. Parameterizing your scripts is built into the core of Airflow using the powerful **Jinja** templating engine.
- **Scalable:** Airflow has a modular architecture and uses a message queue to orchestrate an arbitrary number of workers. Airflow is ready to scale to infinity.

CHAPTER 2

Beyond the Horizon

Airflow **is not** a data streaming solution. Tasks do not move data from one to the other (though tasks can exchange metadata!). Airflow is not in the [Spark Streaming](#) or [Storm](#) space, it is more comparable to [Oozie](#) or [Azkaban](#).

Workflows are expected to be mostly static or slowly changing. You can think of the structure of the tasks in your workflow as slightly more dynamic than a database structure would be. Airflow workflows are expected to look similar from a run to the next, this allows for clarity around unit of work and continuity.

3.1 Project

3.1.1 History

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 incubation program in March 2016.

3.1.2 Committers

- @mistercrunch (Maxime “Max” Beauchemin)
- @r39132 (Siddharth “Sid” Anand)
- @criccomini (Chris Riccomini)
- @bolkedebuin (Bolke de Bruin)
- @artwr (Arthur Wiedmer)
- @jlowin (Jeremiah Lowin)
- @patrickleotardif (Patrick Leo Tardif)
- @aoen (Dan Davydov)
- @syvineckruiyk (Steven Yvinec-Kruiyk)

For the full list of contributors, take a look at [Airflow's Github Contributor page](#):

3.1.3 Resources & links

- [Airflow's official documentation](#)

- Mailing list (send emails to dev-subscribe@airflow.incubator.apache.org and/or commits-subscribe@airflow.incubator.apache.org to subscribe to each)
- [Issues on Apache's Jira](#)
- [Gitter \(chat\) Channel](#)
- [More resources and links to Airflow related content on the Wiki](#)

3.1.4 Roadmap

Please refer to the Roadmap on [the wiki](#)

3.2 License



Apache License
Version 2.0, January 2004
<http://www.apache.org/licenses/>

TERMS AND CONDITIONS FOR USE, REPRODUCTION, AND DISTRIBUTION

1. Definitions.

"License" shall mean the terms and conditions for use, reproduction, and distribution as defined by Sections 1 through 9 of this document.

"Licensor" shall mean the copyright owner or entity authorized by the copyright owner that is granting the License.

"Legal Entity" shall mean the union of the acting entity and all other entities that control, are controlled by, or are under common control with that entity. For the purposes of this definition, "control" means (i) the power, direct or indirect, to cause the direction or management of such entity, whether by contract or otherwise, or (ii) ownership of fifty percent (50%) or more of the outstanding shares, or (iii) beneficial ownership of such entity.

"You" (or "Your") shall mean an individual or Legal Entity exercising permissions granted by this License.

"Source" form shall mean the preferred form for making modifications, including but not limited to software source code, documentation source, and configuration files.

"Object" form shall mean any form resulting from mechanical transformation or translation of a Source form, including but not limited to compiled object code, generated documentation, and conversions to other media types.

"Work" shall mean the work of authorship, whether in Source or Object form, made available under the License, as indicated by a copyright notice that is included in or attached to the work (an example is provided in the Appendix below).

"Derivative Works" shall mean any work, whether in Source or Object form, that is based on (or derived from) the Work and for which the editorial revisions, annotations, elaborations, or other modifications represent, as a whole, an original work of authorship. For the purposes of this License, Derivative Works shall not include works that remain separable from, or merely link (or bind by name) to the interfaces of, the Work and Derivative Works thereof.

"Contribution" shall mean any work of authorship, including the original version of the Work and any modifications or additions to that Work or Derivative Works thereof, that is intentionally submitted to Licensor for inclusion in the Work by the copyright owner or by an individual or Legal Entity authorized to submit on behalf of the copyright owner. For the purposes of this definition, "submitted" means any form of electronic, verbal, or written communication sent to the Licensor or its representatives, including but not limited to communication on electronic mailing lists, source code control systems, and issue tracking systems that are managed by, or on behalf of, the Licensor for the purpose of discussing and improving the Work, but excluding communication that is conspicuously marked or otherwise designated in writing by the copyright owner as "Not a Contribution."

"Contributor" shall mean Licensor and any individual or Legal Entity on behalf of whom a Contribution has been received by Licensor and subsequently incorporated within the Work.

2. Grant of Copyright License. Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable copyright license to reproduce, prepare Derivative Works of, publicly display, publicly perform, sublicense, and distribute the Work and such Derivative Works in Source or Object form.
3. Grant of Patent License. Subject to the terms and conditions of this License, each Contributor hereby grants to You a perpetual, worldwide, non-exclusive, no-charge, royalty-free, irrevocable (except as stated in this section) patent license to make, have made, use, offer to sell, sell, import, and otherwise transfer the Work, where such license applies only to those patent claims licensable by such Contributor that are necessarily infringed by their Contribution(s) alone or by combination of their Contribution(s) with the Work to which such Contribution(s) was submitted. If You institute patent litigation against any entity (including a cross-claim or counterclaim in a lawsuit) alleging that the Work or a Contribution incorporated within the Work constitutes direct or contributory patent infringement, then any patent licenses granted to You under this License for that Work shall terminate as of the date such litigation is filed.
4. Redistribution. You may reproduce and distribute copies of the Work or Derivative Works thereof in any medium, with or without modifications, and in Source or Object form, provided that You meet the following conditions:

- (a) You must give any other recipients of the Work or Derivative Works a copy of this License; and
- (b) You must cause any modified files to carry prominent notices stating that You changed the files; and
- (c) You must retain, in the Source form of any Derivative Works that You distribute, all copyright, patent, trademark, and attribution notices from the Source form of the Work, excluding those notices that do not pertain to any part of the Derivative Works; and
- (d) If the Work includes a "NOTICE" text file as part of its distribution, then any Derivative Works that You distribute must include a readable copy of the attribution notices contained within such NOTICE file, excluding those notices that do not pertain to any part of the Derivative Works, in at least one of the following places: within a NOTICE text file distributed as part of the Derivative Works; within the Source form or documentation, if provided along with the Derivative Works; or, within a display generated by the Derivative Works, if and wherever such third-party notices normally appear. The contents of the NOTICE file are for informational purposes only and do not modify the License. You may add Your own attribution notices within Derivative Works that You distribute, alongside or as an addendum to the NOTICE text from the Work, provided that such additional attribution notices cannot be construed as modifying the License.

You may add Your own copyright statement to Your modifications and may provide additional or different license terms and conditions for use, reproduction, or distribution of Your modifications, or for any such Derivative Works as a whole, provided Your use, reproduction, and distribution of the Work otherwise complies with the conditions stated in this License.

- 5. Submission of Contributions. Unless You explicitly state otherwise, any Contribution intentionally submitted for inclusion in the Work by You to the Licensor shall be under the terms and conditions of this License, without any additional terms or conditions. Notwithstanding the above, nothing herein shall supersede or modify the terms of any separate license agreement you may have executed with Licensor regarding such Contributions.
- 6. Trademarks. This License does not grant permission to use the trade names, trademarks, service marks, or product names of the Licensor, except as required for reasonable and customary use in describing the origin of the Work and reproducing the content of the NOTICE file.
- 7. Disclaimer of Warranty. Unless required by applicable law or agreed to in writing, Licensor provides the Work (and each Contributor provides its Contributions) on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied, including, without limitation, any warranties or conditions of TITLE, NON-INFRINGEMENT, MERCHANTABILITY, or FITNESS FOR A PARTICULAR PURPOSE. You are solely responsible for determining the appropriateness of using or redistributing the Work and assume any

risks associated with Your exercise of permissions under this License.

8. **Limitation of Liability.** In no event and under no legal theory, whether in tort (including negligence), contract, or otherwise, unless required by applicable law (such as deliberate and grossly negligent acts) or agreed to in writing, shall any Contributor be liable to You for damages, including any direct, indirect, special, incidental, or consequential damages of any character arising as a result of this License or out of the use or inability to use the Work (including but not limited to damages for loss of goodwill, work stoppage, computer failure or malfunction, or any and all other commercial damages or losses), even if such Contributor has been advised of the possibility of such damages.
9. **Accepting Warranty or Additional Liability.** While redistributing the Work or Derivative Works thereof, You may choose to offer, and charge a fee for, acceptance of support, warranty, indemnity, or other liability obligations and/or rights consistent with this License. However, in accepting such obligations, You may act only on Your own behalf and on Your sole responsibility, not on behalf of any other Contributor, and only if You agree to indemnify, defend, and hold each Contributor harmless for any liability incurred by, or claims asserted against, such Contributor by reason of your accepting any such warranty or additional liability.

END OF TERMS AND CONDITIONS

APPENDIX: How to apply the Apache License to your work.

To apply the Apache License to your work, attach the following boilerplate notice, with the fields enclosed by brackets "[]" replaced with your own identifying information. (Don't include the brackets!) The text should be enclosed in the appropriate comment syntax for the file format. We also recommend that a file or class name and description of purpose be included on the same "printed page" as the copyright notice for easier identification within third-party archives.

Copyright 2015 Apache Software Foundation

Licensed under the Apache License, Version 2.0 (the "License");
you may not use this file except in compliance with the License.
You may obtain a copy of the License at

<http://www.apache.org/licenses/LICENSE-2.0>

Unless required by applicable law or agreed to in writing, software distributed under the License is distributed on an "AS IS" BASIS, WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied. See the License for the specific language governing permissions and limitations under the License.

Status API Training Shop Blog About

© 2016 GitHub, Inc. Terms Privacy Security Contact Help

3.3 Quick Start

The installation is quick and straightforward.

```
# airflow needs a home, ~/airflow is the default,  
# but you can lay foundation somewhere else if you prefer  
# (optional)  
export AIRFLOW_HOME=~/airflow  
  
# install from pypi using pip  
pip install apache-airflow  
  
# initialize the database  
airflow initdb  
  
# start the web server, default port is 8080  
airflow webserver -p 8080
```

Upon running these commands, Airflow will create the `$AIRFLOW_HOME` folder and lay an “airflow.cfg” file with defaults that get you going fast. You can inspect the file either in `$AIRFLOW_HOME/airflow.cfg`, or through the UI in the Admin->Configuration menu. The PID file for the webserver will be stored in `$AIRFLOW_HOME/airflow-webserver.pid` or in `/run/airflow/webserver.pid` if started by systemd.

Out of the box, Airflow uses a sqlite database, which you should outgrow fairly quickly since no parallelization is possible using this database backend. It works in conjunction with the `SequentialExecutor` which will only run task instances sequentially. While this is very limiting, it allows you to get up and running quickly and take a tour of the UI and the command line utilities.

Here are a few commands that will trigger a few task instances. You should be able to see the status of the jobs change in the `example1` DAG as you run the commands below.

```
# run your first task instance  
airflow run example_bash_operator runme_0 2015-01-01  
# run a backfill over 2 days  
airflow backfill example_bash_operator -s 2015-01-01 -e 2015-01-02
```

3.3.1 What’s Next?

From this point, you can head to the [Tutorial](#) section for further examples or the [Configuration](#) section if you’re ready to get your hands dirty.

3.4 Installation

3.4.1 Getting Airflow

The easiest way to install the latest stable version of Airflow is with `pip`:

```
pip install apache-airflow
```

You can also install Airflow with support for extra features like `s3` or `postgres`:

```
pip install "apache-airflow[s3, postgres]"
```

3.4.2 Extra Packages

The `apache-airflow` PyPI basic package only installs what's needed to get started. Subpackages can be installed depending on what will be useful in your environment. For instance, if you don't need connectivity with Postgres, you won't have to go through the trouble of installing the `postgres-devel` yum package, or whatever equivalent applies on the distribution you are using.

Behind the scenes, Airflow does conditional imports of operators that require these extra dependencies.

Here's the list of the subpackages and what they enable:

sub-pack-age	install command	enables
all	<code>pip install apache-airflow[all]</code>	All Airflow features known to man
all_dbs	<code>pip install apache-airflow[all_dbs]</code>	All databases integrations
async	<code>pip install apache-airflow[async]</code>	Async worker classes for gunicorn
devel	<code>pip install apache-airflow[devel]</code>	Minimum dev tools requirements
devel_hadoop	<code>pip install apache-airflow[devel_hadoop]</code>	Airflow + dependencies on the Hadoop stack
celery	<code>pip install apache-airflow[celery]</code>	CeleryExecutor
crypto	<code>pip install apache-airflow[crypto]</code>	Encrypt connection passwords in metadata db
druid	<code>pip install apache-airflow[druid]</code>	Druid.io related operators & hooks
gcp_api	<code>pip install apache-airflow[gcp_api]</code>	Google Cloud Platform hooks and operators (using google-api-python-client)
jdbc	<code>pip install apache-airflow[jdbc]</code>	JDBC hooks and operators
hdfs	<code>pip install apache-airflow[hdfs]</code>	HDFS hooks and operators
hive	<code>pip install apache-airflow[hive]</code>	All Hive related operators
kerberos	<code>pip install apache-airflow[kerberos]</code>	kerberos integration for kerberized hadoop
ldap	<code>pip install apache-airflow[ldap]</code>	ldap authentication for users
mssql	<code>pip install apache-airflow[mssql]</code>	Microsoft SQL operators and hook, support as an Airflow backend
mysql	<code>pip install apache-airflow[mysql]</code>	MySQL operators and hook, support as an Airflow backend
password	<code>pip install apache-airflow[password]</code>	Password Authentication for users
postgres	<code>pip install apache-airflow[postgres]</code>	Postgres operators and hook, support as an Airflow backend
qds	<code>pip install apache-airflow[qds]</code>	Enable QDS (qubole data services) support
rabbitmq	<code>pip install apache-airflow[rabbitmq]</code>	Rabbitmq support as a Celery backend
s3	<code>pip install apache-airflow[s3]</code>	S3KeySensor, S3PrefixSensor
samba	<code>pip install apache-airflow[samba]</code>	Hive2SambaOperator
slack	<code>pip install apache-airflow[slack]</code>	SlackAPIPostOperator
vertica	<code>pip install apache-airflow[vertica]</code>	Vertica hook support as an Airflow backend
cloudant	<code>pip install apache-airflow[cloudant]</code>	Cloudant hook
redis	<code>pip install apache-airflow[redis]</code>	Redis hooks and sensors
14	<code>apache-airflow[redis]</code>	Chapter 3. Content

3.5 Tutorial

This tutorial walks you through some of the fundamental Airflow concepts, objects, and their usage while writing your first pipeline.

3.5.1 Example Pipeline definition

Here is an example of a basic pipeline definition. Do not worry if this looks complicated, a line by line explanation follows below.

```
"""
Code that goes along with the Airflow tutorial located at:
https://github.com/airbnb/airflow/blob/master/airflow/example_dags/tutorial.py
"""
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
from datetime import datetime, timedelta

default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2015, 6, 1),
    'email': ['airflow@example.com'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5),
    # 'queue': 'bash_queue',
    # 'pool': 'backfill',
    # 'priority_weight': 10,
    # 'end_date': datetime(2016, 1, 1),
}

dag = DAG('tutorial', default_args=default_args)

# t1, t2 and t3 are examples of tasks created by instantiating operators
t1 = BashOperator(
    task_id='print_date',
    bash_command='date',
    dag=dag)

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

templated_command = """
{% for i in range(5) %}
    echo "{{ ds }}"
    echo "{{ macros.ds_add(ds, 7) }}"
    echo "{{ params.my_param }}"
{% endfor %}
"""
```

```
t3 = BashOperator(
    task_id='templated',
    bash_command=templated_command,
    params={'my_param': 'Parameter I passed in'},
    dag=dag)

t2.set_upstream(t1)
t3.set_upstream(t1)
```

3.5.2 It's a DAG definition file

One thing to wrap your head around (it may not be very intuitive for everyone at first) is that this Airflow Python script is really just a configuration file specifying the DAG's structure as code. The actual tasks defined here will run in a different context from the context of this script. Different tasks run on different workers at different points in time, which means that this script cannot be used to cross communicate between tasks. Note that for this purpose we have a more advanced feature called XCom.

People sometimes think of the DAG definition file as a place where they can do some actual data processing - that is not the case at all! The script's purpose is to define a DAG object. It needs to evaluate quickly (seconds, not minutes) since the scheduler will execute it periodically to reflect the changes if any.

3.5.3 Importing Modules

An Airflow pipeline is just a Python script that happens to define an Airflow DAG object. Let's start by importing the libraries we will need.

```
# The DAG object; we'll need this to instantiate a DAG
from airflow import DAG

# Operators; we need this to operate!
from airflow.operators.bash_operator import BashOperator
```

3.5.4 Default Arguments

We're about to create a DAG and some tasks, and we have the choice to explicitly pass a set of arguments to each task's constructor (which would become redundant), or (better!) we can define a dictionary of default parameters that we can use when creating tasks.

```
from datetime import datetime, timedelta

default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2015, 6, 1),
    'email': ['airflow@example.com'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5),
    # 'queue': 'bash_queue',
    # 'pool': 'backfill',
    # 'priority_weight': 10,
```

```
# 'end_date': datetime(2016, 1, 1),
}
```

For more information about the `BaseOperator`'s parameters and what they do, refer to the `:py:class:airflow.models.BaseOperator` documentation.

Also, note that you could easily define different sets of arguments that would serve different purposes. An example of that would be to have different settings between a production and development environment.

3.5.5 Instantiate a DAG

We'll need a DAG object to nest our tasks into. Here we pass a string that defines the `dag_id`, which serves as a unique identifier for your DAG. We also pass the default argument dictionary that we just defined and define a `schedule_interval` of 1 day for the DAG.

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

3.5.6 Tasks

Tasks are generated when instantiating operator objects. An object instantiated from an operator is called a constructor. The first argument `task_id` acts as a unique identifier for the task.

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

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

Notice how we pass a mix of operator specific arguments (`bash_command`) and an argument common to all operators (`retries`) inherited from `BaseOperator` to the operator's constructor. This is simpler than passing every argument for every constructor call. Also, notice that in the second task we override the `retries` parameter with 3.

The precedence rules for a task are as follows:

1. Explicitly passed arguments
2. Values that exist in the `default_args` dictionary
3. The operator's default value, if one exists

A task must include or inherit the arguments `task_id` and `owner`, otherwise Airflow will raise an exception.

3.5.7 Templating with Jinja

Airflow leverages the power of [Jinja Templating](#) and provides the pipeline author with a set of built-in parameters and macros. Airflow also provides hooks for the pipeline author to define their own parameters, macros and templates.

This tutorial barely scratches the surface of what you can do with templating in Airflow, but the goal of this section is to let you know this feature exists, get you familiar with double curly brackets, and point to the most common template variable: `{{ ds }}`.

```
templated_command = """
    {% for i in range(5) %}
        echo "{{ ds }}"
        echo "{{ macros.ds_add(ds, 7) }}"
        echo "{{ params.my_param }}"
    {% endfor %}
"""

t3 = BashOperator(
    task_id='templated',
    bash_command=templated_command,
    params={'my_param': 'Parameter I passed in'},
    dag=dag)
```

Notice that the `templated_command` contains code logic in `{% %}` blocks, references parameters like `{{ ds }}`, calls a function as in `{{ macros.ds_add(ds, 7) }}`, and references a user-defined parameter in `{{ params.my_param }}`.

The `params` hook in `BaseOperator` allows you to pass a dictionary of parameters and/or objects to your templates. Please take the time to understand how the parameter `my_param` makes it through to the template.

Files can also be passed to the `bash_command` argument, like `bash_command='templated_command.sh'`, where the file location is relative to the directory containing the pipeline file (`tutorial.py` in this case). This may be desirable for many reasons, like separating your script's logic and pipeline code, allowing for proper code highlighting in files composed in different languages, and general flexibility in structuring pipelines. It is also possible to define your `template_searchpath` as pointing to any folder locations in the DAG constructor call.

Using that same DAG constructor call, it is possible to define `user_defined_macros` which allow you to specify your own variables. For example, passing `dict(foo='bar')` to this argument allows you to use `{{ foo }}` in your templates. Moreover, specifying `user_defined_filters` allow you to register you own filters. For example, passing `dict(hello=lambda name: 'Hello %s' % name)` to this argument allows you to use `{{ 'world' | hello }}` in your templates. For more information regarding custom filters have a look at the [Jinja Documentation](#)

For more information on the variables and macros that can be referenced in templates, make sure to read through the [Macros](#) section

3.5.8 Setting up Dependencies

We have two simple tasks that do not depend on each other. Here's a few ways you can define dependencies between them:

```
t2.set_upstream(t1)

# This means that t2 will depend on t1
# running successfully to run
# It is equivalent to
# t1.set_downstream(t2)

t3.set_upstream(t1)

# all of this is equivalent to
# dag.set_dependency('print_date', 'sleep')
# dag.set_dependency('print_date', 'templated')
```

Note that when executing your script, Airflow will raise exceptions when it finds cycles in your DAG or when a dependency is referenced more than once.

3.5.9 Recap

Alright, so we have a pretty basic DAG. At this point your code should look something like this:

```
"""
Code that goes along with the Airflow located at:
http://airflow.readthedocs.org/en/latest/tutorial.html
"""

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

default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2015, 6, 1),
    'email': ['airflow@example.com'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5),
    # 'queue': 'bash_queue',
    # 'pool': 'backfill',
    # 'priority_weight': 10,
    # 'end_date': datetime(2016, 1, 1),
}

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

# t1, t2 and t3 are examples of tasks created by instantiating operators
t1 = BashOperator(
    task_id='print_date',
    bash_command='date',
    dag=dag)

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

templated_command = """
{% for i in range(5) %}
    echo "{{ ds }}"
    echo "{{ macros.ds_add(ds, 7) }}"
    echo "{{ params.my_param }}"
{% endfor %}
"""

t3 = BashOperator(
    task_id='templated',
    bash_command=templated_command,
    params={'my_param': 'Parameter I passed in'},
    dag=dag)

t2.set_upstream(t1)
```

```
t3.set_upstream(t1)
```

3.5.10 Testing

3.5.10.1 Running the Script

Time to run some tests. First let's make sure that the pipeline parses. Let's assume we're saving the code from the previous step in `tutorial.py` in the DAGs folder referenced in your `airflow.cfg`. The default location for your DAGs is `~/airflow/dags`.

```
python ~/airflow/dags/tutorial.py
```

If the script does not raise an exception it means that you haven't done anything horribly wrong, and that your Airflow environment is somewhat sound.

3.5.10.2 Command Line Metadata Validation

Let's run a few commands to validate this script further.

```
# print the list of active DAGs
airflow list_dags

# prints the list of tasks the "tutorial" dag_id
airflow list_tasks tutorial

# prints the hierarchy of tasks in the tutorial DAG
airflow list_tasks tutorial --tree
```

3.5.10.3 Testing

Let's test by running the actual task instances on a specific date. The date specified in this context is an `execution_date`, which simulates the scheduler running your task or dag at a specific date + time:

```
# command layout: command subcommand dag_id task_id date

# testing print_date
airflow test tutorial print_date 2015-06-01

# testing sleep
airflow test tutorial sleep 2015-06-01
```

Now remember what we did with templating earlier? See how this template gets rendered and executed by running this command:

```
# testing templated
airflow test tutorial templated 2015-06-01
```

This should result in displaying a verbose log of events and ultimately running your bash command and printing the result.

Note that the `airflow test` command runs task instances locally, outputs their log to stdout (on screen), doesn't bother with dependencies, and doesn't communicate state (running, success, failed, ...) to the database. It simply allows testing a single task instance.

3.5.10.4 Backfill

Everything looks like it's running fine so let's run a backfill. `backfill` will respect your dependencies, emit logs into files and talk to the database to record status. If you do have a webserver up, you'll be able to track the progress. `airflow webserver` will start a web server if you are interested in tracking the progress visually as your backfill progresses.

Note that if you use `depends_on_past=True`, individual task instances will depend on the success of the preceding task instance, except for the `start_date` specified itself, for which this dependency is disregarded.

The date range in this context is a `start_date` and optionally an `end_date`, which are used to populate the run schedule with task instances from this dag.

```
# optional, start a web server in debug mode in the background
# airflow webserver --debug &

# start your backfill on a date range
airflow backfill tutorial -s 2015-06-01 -e 2015-06-07
```

3.5.11 What's Next?

That's it, you've written, tested and backfilled your very first Airflow pipeline. Merging your code into a code repository that has a master scheduler running against it should get it to get triggered and run every day.

Here's a few things you might want to do next:

- Take an in-depth tour of the UI - click all the things!
- Keep reading the docs! Especially the sections on:
 - Command line interface
 - Operators
 - Macros
- Write your first pipeline!

3.6 Configuration

Setting up the sandbox in the [Quick Start](#) section was easy; building a production-grade environment requires a bit more work!

3.6.1 Setting Configuration Options

The first time you run Airflow, it will create a file called `airflow.cfg` in your `$AIRFLOW_HOME` directory (`~/airflow` by default). This file contains Airflow's configuration and you can edit it to change any of the settings. You can also set options with environment variables by using this format: `$AIRFLOW__{SECTION}__{KEY}` (note the double underscores).

For example, the metadata database connection string can either be set in `airflow.cfg` like this:

```
[core]
sql_alchemy_conn = my_conn_string
```

or by creating a corresponding environment variable:

```
AIRFLOW__CORE__SQL_ALCHEMY_CONN=my_conn_string
```

You can also derive the connection string at run time by appending `_cmd` to the key like this:

```
[core]
sql_alchemy_conn_cmd = bash_command_to_run
```

But only three such configuration elements namely `sql_alchemy_conn`, `broker_url` and `celery_result_backend` can be fetched as a command. The idea behind this is to not store passwords on boxes in plain text files. The order of precedence is as follows -

1. environment variable
2. configuration in `airflow.cfg`
3. command in `airflow.cfg`
4. default

3.6.2 Setting up a Backend

If you want to take a real test drive of Airflow, you should consider setting up a real database backend and switching to the LocalExecutor.

As Airflow was built to interact with its metadata using the great SQLAlchemy library, you should be able to use any database backend supported as a SQLAlchemy backend. We recommend using **MySQL** or **Postgres**.

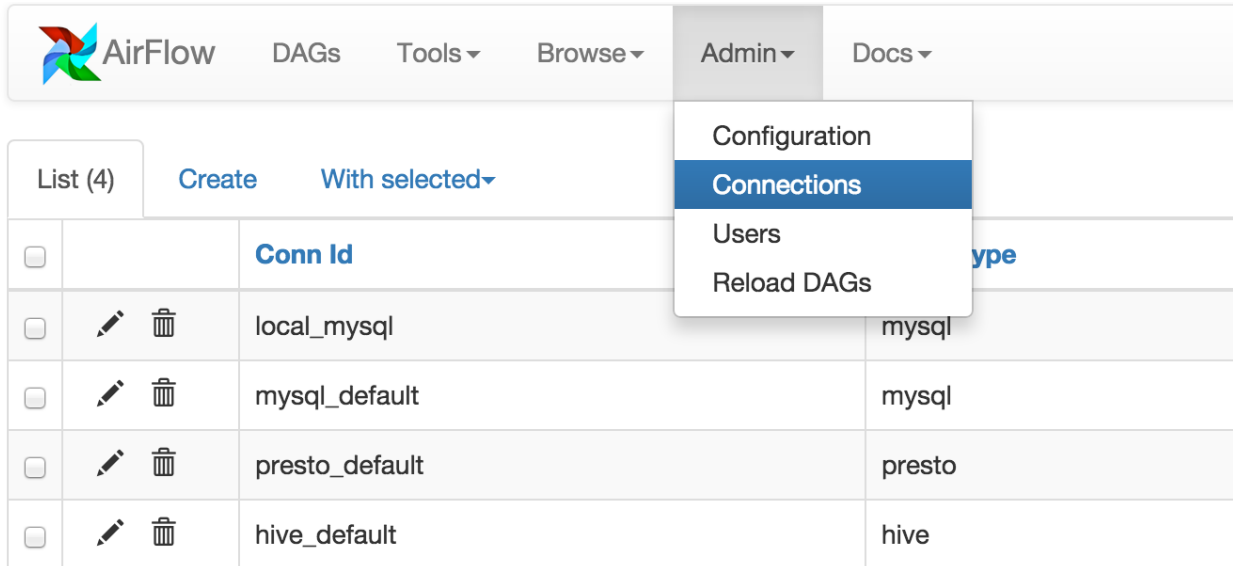
Note: If you decide to use **Postgres**, we recommend using the `psycopg2` driver and specifying it in your SQLAlchemy connection string. Also note that since SQLAlchemy does not expose a way to target a specific schema in the Postgres connection URI, you may want to set a default schema for your role with a command similar to `ALTER ROLE username SET search_path = airflow, foobar;`

Once you’ve setup your database to host Airflow, you’ll need to alter the SQLAlchemy connection string located in your configuration file `$AIRFLOW_HOME/airflow.cfg`. You should then also change the “executor” setting to use “LocalExecutor”, an executor that can parallelize task instances locally.









```
# initialize the database
airflow initdb
```

3.6.3 Connections

Airflow needs to know how to connect to your environment. Information such as hostname, port, login and passwords to other systems and services is handled in the Admin->Connection section of the UI. The pipeline code you will author will reference the ‘`conn_id`’ of the Connection objects.



The screenshot shows the Airflow Admin interface. At the top, there is a navigation bar with the AirFlow logo and tabs for DAGs, Tools, Browse, Admin, and Docs. The Admin tab is selected, and a dropdown menu is open, showing options: Configuration, Connections (highlighted in blue), Users, and Reload DAGs. Below the navigation bar, there is a table of connections. The table has columns for a checkbox, edit/delete icons, Conn Id, and type. The table contains five rows of connections: local_mysql, mysql_default, presto_default, and hive_default.

<input type="checkbox"/>		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

By default, Airflow will save the passwords for the connection in plain text within the metadata database. The `crypto` package is highly recommended during installation. The `crypto` package does require that your operating system have `libffi-dev` installed.

If `crypto` package was not installed initially, you can still enable encryption for connections by following steps below:

1. Install `crypto` package `pip install apache-airflow[crypto]`
2. Generate `fernet_key`, using this code snippet below. `fernet_key` must be a base64-encoded 32-byte key.

```
from cryptography.fernet import Fernet
fernet_key= Fernet.generate_key()
print(fernet_key) # your fernet_key, keep it in secured place!
```

3. Replace `airflow.cfg` `fernet_key` value with the one from step 2. Alternatively, you can store your `fernet_key` in OS environment variable. You do not need to change `airflow.cfg` in this case as AirFlow will use environment variable over the value in `airflow.cfg`:

```
# Note the double underscores
EXPORT AIRFLOW__CORE__FERNET_KEY = your_fernet_key
```

4. Restart AirFlow webserver.
5. For existing connections (the ones that you had defined before installing `airflow[crypto]` and creating a Fernet key), you need to open each connection in the connection admin UI, re-type the password, and save it.

Connections in Airflow pipelines can be created using environment variables. The environment variable needs to have a prefix of `AIRFLOW_CONN_` for Airflow with the value in a URI format to use the connection properly. Please see the [Concepts](#) documentation for more information on environment variables and connections.

3.6.4 Scaling Out with Celery

`CeleryExecutor` is one of the ways you can scale out the number of workers. For this to work, you need to setup a Celery backend (**RabbitMQ**, **Redis**, ...) and change your `airflow.cfg` to point the executor parameter to `CeleryExecutor` and provide the related Celery settings.

For more information about setting up a Celery broker, refer to the exhaustive [Celery documentation on the topic](#).

Here are a few imperative requirements for your workers:

- `airflow` needs to be installed, and the CLI needs to be in the path
- Airflow configuration settings should be homogeneous across the cluster
- Operators that are executed on the worker need to have their dependencies met in that context. For example, if you use the `HiveOperator`, the `hive` CLI needs to be installed on that box, or if you use the `MySqlOperator`, the required Python library needs to be available in the `PYTHONPATH` somehow
- The worker needs to have access to its `DAGS_FOLDER`, and you need to synchronize the filesystems by your own means. A common setup would be to store your `DAGS_FOLDER` in a Git repository and sync it across machines using Chef, Puppet, Ansible, or whatever you use to configure machines in your environment. If all your boxes have a common mount point, having your pipelines files shared there should work as well

To kick off a worker, you need to setup Airflow and kick off the worker subcommand

```
airflow worker
```

Your worker should start picking up tasks as soon as they get fired in its direction.

Note that you can also run “Celery Flower”, a web UI built on top of Celery, to monitor your workers. You can use the shortcut command `airflow flower` to start a Flower web server.

3.6.5 Scaling Out with Dask

`DaskExecutor` allows you to run Airflow tasks in a Dask Distributed cluster.

Dask clusters can be run on a single machine or on remote networks. For complete details, consult the [Distributed documentation](#).

To create a cluster, first start a Scheduler:

```
# default settings for a local cluster
DASK_HOST=127.0.0.1
DASK_PORT=8786

dask-scheduler --host $DASK_HOST --port $DASK_PORT
```

Next start at least one Worker on any machine that can connect to the host:

```
dask-worker $DASK_HOST:$DASK_PORT
```

Edit your `airflow.cfg` to set your executor to `DaskExecutor` and provide the Dask Scheduler address in the `[dask]` section.

Please note:

- Each Dask worker must be able to import Airflow and any dependencies you require.
- Dask does not support queues. If an Airflow task was created with a queue, a warning will be raised but the task will be submitted to the cluster.

3.6.6 Logs

Users can specify a logs folder in `airflow.cfg`. By default, it is in the `AIRFLOW_HOME` directory.

In addition, users can supply a remote location for storing logs and log backups in cloud storage. At this time, Amazon S3 and Google Cloud Storage are supported. To enable this feature, `airflow.cfg` must be configured as in this example:

```
[core]
# Airflow can store logs remotely in AWS S3 or Google Cloud Storage. Users
# must supply a remote location URL (starting with either 's3://...' or
# 'gs://...') and an Airflow connection id that provides access to the storage
# location.
remote_base_log_folder = s3://my-bucket/path/to/logs
remote_log_conn_id = MyS3Conn
# Use server-side encryption for logs stored in S3
encrypt_s3_logs = False
```

Remote logging uses an existing Airflow connection to read/write logs. If you don't have a connection properly setup, this will fail. In the above example, Airflow will try to use `S3Hook('MyS3Conn')`.

In the Airflow Web UI, local logs take precedence over remote logs. If local logs can not be found or accessed, the remote logs will be displayed. Note that logs are only sent to remote storage once a task completes (including failure). In other words, remote logs for running tasks are unavailable.

3.6.7 Scaling Out on Mesos (community contributed)

`MesosExecutor` allows you to schedule airflow tasks on a Mesos cluster. For this to work, you need a running mesos cluster and you must perform the following steps -

1. Install airflow on a machine where web server and scheduler will run, let's refer to this as the "Airflow server".
2. On the Airflow server, install mesos python eggs from [mesos downloads](#).
3. On the Airflow server, use a database (such as mysql) which can be accessed from mesos slave machines and add configuration in `airflow.cfg`.
4. Change your `airflow.cfg` to point executor parameter to `MesosExecutor` and provide related Mesos settings.
5. On all mesos slaves, install airflow. Copy the `airflow.cfg` from Airflow server (so that it uses same sql alchemy connection).
6. On all mesos slaves, run the following for serving logs:

```
airflow serve_logs
```

7. On Airflow server, to start processing/scheduling DAGs on mesos, run:

```
airflow scheduler -p
```

Note: We need `-p` parameter to pickle the DAGs.

You can now see the airflow framework and corresponding tasks in mesos UI. The logs for airflow tasks can be seen in airflow UI as usual.

For more information about mesos, refer to [mesos documentation](#). For any queries/bugs on `MesosExecutor`, please contact [@kapil-malik](#).

3.6.8 Integration with systemd

Airflow can integrate with systemd based systems. This makes watching your daemons easy as systemd can take care of restarting a daemon on failure. In the `scripts/systemd` directory you can find unit files that have been tested

on Redhat based systems. You can copy those to `/usr/lib/systemd/system`. It is assumed that Airflow will run under `airflow:airflow`. If not (or if you are running on a non Redhat based system) you probably need to adjust the unit files.

Environment configuration is picked up from `/etc/sysconfig/airflow`. An example file is supplied. Make sure to specify the `SCHEDULER_RUNS` variable in this file when you run the scheduler. You can also define here, for example, `AIRFLOW_HOME` or `AIRFLOW_CONFIG`.

3.6.9 Integration with upstart

Airflow can integrate with upstart based systems. Upstart automatically starts all airflow services for which you have a corresponding `*.conf` file in `/etc/init` upon system boot. On failure, upstart automatically restarts the process (until it reaches re-spawn limit set in a `*.conf` file).

You can find sample upstart job files in the `scripts/upstart` directory. These files have been tested on Ubuntu 14.04 LTS. You may have to adjust `start on` and `stop on` stanzas to make it work on other upstart systems. Some of the possible options are listed in `scripts/upstart/README`.

Modify `*.conf` files as needed and copy to `/etc/init` directory. It is assumed that airflow will run under `airflow:airflow`. Change `setuid` and `setgid` in `*.conf` files if you use other user/group

You can use `initctl` to manually start, stop, view status of the airflow process that has been integrated with upstart

```
initctl airflow-webserver status
```

3.6.10 Test Mode

Airflow has a fixed set of “test mode” configuration options. You can load these at any time by calling `airflow.configuration.load_test_config()` (note this operation is not reversible!). However, some options (like the `DAG_FOLDER`) are loaded before you have a chance to call `load_test_config()`. In order to eagerly load the test configuration, set `test_mode` in `airflow.cfg`:

```
[tests]
unit_test_mode = True
```

Due to Airflow’s automatic environment variable expansion (see [Setting Configuration Options](#)), you can also set the env var `AIRFLOW__CORE__UNIT_TEST_MODE` to temporarily overwrite `airflow.cfg`.

3.7 UI / Screenshots

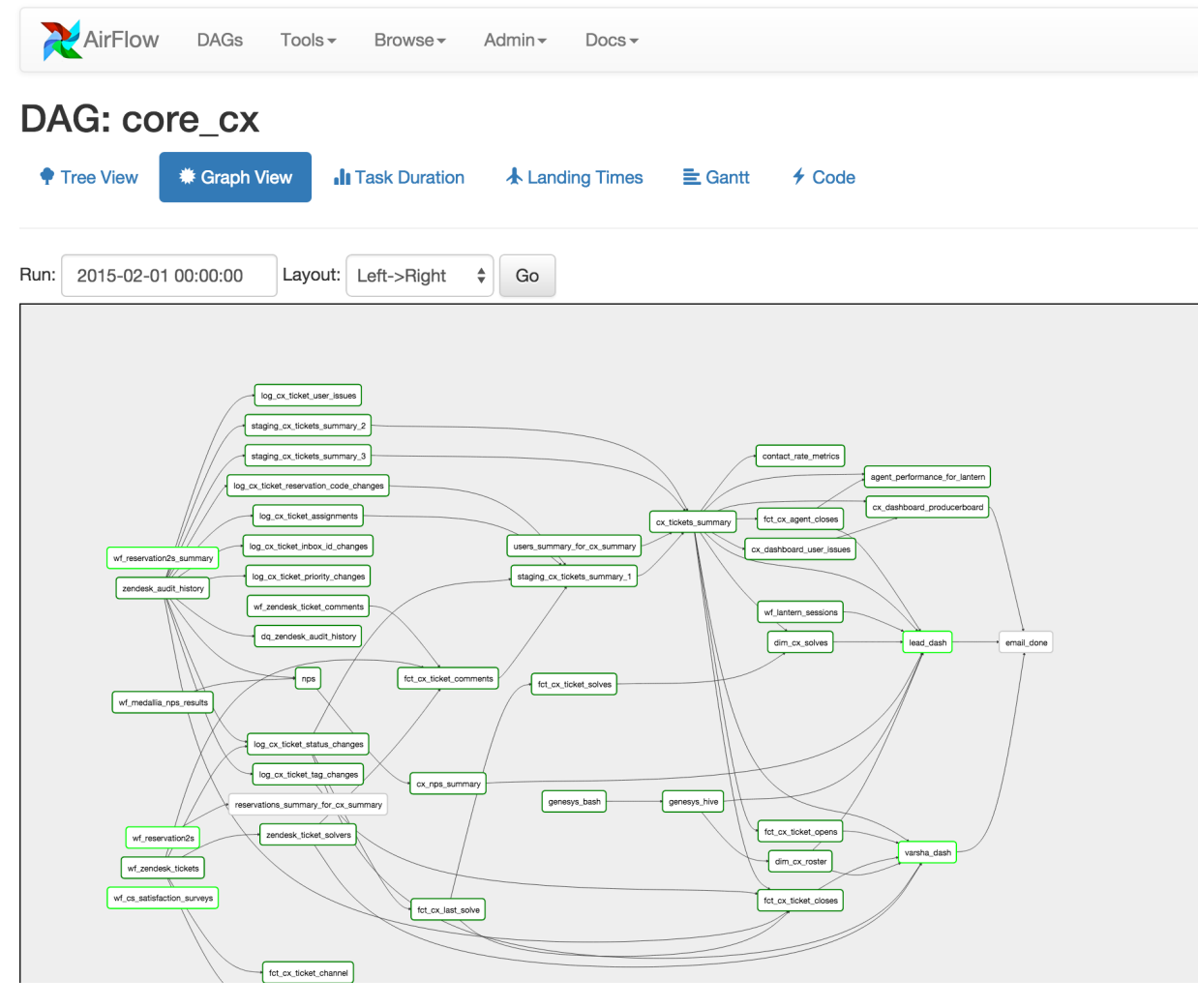
The Airflow UI make it easy to monitor and troubleshoot your data pipelines. Here’s a quick overview of some of the features and visualizations you can find in the Airflow UI.

3.7.1 DAGs View

List of the DAGs in your environment, and a set of shortcuts to useful pages. You can see exactly how many tasks succeeded, failed, or are currently running at a glance.

3.7.3 Graph View

The graph view is perhaps the most comprehensive. Visualize your DAG's dependencies and their current status for a specific run.



3.7.4 Variable View

The variable view allows you to list, create, edit or delete the key-value pair of a variable used during jobs. Value of a variable will be hidden if the key contains any words in ('password', 'secret', 'passwd', 'authorization', 'api_key', 'apikey', 'access_token') by default, but can be configured to show in clear-text.



Variables



















List (9)

Create

Add Filter

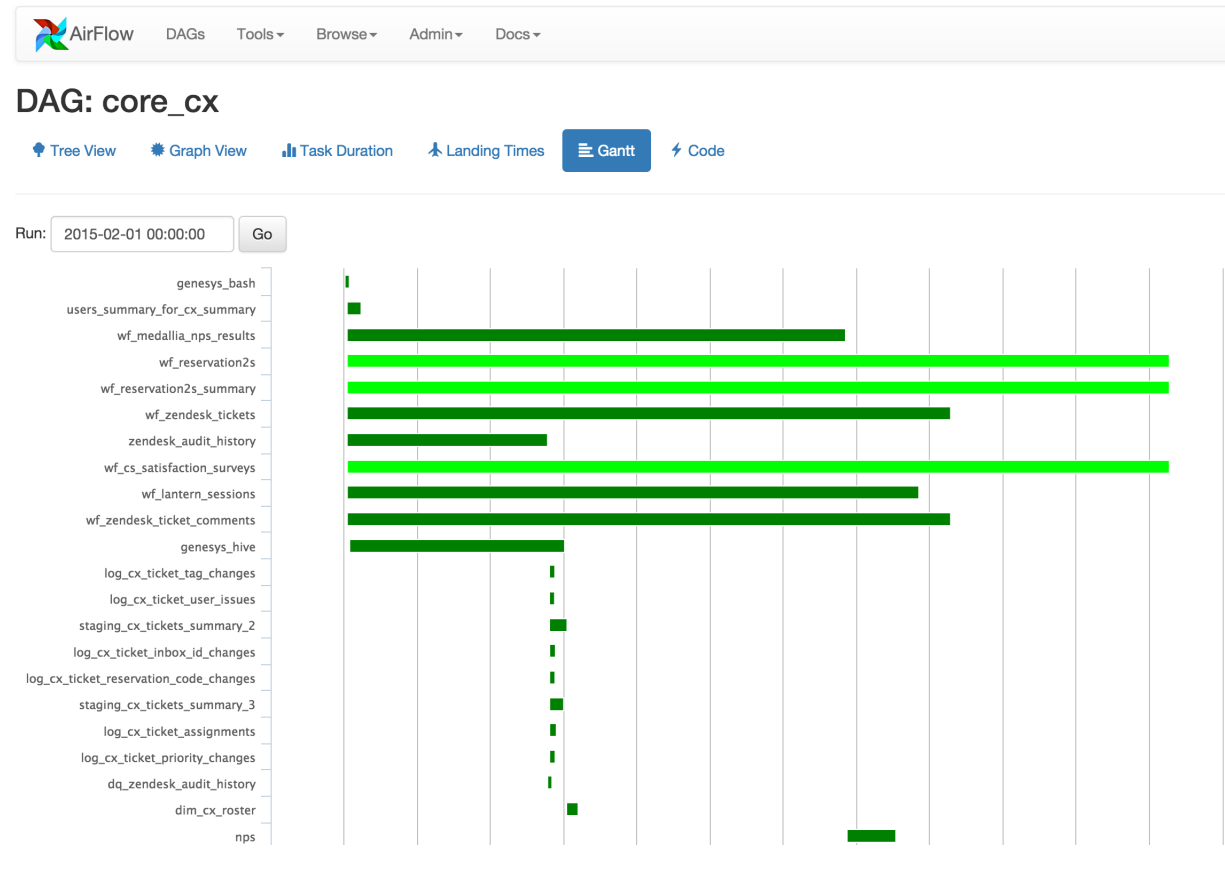
With selected

Search

<input type="checkbox"/>		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	*****

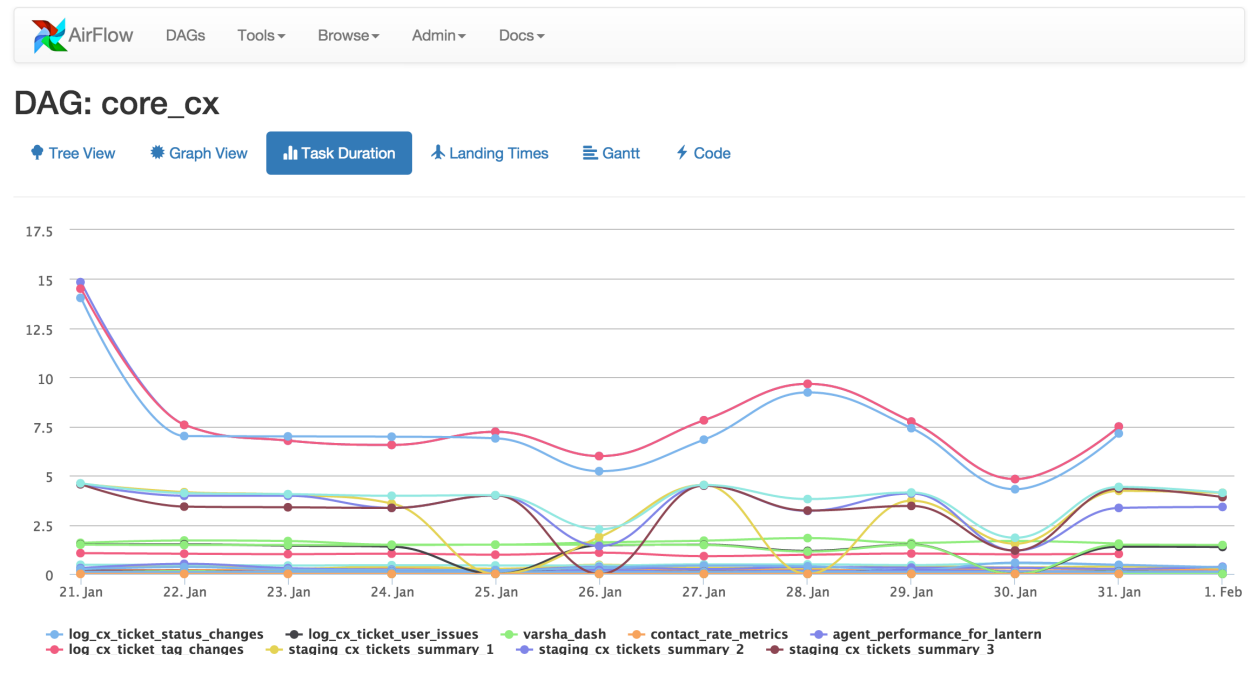
3.7.5 Gantt Chart

The Gantt chart lets you analyse task duration and overlap. You can quickly identify bottlenecks and where the bulk of the time is spent for specific DAG runs.



3.7.6 Task Duration

The duration of your different tasks over the past N runs. This view lets you find outliers and quickly understand where the time is spent in your DAG over many runs.



3.7.7 Code View

Transparency is everything. While the code for your pipeline is in source control, this is a quick way to get to the code that generates the DAG and provide yet more context.



DAG: example1

[Tree View](#)[Graph View](#)[Task Duration](#)[Landing Times](#)[Gantt](#)[Code](#)

example_dags/example1.py

```
from airflow.operators import BashOperator, DummyOperator
from airflow.models import DAG
from datetime import datetime

args = {
    'owner': 'airflow',
    'start_date': datetime(2015, 1, 1),
}

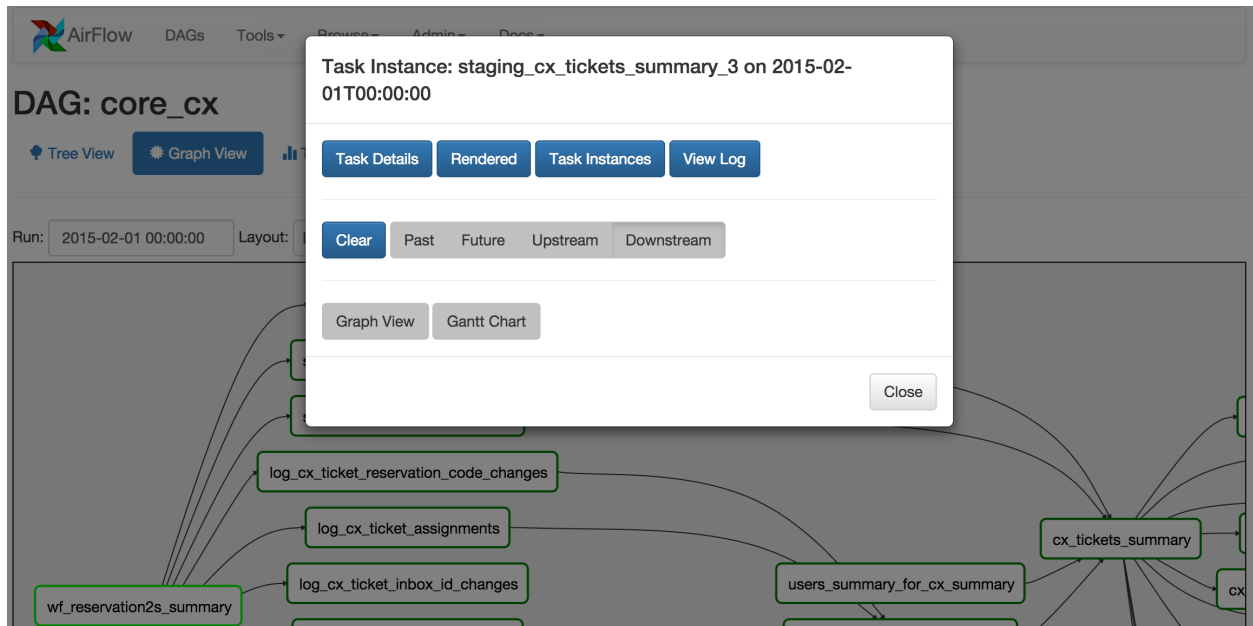
dag = DAG(dag_id='example1')

cmd = 'ls -l'
run_this_last = DummyOperator(
    task_id='run_this_last',
    default_args=args)
dag.add_task(run_this_last)

run_this = BashOperator(
    task_id='run_after_loop', bash_command='echo 1',
    default_args=args)
dag.add_task(run_this)
run_this.set_downstream(run_this_last)
for i in range(9):
    i = str(i)
    task = BashOperator(
```

3.7.8 Task Instance Context Menu

From the pages seen above (tree view, graph view, gantt, ...), it is always possible to click on a task instance, and get to this rich context menu that can take you to more detailed metadata, and perform some actions.



3.8 Concepts

The Airflow Platform is a tool for describing, executing, and monitoring workflows.

3.8.1 Core Ideas

3.8.1.1 DAGs

In Airflow, a DAG – or a Directed Acyclic Graph – is a collection of all the tasks you want to run, organized in a way that reflects their relationships and dependencies.

For example, a simple DAG could consist of three tasks: A, B, and C. It could say that A has to run successfully before B can run, but C can run anytime. It could say that task A times out after 5 minutes, and B can be restarted up to 5 times in case it fails. It might also say that the workflow will run every night at 10pm, but shouldn't start until a certain date.

In this way, a DAG describes *how* you want to carry out your workflow; but notice that we haven't said anything about *what* we actually want to do! A, B, and C could be anything. Maybe A prepares data for B to analyze while C sends an email. Or perhaps A monitors your location so B can open your garage door while C turns on your house lights. The important thing is that the DAG isn't concerned with what its constituent tasks do; its job is to make sure that whatever they do happens at the right time, or in the right order, or with the right handling of any unexpected issues.

DAGs are defined in standard Python files that are placed in Airflow's `DAG_FOLDER`. Airflow will execute the code in each file to dynamically build the DAG objects. You can have as many DAGs as you want, each describing an arbitrary number of tasks. In general, each one should correspond to a single logical workflow.

Scope

Airflow will load any DAG object it can import from a DAGfile. Critically, that means the DAG must appear in `globals()`. Consider the following two DAGs. Only `dag_1` will be loaded; the other one only appears in a local scope.

```
dag_1 = DAG('this_dag_will_be_discovered')

def my_function():
    dag_2 = DAG('but_this_dag_will_not')

my_function()
```

Sometimes this can be put to good use. For example, a common pattern with `SubDagOperator` is to define the subdag inside a function so that Airflow doesn't try to load it as a standalone DAG.

Default Arguments

If a dictionary of `default_args` is passed to a DAG, it will apply them to any of its operators. This makes it easy to apply a common parameter to many operators without having to type it many times.

```
default_args=dict(
    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
```

Context Manager

Added in Airflow 1.8

DAGs can be used as context managers to automatically assign new operators to that DAG.

```
with DAG('my_dag', start_date=datetime(2016, 1, 1)) as dag:
    op = DummyOperator('op')

op.dag is dag # True
```

3.8.1.2 Operators

While DAGs describe *how* to run a workflow, `Operators` determine what actually gets done.

An operator describes a single task in a workflow. Operators are usually (but not always) atomic, meaning they can stand on their own and don't need to share resources with any other operators. The DAG will make sure that operators run in the correct certain order; other than those dependencies, operators generally run independently. In fact, they may run on two completely different machines.

This is a subtle but very important point: in general, if two operators need to share information, like a filename or small amount of data, you should consider combining them into a single operator. If it absolutely can't be avoided, Airflow does have a feature for operator cross-communication called XCom that is described elsewhere in this document.

Airflow provides operators for many common tasks, including:

- `BashOperator` - executes a bash command
- `PythonOperator` - calls an arbitrary Python function
- `EmailOperator` - sends an email
- `HTTPOperator` - sends an HTTP request

- `MySqlOperator`, `SqliteOperator`, `PostgresOperator`, `MsSqlOperator`, `OracleOperator`, `JdbcOperator`, etc. - executes a SQL command
- `Sensor` - waits for a certain time, file, database row, S3 key, etc...

In addition to these basic building blocks, there are many more specific operators: `DockerOperator`, `HiveOperator`, `S3FileTransferOperator`, `PrestoToMysqlOperator`, `SlackOperator`... you get the idea!

The `airflow/contrib/` directory contains yet more operators built by the community. These operators aren't always as complete or well-tested as those in the main distribution, but allow users to more easily add new functionality to the platform.

Operators are only loaded by Airflow if they are assigned to a DAG.

DAG Assignment

Added in Airflow 1.8

Operators do not have to be assigned to DAGs immediately (previously `dag` was a required argument). However, once an operator is assigned to a DAG, it can not be transferred or unassigned. DAG assignment can be done explicitly when the operator is created, through deferred assignment, or even inferred from other operators.

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

# sets the DAG explicitly
explicit_op = DummyOperator(task_id='op1', dag=dag)

# deferred DAG assignment
deferred_op = DummyOperator(task_id='op2')
deferred_op.dag = dag

# inferred DAG assignment (linked operators must be in the same DAG)
inferred_op = DummyOperator(task_id='op3')
inferred_op.set_upstream(deferred_op)
```

Bitshift Composition

Added in Airflow 1.8

Traditionally, operator relationships are set with the `set_upstream()` and `set_downstream()` methods. In Airflow 1.8, this can be done with the Python bitshift operators `>>` and `<<`. The following four statements are all functionally equivalent:

```
op1 >> op2
op1.set_downstream(op2)

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

When using the bitshift to compose operators, the relationship is set in the direction that the bitshift operator points. For example, `op1 >> op2` means that `op1` runs first and `op2` runs second. Multiple operators can be composed – keep in mind the chain is executed left-to-right and the rightmost object is always returned. For example:

```
op1 >> op2 >> op3 << op4
```

is equivalent to:

```
op1.set_downstream(op2)
op2.set_downstream(op3)
op3.set_upstream(op4)
```

For convenience, the bitshift operators can also be used with DAGs. For example:

```
dag >> op1 >> op2
```

is equivalent to:

```
op1.dag = dag
op1.set_downstream(op2)
```

We can put this all together to build a simple pipeline:

```
with DAG('my_dag', start_date=datetime(2016, 1, 1)) as dag:
    (
        DummyOperator(task_id='dummy_1')
        >> BashOperator(
            task_id='bash_1',
            bash_command='echo "HELLO!"')
        >> PythonOperator(
            task_id='python_1',
            python_callable=lambda: print("GOODBYE!"))
    )
```

3.8.1.3 Tasks

Once an operator is instantiated, it is referred to as a “task”. The instantiation defines specific values when calling the abstract operator, and the parameterized task becomes a node in a DAG.

3.8.1.4 Task Instances

A task instance represents a specific run of a task and is characterized as the combination of a dag, a task, and a point in time. Task instances also have an indicative state, which could be “running”, “success”, “failed”, “skipped”, “up for retry”, etc.

3.8.1.5 Workflows

You’re now familiar with the core building blocks of Airflow. Some of the concepts may sound very similar, but the vocabulary can be conceptualized like this:

- 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 1) has been assigned to a DAG and 2) has a state associated with a specific run of the DAG

By combining DAGs and Operators to create TaskInstances, you can build complex workflows.

3.8.2 Additional Functionality

In addition to the core Airflow objects, there are a number of more complex features that enable behaviors like limiting simultaneous access to resources, cross-communication, conditional execution, and more.

3.8.2.1 Hooks

Hooks are interfaces to external platforms and databases like Hive, S3, MySQL, Postgres, HDFS, and Pig. Hooks implement a common interface when possible, and act as a building block for operators. They also use the `airflow.models.Connection` model to retrieve hostnames and authentication information. Hooks keep authentication code and information out of pipelines, centralized in the metadata database.

Hooks are also very useful on their own to use in Python scripts, Airflow `airflow.operators.PythonOperator`, and in interactive environments like iPython or Jupyter Notebook.

3.8.2.2 Pools

Some systems can get overwhelmed when too many processes hit them at the same time. Airflow pools can be used to **limit the execution parallelism** on arbitrary sets of tasks. The list of pools is managed in the UI (Menu -> Admin -> Pools) by giving the pools a name and assigning it a number of worker slots. Tasks can then be associated with one of the existing pools by using the `pool` parameter when creating tasks (i.e., instantiating operators).

```
aggregate_db_message_job = BashOperator(
    task_id='aggregate_db_message_job',
    execution_timeout=timedelta(hours=3),
    pool='ep_data_pipeline_db_msg_agg',
    bash_command=aggregate_db_message_job_cmd,
    dag=dag)
aggregate_db_message_job.set_upstream(wait_for_empty_queue)
```

The `pool` parameter can be used in conjunction with `priority_weight` to define priorities in the queue, and which tasks get executed first as slots open up in the pool. The default `priority_weight` is 1, and can be bumped to any number. When sorting the queue to evaluate which task should be executed next, we use the `priority_weight`, summed up with all of the `priority_weight` values from tasks downstream from this task. You can use this to bump a specific important task and the whole path to that task gets prioritized accordingly.

Tasks will be scheduled as usual while the slots fill up. Once capacity is reached, runnable tasks get queued and their state will show as such in the UI. As slots free up, queued tasks start running based on the `priority_weight` (of the task and its descendants).

Note that by default tasks aren't assigned to any pool and their execution parallelism is only limited to the executor's setting.

3.8.2.3 Connections

The connection information to external systems is stored in the Airflow metadata database and managed in the UI (Menu -> Admin -> Connections). A `conn_id` is defined there and hostname / login / password / schema information attached to it. Airflow pipelines can simply refer to the centrally managed `conn_id` without having to hard code any of this information anywhere.

Many connections with the same `conn_id` can be defined and when that is the case, and when the **hooks** uses the `get_connection` method from `BaseHook`, Airflow will choose one connection randomly, allowing for some basic load balancing and fault tolerance when used in conjunction with retries.

Airflow also has the ability to reference connections via environment variables from the operating system. The environment variable needs to be prefixed with `AIRFLOW_CONN_` to be considered a connection. When referencing the connection in the Airflow pipeline, the `conn_id` should be the name of the variable without the prefix. For example, if the `conn_id` is named `postgres_master` the environment variable should be named `AIRFLOW_CONN_POSTGRES_MASTER` (note that the environment variable must be all uppercase). Airflow assumes the value returned from the environment variable to be in a URI format (e.g. `postgres://user:password@localhost:5432/master` or `s3://accesskey:secretkey@S3`).

3.8.2.4 Queues

When using the CeleryExecutor, the celery queues that tasks are sent to can be specified. `queue` is an attribute of `BaseOperator`, so any task can be assigned to any queue. The default queue for the environment is defined in the `airflow.cfg`'s `celery -> default_queue`. This defines the queue that tasks get assigned to when not specified, as well as which queue Airflow workers listen to when started.

Workers can listen to one or multiple queues of tasks. When a worker is started (using the command `airflow worker`), a set of comma delimited queue names can be specified (e.g. `airflow worker -q spark`). This worker will then only pick up tasks wired to the specified queue(s).

This can be useful if you need specialized workers, either from a resource perspective (for say very lightweight tasks where one worker could take thousands of tasks without a problem), or from an environment perspective (you want a worker running from within the Spark cluster itself because it needs a very specific environment and security rights).

3.8.2.5 XComs

XComs let tasks exchange messages, allowing more nuanced forms of control and shared state. The name is an abbreviation of “cross-communication”. XComs are principally defined by a key, value, and timestamp, but 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, so users should make sure to use objects of appropriate size.

XComs can be “pushed” (sent) or “pulled” (received). When a task pushes an XCom, it makes it generally available to other tasks. Tasks can push XComs at any time by calling the `xcom_push()` method. In addition, if a task returns a value (either from its `Operator`'s `execute()` method, or from a `PythonOperator`'s `python_callable` function), then an XCom containing that value is automatically pushed.

Tasks call `xcom_pull()` to retrieve XComs, optionally applying filters based on criteria like `key`, `source task_ids`, and `source dag_id`. By default, `xcom_pull()` filters for the keys that are automatically given to XComs when they are pushed by being returned from `execute` functions (as opposed to XComs that are pushed manually).

If `xcom_pull` is passed a single string for `task_ids`, then the most recent XCom value from that task is returned; if a list of `task_ids` is passed, then a corresponding list of XCom values is returned.

```
# 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')
```

It is also possible to pull XCom directly in a template, here's an example of what this may look like:

```
SELECT * FROM {{ task_instance.xcom_pull(task_ids='foo', key='table_name') }}
```

Note that XComs are similar to *Variables*, but are specifically designed for inter-task communication rather than global settings.

3.8.2.6 Variables

Variables are 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.

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

The second call assumes json content and will be deserialized into bar. Note that Variable is a sqlalchemy model and can be used as such.

3.8.2.7 Branching

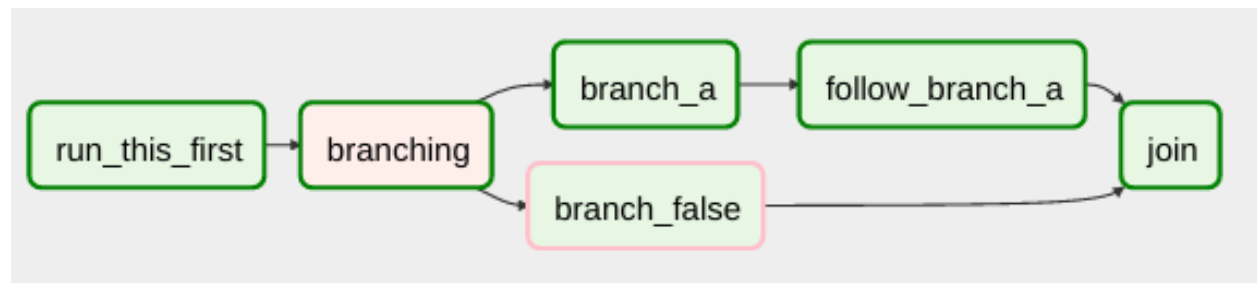
Sometimes you need a workflow to branch, or only go down a certain path based on an arbitrary condition which is typically related to something that happened in an upstream task. One way to do this is by using the BranchPythonOperator.

The BranchPythonOperator is much like the PythonOperator except that it expects a python_callable that returns a task_id. The task_id returned is followed, and all of the other paths are skipped. The task_id returned by the Python function has to be referencing a task directly downstream from the BranchPythonOperator task.

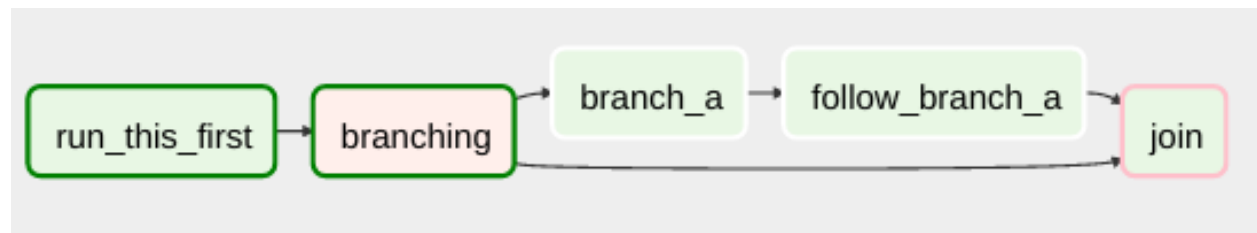
Note that using tasks with depends_on_past=True downstream from BranchPythonOperator is logically unsound as skipped status will invariably lead to block tasks that depend on their past successes. skipped states propagates where all directly upstream tasks are skipped.

If you want to skip some tasks, keep in mind that you can't have an empty path, if so make a dummy task.

like this, the dummy task "branch_false" is skipped



Not like this, where the join task is skipped

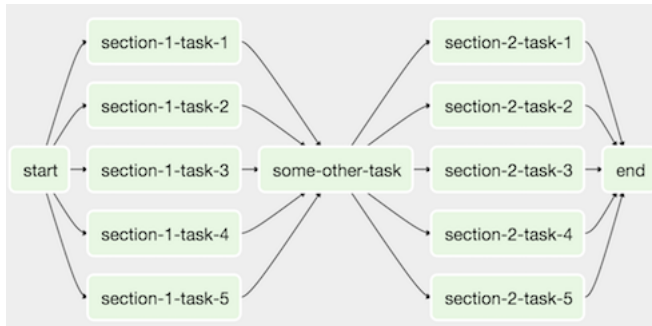


3.8.2.8 SubDAGs

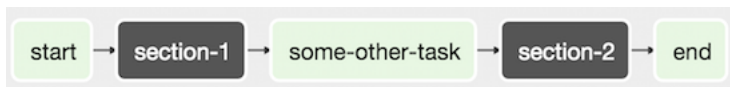
SubDAGs are perfect for repeating patterns. Defining a function that returns a DAG object is a nice design pattern when using Airflow.

Airbnb uses the *stage-check-exchange* pattern when loading data. Data is staged in a temporary table, after which data quality checks are performed against that table. Once the checks all pass the partition is moved into the production table.

As another example, consider the following DAG:



We can combine all of the parallel `task-*` operators into a single SubDAG, so that the resulting DAG resembles the following:



Note that SubDAG operators should contain a factory method that returns a DAG object. This will prevent the SubDAG from being treated like a separate DAG in the main UI. For example:

```
# dags/subdag.py
from airflow.models import DAG
from airflow.operators.dummy_operator import DummyOperator

# Dag is returned by a factory method
def sub_dag(parent_dag_name, child_dag_name, start_date, schedule_interval):
    dag = DAG(
        '%s.%s' % (parent_dag_name, child_dag_name),
        schedule_interval=schedule_interval,
        start_date=start_date,
    )

    dummy_operator = DummyOperator(
        task_id='dummy_task',
        dag=dag,
    )

    return dag
```

This SubDAG can then be referenced in your main DAG file:

```
# main_dag.py
from datetime import datetime, timedelta
from airflow.models import DAG
from airflow.operators.subdag_operator import SubDagOperator
from dags.subdag import sub_dag
```

```

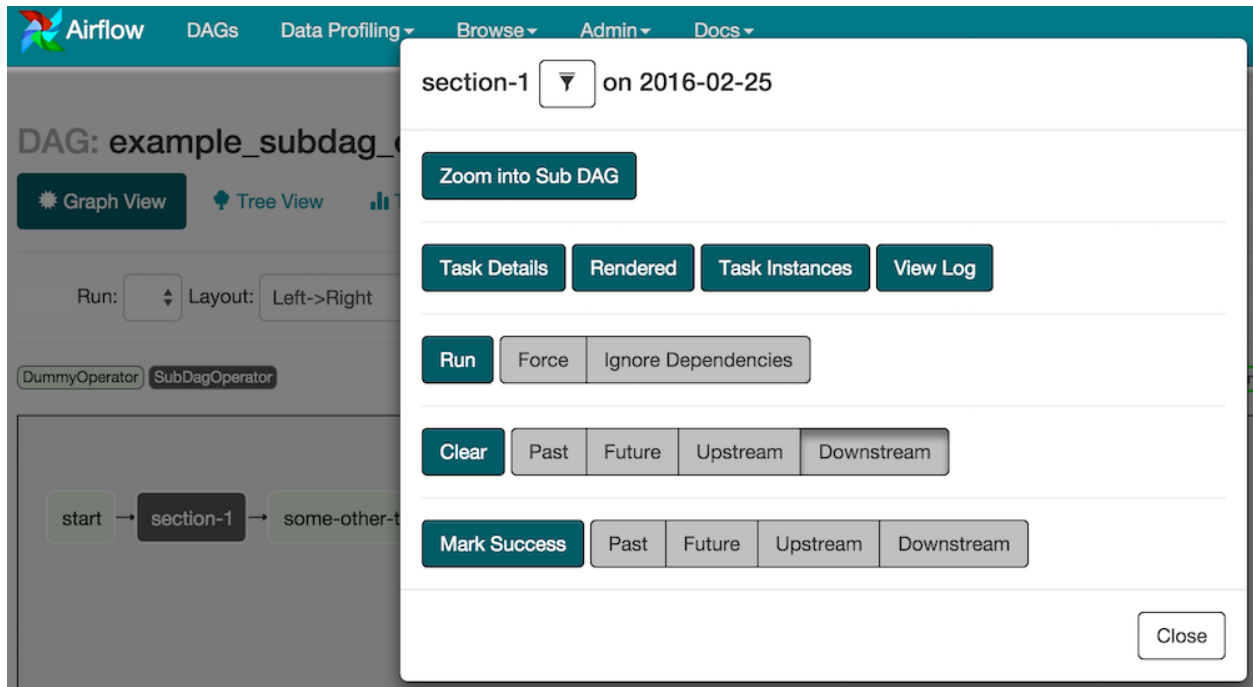
PARENT_DAG_NAME = 'parent_dag'
CHILD_DAG_NAME = 'child_dag'

main_dag = DAG(
    dag_id=PARENT_DAG_NAME,
    schedule_interval=timedelta(hours=1),
    start_date=datetime(2016, 1, 1)
)

sub_dag = SubDagOperator(
    subdag=sub_dag(PARENT_DAG_NAME, CHILD_DAG_NAME, main_dag.start_date,
                    main_dag.schedule_interval),
    task_id=CHILD_DAG_NAME,
    dag=main_dag,
)

```

You can zoom into a SubDagOperator from the graph view of the main DAG to show the tasks contained within the SubDAG:



Some other tips when using SubDAGs:

- by convention, a SubDAG's `dag_id` should be prefixed by its parent and a dot. As in `parent.child`
- share arguments between the main DAG and the SubDAG by passing arguments to the SubDAG operator (as demonstrated above)
- SubDAGs must have a schedule and be enabled. If the SubDAG's schedule is set to `None` or `@once`, the SubDAG will succeed without having done anything
- clearing a SubDagOperator also clears the state of the tasks within
- marking success on a SubDagOperator does not affect the state of the tasks within
- refrain from using `depends_on_past=True` in tasks within the SubDAG as this can be confusing

- it is possible to specify an executor for the SubDAG. It is common to use the `SequentialExecutor` if you want to run the SubDAG in-process and effectively limit its parallelism to one. Using `LocalExecutor` can be problematic as it may over-subscribe your worker, running multiple tasks in a single slot

See `airflow/example_dags` for a demonstration.

3.8.2.9 SLAs

Service Level Agreements, or time by which a task or DAG should have succeeded, can be set at a task level as a `timedelta`. If one or many instances have not succeeded by that time, an alert email is sent detailing the list of tasks that missed their SLA. The event is also recorded in the database and made available in the web UI under `Browse->Missed SLAs` where events can be analyzed and documented.

3.8.2.10 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

Note that these can be used in conjunction with `depends_on_past` (boolean) that, when set to `True`, keeps a task from getting triggered if the previous schedule for the task hasn't succeeded.

3.8.2.11 Latest Run Only

Standard workflow behavior involves running a series of tasks for a particular date/time range. Some workflows, however, perform tasks that are independent of run time but need to be run on a schedule, much like a standard cron job. In these cases, backfills or running jobs missed during a pause just wastes CPU cycles.

For situations like this, you can use the `LatestOnlyOperator` to skip tasks that are not being run during the most recent scheduled run for a DAG. The `LatestOnlyOperator` skips all immediate downstream tasks, and itself, if the time right now is not between its `execution_time` and the next scheduled `execution_time`.

One must be aware of the interaction between skipped tasks and trigger rules. Skipped tasks will cascade through trigger rules `all_success` and `all_failed` but not `all_done`, `one_failed`, `one_success`, and `dummy`. If you would like to use the `LatestOnlyOperator` with trigger rules that do not cascade skips, you will need to ensure that the `LatestOnlyOperator` is **directly** upstream of the task you would like to skip.

It is possible, through use of trigger rules to mix tasks that should run in the typical date/time dependent mode and those using the `LatestOnlyOperator`.

For example, consider the following dag:


```
#dags/latest_only_with_trigger.py
import datetime as dt

from airflow.models import DAG
from airflow.operators.dummy_operator import DummyOperator
from airflow.operators.latest_only_operator import LatestOnlyOperator
from airflow.utils.trigger_rule import TriggerRule

dag = DAG(
    dag_id='latest_only_with_trigger',
    schedule_interval=dt.timedelta(hours=4),
    start_date=dt.datetime(2016, 9, 20),
)

latest_only = LatestOnlyOperator(task_id='latest_only', dag=dag)

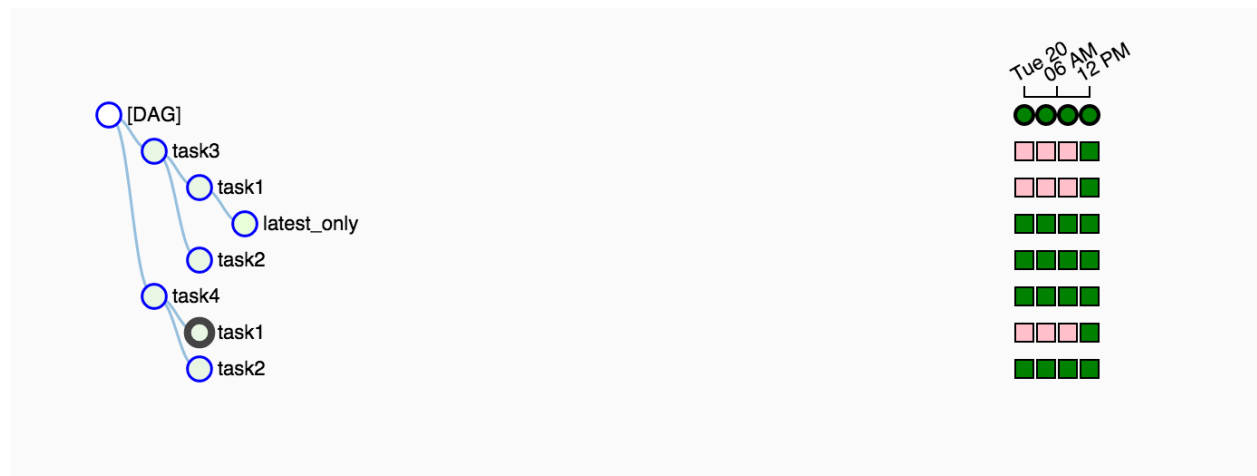
task1 = DummyOperator(task_id='task1', dag=dag)
task1.set_upstream(latest_only)

task2 = DummyOperator(task_id='task2', dag=dag)

task3 = DummyOperator(task_id='task3', dag=dag)
task3.set_upstream([task1, task2])

task4 = DummyOperator(task_id='task4', dag=dag,
                       trigger_rule=TriggerRule.ALL_DONE)
task4.set_upstream([task1, task2])
```

In the case of this dag, the `latest_only` task will show up as skipped for all runs except the latest run. `task1` is directly downstream of `latest_only` and will also skip for all runs except the latest. `task2` is entirely independent of `latest_only` and will run in all scheduled periods. `task3` is downstream of `task1` and `task2` and because of the default `trigger_rule` being `all_success` will receive a cascaded skip from `task1`. `task4` is downstream of `task1` and `task2` but since its `trigger_rule` is set to `all_done` it will trigger as soon as `task1` has been skipped (a valid completion state) and `task2` has succeeded.



3.8.2.12 Zombies & Undeads

Task instances die all the time, usually as part of their normal life cycle, but sometimes unexpectedly.

Zombie tasks are characterized by the absence of an heartbeat (emitted by the job periodically) and a `running` status in the database. They can occur when a worker node can't reach the database, when Airflow processes are killed externally, or when a node gets rebooted for instance. Zombie killing is performed periodically by the scheduler's process.

Undead processes are characterized by the existence of a process and a matching heartbeat, but Airflow isn't aware of this task as `running` in the database. This mismatch typically occurs as the state of the database is altered, most likely by deleting rows in the "Task Instances" view in the UI. Tasks are instructed to verify their state as part of the heartbeat routine, and terminate themselves upon figuring out that they are in this "undead" state.

3.8.2.13 Cluster Policy

Your local airflow settings file can define a `policy` function that has the ability to mutate task attributes based on other task or DAG attributes. It receives a single argument as a reference to task objects, and is expected to alter its attributes.

For example, this function could apply a specific queue property when using a specific operator, or enforce a task timeout policy, making sure that no tasks run for more than 48 hours. Here's an example of what this may look like inside your `airflow_settings.py`:

```
def policy(task):
    if task.__class__.__name__ == 'HivePartitionSensor':
        task.queue = "sensor_queue"
    if task.timeout > timedelta(hours=48):
        task.timeout = timedelta(hours=48)
```

3.8.2.14 Documentation & Notes

It's possible to add documentation or notes to your dags & task objects that become visible in the web interface ("Graph View" for dags, "Task Details" for tasks). There are a set of special task attributes that get rendered as rich content if defined:

attribute	rendered to
<code>doc</code>	monospace
<code>doc_json</code>	json
<code>doc_yaml</code>	yaml
<code>doc_md</code>	markdown
<code>doc_rst</code>	reStructuredText

Please note that for dags, `dag_md` is the only attribute interpreted.

This is especially useful if your tasks are built dynamically from configuration files, it allows you to expose the configuration that led to the related tasks in Airflow.

```
"""
### My great DAG
"""

dag = DAG('my_dag', default_args=default_args)
dag.doc_md = __doc__

t = BashOperator("foo", dag=dag)
t.doc_md = """\
#Title"
```

```
Here's a [url] (www.airbnb.com)
"""
```

This content will get rendered as markdown respectively in the “Graph View” and “Task Details” pages.

3.8.2.15 Jinja Templating

Airflow leverages the power of [Jinja Templating](#) and this can be a powerful tool to use in combination with macros (see the [Macros](#) section).

For example, say you want to pass the execution date as an environment variable to a Bash script using the `BashOperator`.

```
# The execution date as YYYY-MM-DD
date = "{{ ds }}"
t = BashOperator(
    task_id='test_env',
    bash_command='/tmp/test.sh ',
    dag=dag,
    env={'EXECUTION_DATE': date})
```

Here, `{{ ds }}` is a macro, and because the `env` parameter of the `BashOperator` is templated with Jinja, the execution date will be available as an environment variable named `EXECUTION_DATE` in your Bash script.

You can use Jinja templating with every parameter that is marked as “templated” in the documentation.

3.8.3 Packaged dags

While often you will specify dags in a single `.py` file it might sometimes be required to combine dag and its dependencies. For example, you might want to combine several dags together to version them together or you might want to manage them together or you might need an extra module that is not available by default on the system you are running airflow on. To allow this you can create a zip file that contains the dag(s) in the root of the zip file and have the extra modules unpacked in directories.

For instance you can create a zip file that looks like this:

```
my_dag1.py
my_dag2.py
package1/__init__.py
package1/functions.py
```

Airflow will scan the zip file and try to load `my_dag1.py` and `my_dag2.py`. It will not go into subdirectories as these are considered to be potential packages.

In case you would like to add module dependencies to your DAG you basically would do the same, but then it is more to use a `virtualenv` and `pip`.

```
virtualenv zip_dag
source zip_dag/bin/activate

mkdir zip_dag_contents
cd zip_dag_contents

pip install --install-option="--install-lib=$PWD" my_useful_package
cp ~/my_dag.py .
```

```
zip -r zip_dag.zip *
```

Note: the zip file will be inserted at the beginning of module search list (sys.path) and as such it will be available to any other code that resides within the same interpreter.

Note: packaged dags cannot be used with pickling turned on.

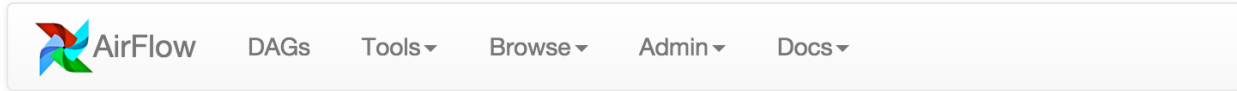
Note: packaged dags cannot contain dynamic libraries (eg. libz.so) these need to be available on the system if a module needs those. In other words only pure python modules can be packaged.

3.9 Data Profiling

Part of being productive with data is having the right weapons to profile the data you are working with. Airflow provides a simple query interface to write SQL and get results quickly, and a charting application letting you visualize data.

3.9.1 Adhoc Queries

The adhoc query UI allows for simple SQL interactions with the database connections registered in Airflow.



Ad Hoc Query

airflow_db

1 `SELECT * FROM task_instance LIMIT 1000`

Show entries Search:

task_id ▲	dag_id ◆	execution_date ◆	start_date ◆	end_date ◆	duration ◆	state
agent_performance_for_lantern	core_cx	2014-11-22 00:00:00	2014-11-23 22:50:51	2014-11-23 22:54:54	243	success
agent_performance_for_lantern	core_cx	2014-11-23 00:00:00	2014-11-24 23:04:53	2014-11-24 23:08:58	245	success
agent_performance_for_lantern	core_cx	2014-11-24 00:00:00	2014-11-26 00:25:46	2014-11-26 00:29:27	220	success
agent_performance_for_lantern	core_cx	2014-11-25 00:00:00	2014-11-29 00:05:02	2014-11-29 00:09:07	244	success
agent_performance_for_lantern	core_cx	2014-11-26 00:00:00	2014-11-29 01:46:23	2014-11-29 02:05:50	1167	success
agent_performance_for_lantern	core_cx	2014-11-27 00:00:00	2014-11-29 18:06:04	2014-11-29 18:10:04	239	success
agent_performance_for_lantern	core_cx	2014-11-28 00:00:00	2014-11-29 18:20:12	2014-11-29 18:23:45	212	success
agent_performance_for_lantern	core_cx	2014-11-29 00:00:00	2014-12-01 05:46:37	2014-12-01 05:50:32	234	success

3.9.2 Charts

A simple UI built on top of flask-admin and highcharts allows building data visualizations and charts easily. Fill in a form with a label, SQL, chart type, pick a source database from your environment's connections, select a few other options, and save it for later use.

You can even use the same templating and macros available when writing airflow pipelines, parameterizing your queries and modifying parameters directly in the URL.

These charts are basic, but they're easy to create, modify and share.

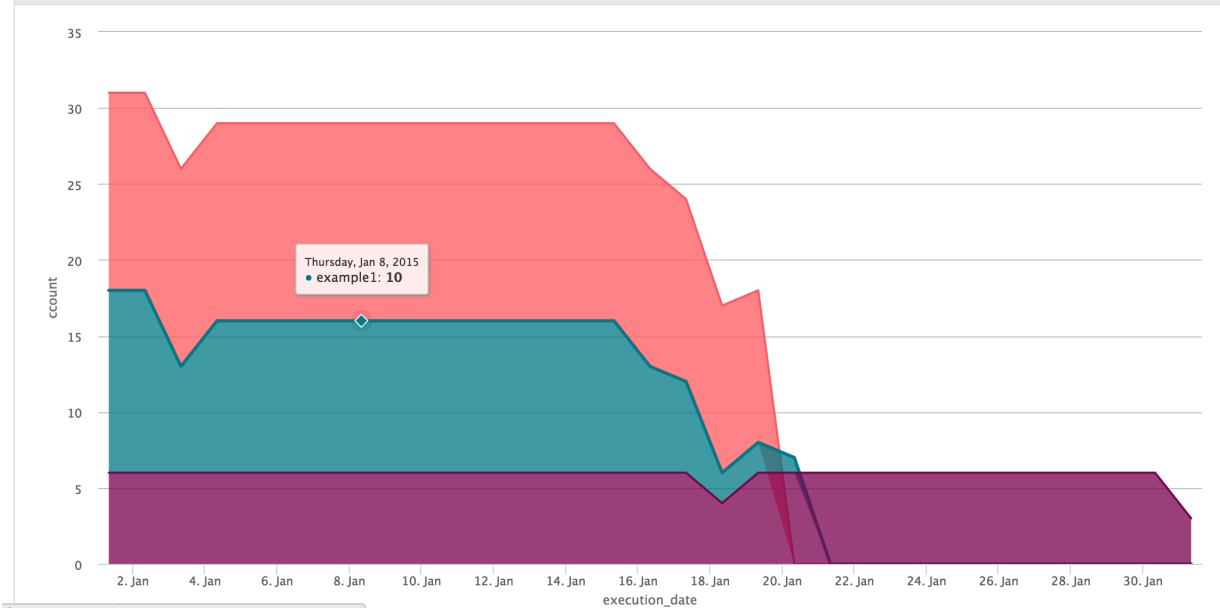
3.9.2.1 Chart Screenshot

Tasks

SQL

```
SELECT dag_id, execution_date, count(*) as ccount
FROM task_instance
GROUP BY dag_id, execution_date
```

Chart



3.9.2.2 Chart Form Screenshot

Label

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

Owner

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

Source Database

Chart Type

Line 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 prop

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 t

SQL

1

SELECT series, x, y FROM table

3.10 Command Line Interface

Airflow has a very rich command line interface that allows for many types of operation on a DAG, starting services, and supporting development and testing.

```
usage: airflow [-h]
               {resetdb,render,variables,connections,pause,task_failed_deps,version,
↪trigger_dag,initdb,test,unpause,dag_state,run,list_tasks,backfill,list_dags,
↪kerberos,worker,webserver,flower,scheduler,task_state,pool,serve_logs,clear,
↪upgradedb}
               ...
```

3.10.1 Positional Arguments

subcommand

Possible choices: resetdb, render, variables, connections, pause, task_failed_deps, version, trigger_dag, initdb, test, unpause, dag_state,

run, list_tasks, backfill, list_dags, kerberos, worker, webserver, flower, scheduler,
task_state, pool, serve_logs, clear, upgradedb
sub-command help

3.10.2 Sub-commands:

3.10.2.1 resetdb

Burn down and rebuild the metadata database

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

Named Arguments

-y, --yes	Do not prompt to confirm reset. Use with care!
	Default: False

3.10.2.2 render

Render a task instance's template(s)

```
airflow render [-h] [-sd SUBDIR] dag_id task_id execution_date
```

Positional Arguments

dag_id	The id of the dag
task_id	The id of the task
execution_date	The execution date of the DAG

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag
	Default: /home/docs/airflow/dags

3.10.2.3 variables

CRUD operations on variables

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


Named Arguments

-s, --set	Set a variable
-g, --get	Get value of a variable
-j, --json	Deserialize JSON variable
	Default: False
-d, --default	Default value returned if variable does not exist
-i, --import	Import variables from JSON file
-e, --export	Export variables to JSON file
-x, --delete	Delete a variable

3.10.2.4 connections

List/Add/Delete connections

```
airflow connections [-h] [-l] [-a] [-d] [--conn_id CONN_ID]
                   [--conn_uri CONN_URI] [--conn_extra CONN_EXTRA]
```

Named Arguments

-l, --list	List all connections
	Default: False
-a, --add	Add a connection
	Default: False
-d, --delete	Delete a connection
	Default: False
--conn_id	Connection id, required to add/delete a connection
--conn_uri	Connection URI, required to add a connection
--conn_extra	Connection <i>Extra</i> field, optional when adding a connection

3.10.2.5 pause

Pause a DAG

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

Positional Arguments

dag_id	The id of the dag
---------------	-------------------

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag
	Default: /home/docs/airflow/dags

3.10.2.6 task_failed_deps

Returns the unmet dependencies for a task instance from the perspective of the scheduler. In other words, why a task instance doesn't get scheduled and then queued by the scheduler, and then run by an executor).

```
airflow task_failed_deps [-h] [-sd SUBDIR] dag_id task_id execution_date
```

Positional Arguments

dag_id	The id of the dag
task_id	The id of the task
execution_date	The execution date of the DAG

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag
	Default: /home/docs/airflow/dags

3.10.2.7 version

Show the version

```
airflow version [-h]
```

3.10.2.8 trigger_dag

Trigger a DAG run

```
airflow trigger_dag [-h] [-sd SUBDIR] [-r RUN_ID] [-c CONF] [-e EXEC_DATE]
                        dag_id
```

Positional Arguments

dag_id	The id of the dag
---------------	-------------------

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
-r, --run_id	Helps to identify this run
-c, --conf	JSON string that gets pickled into the DagRun's conf attribute
-e, --exec_date	The execution date of the DAG

3.10.2.9 initdb

Initialize the metadata database

```
airflow initdb [-h]
```

3.10.2.10 test

Test a task instance. This will run a task without checking for dependencies or recording it's state in the database.

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

Positional Arguments

dag_id	The id of the dag
task_id	The id of the task
execution_date	The execution date of the DAG

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
-dr, --dry_run	Perform a dry run Default: False
-tp, --task_params	Sends a JSON params dict to the task

3.10.2.11 unpause

Resume a paused DAG

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

Positional Arguments

dag_id The id of the dag

Named Arguments

-sd, --subdir File location or directory from which to look for the dag
Default: /home/docs/airflow/dags

3.10.2.12 dag_state

Get the status of a dag run

```
airflow dag_state [-h] [-sd SUBDIR] dag_id execution_date
```

Positional Arguments

dag_id The id of the dag
execution_date The execution date of the DAG

Named Arguments

-sd, --subdir File location or directory from which to look for the dag
Default: /home/docs/airflow/dags

3.10.2.13 run

Run a single task instance

```
airflow run [-h] [-sd SUBDIR] [-m] [-f] [--pool POOL] [--cfg_path CFG_PATH]
            [-l] [-A IGNORE_ALL_DEPENDENCIES] [-i] [-I] [--ship_dag]
            [-p PICKLE]
            dag_id task_id execution_date
```

Positional Arguments

dag_id The id of the dag
task_id The id of the task
execution_date The execution date of the DAG

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
-m, --mark_success	Mark jobs as succeeded without running them Default: False
-f, --force	Ignore previous task instance state, rerun regardless if task already succeeded/failed Default: False
--pool	Resource pool to use
--cfg_path	Path to config file to use instead of airflow.cfg
-l, --local	Run the task using the LocalExecutor Default: False
-A, --ignore_all_dependencies	Ignores all non-critical dependencies, including ignore_ti_state and ignore_task_depsstore_true
-i, --ignore_dependencies	Ignore task-specific dependencies, e.g. upstream, depends_on_past, and retry delay dependencies Default: False
-I, --ignore_depends_on_past	Ignore depends_on_past dependencies (but respect upstream dependencies) Default: False
--ship_dag	Pickles (serializes) the DAG and ships it to the worker Default: False
-p, --pickle	Serialized pickle object of the entire dag (used internally)

3.10.2.14 list_tasks

List the tasks within a DAG

```
airflow list_tasks [-h] [-t] [-sd SUBDIR] dag_id
```

Positional Arguments

dag_id	The id of the dag
---------------	-------------------

Named Arguments

-t, --tree	Tree view Default: False
-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags

3.10.2.15 backfill

Run subsections of a DAG for a specified date range

```
airflow backfill [-h] [-t TASK_REGEX] [-s START_DATE] [-e END_DATE] [-m] [-l]
                 [-x] [-a] [-i] [-I] [-sd SUBDIR] [--pool POOL] [-dr]
                 dag_id
```

Positional Arguments

dag_id	The id of the dag
---------------	-------------------

Named Arguments

-t, --task_regex	The regex to filter specific task_ids to backfill (optional)
-s, --start_date	Override start_date YYYY-MM-DD
-e, --end_date	Override end_date YYYY-MM-DD
-m, --mark_success	Mark jobs as succeeded without running them Default: False
-l, --local	Run the task using the LocalExecutor Default: False
-x, --donot_pickle	Do not attempt to pickle the DAG object to send over to the workers, just tell the workers to run their version of the code. Default: False
-a, --include_adhoc	Include dags with the adhoc parameter. Default: False
-i, --ignore_dependencies	Skip upstream tasks, run only the tasks matching the regexp. Only works in conjunction with task_regex Default: False
-I, --ignore_first_depends_on_past	Ignores depends_on_past dependencies for the first set of tasks only (subsequent executions in the backfill DO respect depends_on_past). Default: False
-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
--pool	Resource pool to use
-dr, --dry_run	Perform a dry run Default: False

3.10.2.16 list_dags

List all the DAGs

```
airflow list_dags [-h] [-sd SUBDIR] [-r]
```

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
-r, --report	Show DagBag loading report Default: False

3.10.2.17 kerberos

Start a kerberos ticket renewer

```
airflow kerberos [-h] [-kt [KEYTAB]] [--pid [PID]] [-D] [--stdout STDOUT]
                  [--stderr STDERR] [-l LOG_FILE]
                  [principal]
```

Positional Arguments

principal	kerberos principal Default: airflow
------------------	--

Named Arguments

-kt, --keytab	keytab Default: airflow.keytab
--pid	PID file location
-D, --daemon	Daemonize instead of running in the foreground Default: False
--stdout	Redirect stdout to this file
--stderr	Redirect stderr to this file
-l, --log-file	Location of the log file

3.10.2.18 worker

Start a Celery worker node

```
airflow worker [-h] [-p] [-q QUEUES] [-c CONCURRENCY] [--pid [PID]] [-D]
               [--stdout STDOUT] [--stderr STDERR] [-l LOG_FILE]
```

Named Arguments

-p, --do_pickle	Attempt to pickle the DAG object to send over to the workers, instead of letting workers run their version of the code. Default: False
-q, --queues	Comma delimited list of queues to serve Default: default
-c, --concurrency	The number of worker processes Default: 16
--pid	PID file location
-D, --daemon	Daemonize instead of running in the foreground Default: False
--stdout	Redirect stdout to this file
--stderr	Redirect stderr to this file
-l, --log-file	Location of the log file

3.10.2.19 webserver

Start a Airflow webserver instance

```
airflow webserver [-h] [-p PORT] [-w WORKERS]
                  [-k {sync,eventlet,gevent,tornado}] [-t WORKER_TIMEOUT]
                  [-hn HOSTNAME] [--pid [PID]] [-D] [--stdout STDOUT]
                  [--stderr STDERR] [-A ACCESS_LOGFILE] [-E ERROR_LOGFILE]
                  [-l LOG_FILE] [--ssl_cert SSL_CERT] [--ssl_key SSL_KEY] [-d]
```

Named Arguments

-p, --port	The port on which to run the server Default: 8080
-w, --workers	Number of workers to run the webserver on Default: 4
-k, --workerclass	Possible choices: sync, eventlet, gevent, tornado The worker class to use for Gunicorn Default: sync
-t, --worker_timeout	The timeout for waiting on webserver workers Default: 120
-hn, --hostname	Set the hostname on which to run the web server Default: 0.0.0.0
--pid	PID file location

-D, --daemon	Daemonize instead of running in the foreground Default: False
--stdout	Redirect stdout to this file
--stderr	Redirect stderr to this file
-A, --access_logfile	The logfile to store the webserver access log. Use '-' to print to stderr. Default: -
-E, --error_logfile	The logfile to store the webserver error log. Use '-' to print to stderr. Default: -
-l, --log-file	Location of the log file
--ssl_cert	Path to the SSL certificate for the webserver
--ssl_key	Path to the key to use with the SSL certificate
-d, --debug	Use the server that ships with Flask in debug mode Default: False

3.10.2.20 flower

Start a Celery Flower

```
airflow flower [-h] [-hn HOSTNAME] [-p PORT] [-fc FLOWER_CONF] [-a BROKER_API]
               [--pid [PID]] [-D] [--stdout STDOUT] [--stderr STDERR]
               [-l LOG_FILE]
```

Named Arguments

-hn, --hostname	Set the hostname on which to run the server Default: 0.0.0.0
-p, --port	The port on which to run the server Default: 5555
-fc, --flower_conf	Configuration file for flower
-a, --broker_api	Broker api
--pid	PID file location
-D, --daemon	Daemonize instead of running in the foreground Default: False
--stdout	Redirect stdout to this file
--stderr	Redirect stderr to this file
-l, --log-file	Location of the log file

3.10.2.21 scheduler

Start a scheduler instance

```
airflow scheduler [-h] [-d DAG_ID] [-sd SUBDIR] [-r RUN_DURATION]
                  [-n NUM_RUNS] [-p] [--pid [PID]] [-D] [--stdout STDOUT]
                  [--stderr STDERR] [-l LOG_FILE]
```

Named Arguments

-d, --dag_id	The id of the dag to run
-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
-r, --run-duration	Set number of seconds to execute before exiting
-n, --num_runs	Set the number of runs to execute before exiting Default: -1
-p, --do_pickle	Attempt to pickle the DAG object to send over to the workers, instead of letting workers run their version of the code. Default: False
--pid	PID file location
-D, --daemon	Daemonize instead of running in the foreground Default: False
--stdout	Redirect stdout to this file
--stderr	Redirect stderr to this file
-l, --log-file	Location of the log file

3.10.2.22 task_state

Get the status of a task instance

```
airflow task_state [-h] [-sd SUBDIR] dag_id task_id execution_date
```

Positional Arguments

dag_id	The id of the dag
task_id	The id of the task
execution_date	The execution date of the DAG

Named Arguments

-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
----------------------	---

3.10.2.23 pool

CRUD operations on pools

```
airflow pool [-h] [-s NAME SLOT_COUNT POOL_DESCRIPTION] [-g NAME] [-x NAME]
```

Named Arguments

-s, --set	Set pool slot count and description, respectively
-g, --get	Get pool info
-x, --delete	Delete a pool

3.10.2.24 serve_logs

Serve logs generate by worker

```
airflow serve_logs [-h]
```

3.10.2.25 clear

Clear a set of task instance, as if they never ran

```
airflow clear [-h] [-t TASK_REGEX] [-s START_DATE] [-e END_DATE] [-sd SUBDIR]
               [-u] [-d] [-c] [-f] [-r] [-x]
               dag_id
```

Positional Arguments

dag_id	The id of the dag
---------------	-------------------

Named Arguments

-t, --task_regex	The regex to filter specific task_ids to backfill (optional)
-s, --start_date	Override start_date YYYY-MM-DD
-e, --end_date	Override end_date YYYY-MM-DD
-sd, --subdir	File location or directory from which to look for the dag Default: /home/docs/airflow/dags
-u, --upstream	Include upstream tasks Default: False
-d, --downstream	Include downstream tasks Default: False
-c, --no_confirm	Do not request confirmation Default: False

- f, --only_failed** Only failed jobs
Default: False
- r, --only_running** Only running jobs
Default: False
- x, --exclude_subdags** Exclude subdags
Default: False

3.10.2.26 upgradedb

Upgrade the metadata database to latest version

```
airflow upgradedb [-h]
```

3.11 Scheduling & Triggers

The Airflow scheduler monitors all tasks and all DAGs, and triggers the task instances whose dependencies have been met. Behind the scenes, it monitors and stays in sync with a folder for all DAG objects it may contain, and periodically (every minute or so) inspects active tasks to see whether they can be triggered.

The Airflow scheduler is designed to run as a persistent service in an Airflow production environment. To kick it off, all you need to do is execute `airflow scheduler`. It will use the configuration specified in `airflow.cfg`.

Note that if you run a DAG on a `schedule_interval` of one day, the run stamped 2016-01-01 will be trigger soon after 2016-01-01T23:59. In other words, the job instance is started once the period it covers has ended.

Let's Repeat That The scheduler runs your job one `schedule_interval` AFTER the start date, at the END of the period.

The scheduler starts an instance of the executor specified in the your `airflow.cfg`. If it happens to be the `LocalExecutor`, tasks will be executed as subprocesses; in the case of `CeleryExecutor` and `MesosExecutor`, tasks are executed remotely.

To start a scheduler, simply run the command:

```
airflow scheduler
```

3.11.1 DAG Runs

A DAG Run is an object representing an instantiation of the DAG in time.

Each DAG may or may not have a schedule, which informs how DAG Runs are created. `schedule_interval` is defined as a DAG arguments, and receives preferably a [cron expression](#) as a `str`, or a `datetime.timedelta` object. Alternatively, you can also use one of these cron “preset”:

preset	Run once a year at midnight of January 1	cron
None	Don't schedule, use for exclusively "externally triggered" DAGs	
@once	Schedule once and only once	
@hourly	Run once an hour at the beginning of the hour	0 * * * *
@daily	Run once a day at midnight	0 0 * * *
@weekly	Run once a week at midnight on Sunday morning	0 0 * * 0
@monthly	Run once a month at midnight of the first day of the month	0 0 1 * *
@yearly	Run once a year at midnight of January 1	0 0 1 1 *

Your DAG will be instantiated for each schedule, while creating a DAG Run entry for each schedule.

DAG runs have a state associated to them (running, failed, success) and informs the scheduler on which set of schedules should be evaluated for task submissions. Without the metadata at the DAG run level, the Airflow scheduler would have much more work to do in order to figure out what tasks should be triggered and come to a crawl. It might also create undesired processing when changing the shape of your DAG, by say adding in new tasks.

3.11.2 Backfill and Catchup

An Airflow DAG with a `start_date`, possibly an `end_date`, and a `schedule_interval` defines a series of intervals which the scheduler turn into individual Dag Runs and execute. A key capability of Airflow is that these DAG Runs are atomic, idempotent items, and the scheduler, by default, will examine the lifetime of the DAG (from start to end/now, one interval at a time) and kick off a DAG Run for any interval that has not been run (or has been cleared). This concept is called Catchup.

If your DAG is written to handle it's own catchup (IE not limited to the interval, but instead to "Now" for instance.), then you will want to turn catchup off (Either on the DAG itself with `dag.catchup = False`) or by default at the configuration file level with `catchup_by_default = False`. What this will do, is to instruct the scheduler to only create a DAG Run for the most current instance of the DAG interval series.

```
"""
Code that goes along with the Airflow tutorial located at:
https://github.com/airbnb/airflow/blob/master/airflow/example_dags/tutorial.py
"""
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
from datetime import datetime, timedelta

default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date': datetime(2015, 12, 1),
    'email': ['airflow@example.com'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=5),
    'schedule_interval': '@hourly',
}

dag = DAG('tutorial', catchup=False, default_args=default_args)
```

In the example above, if the DAG is picked up by the scheduler daemon on 2016-01-02 at 6 AM, (or from the command line), a single DAG Run will be created, with an `execution_date` of 2016-01-01, and the next one will be created just after midnight on the morning of 2016-01-03 with an execution date of 2016-01-02.

If the `dag.catchup` value had been `True` instead, the scheduler would have created a DAG Run for each completed interval between 2015-12-01 and 2016-01-02 (but not yet one for 2016-01-02, as that interval hasn't completed) and the scheduler will execute them sequentially. This behavior is great for atomic datasets that can easily be split into periods. Turning catchup off is great if your DAG Runs perform backfill internally.

3.11.3 External Triggers

Note that DAG Runs can also be created manually through the CLI while running an `airflow trigger_dag` command, where you can define a specific `run_id`. The DAG Runs created externally to the scheduler get associated to the trigger's timestamp, and will be displayed in the UI alongside scheduled DAG runs.

3.11.4 To Keep in Mind

- The first DAG Run is created based on the minimum `start_date` for the tasks in your DAG.
- Subsequent DAG Runs are created by the scheduler process, based on your DAG's `schedule_interval`, sequentially.
- When clearing a set of tasks' state in hope of getting them to re-run, it is important to keep in mind the DAG Run's state too as it defines whether the scheduler should look into triggering tasks for that run.

Here are some of the ways you can **unblock tasks**:

- From the UI, you can **clear** (as in delete the status of) individual task instances from the task instances dialog, while defining whether you want to include the past/future and the upstream/downstream dependencies. Note that a confirmation window comes next and allows you to see the set you are about to clear.
- The CLI command `airflow clear -h` has lots of options when it comes to clearing task instance states, including specifying date ranges, targeting `task_ids` by specifying a regular expression, flags for including upstream and downstream relatives, and targeting task instances in specific states (`failed`, or `success`)
- Marking task instances as successful can be done through the UI. This is mostly to fix false negatives, or for instance when the fix has been applied outside of Airflow.
- The `airflow backfill` CLI subcommand has a flag to `--mark_success` and allows selecting subsections of the DAG as well as specifying date ranges.

3.12 Plugins

Airflow has a simple plugin manager built-in that can integrate external features to its core by simply dropping files in your `$AIRFLOW_HOME/plugins` folder.

The python modules in the `plugins` folder get imported, and **hooks**, **operators**, **macros**, **executors** and **web views** get integrated to Airflow's main collections and become available for use.

3.12.1 What for?

Airflow offers a generic toolbox for working with data. Different organizations have different stacks and different needs. Using Airflow plugins can be a way for companies to customize their Airflow installation to reflect their ecosystem.

Plugins can be used as an easy way to write, share and activate new sets of features.

There's also a need for a set of more complex applications to interact with different flavors of data and metadata.

Examples:

- A set of tools to parse Hive logs and expose Hive metadata (CPU /IO / phases/ skew /...)
- An anomaly detection framework, allowing people to collect metrics, set thresholds and alerts
- An auditing tool, helping understand who accesses what
- A config-driven SLA monitoring tool, allowing you to set monitored tables and at what time they should land, alert people, and expose visualizations of outages
- ...

3.12.2 Why build on top of Airflow?

Airflow has many components that can be reused when building an application:

- A web server you can use to render your views
- A metadata database to store your models
- Access to your databases, and knowledge of how to connect to them
- An array of workers that your application can push workload to
- Airflow is deployed, you can just piggy back on it's deployment logistics
- Basic charting capabilities, underlying libraries and abstractions

3.12.3 Interface

To create a plugin you will need to derive the `airflow.plugins_manager.AirflowPlugin` class and reference the objects you want to plug into Airflow. Here's what the class you need to derive looks like:

```
class AirflowPlugin(object):
    # The name of your plugin (str)
    name = None
    # A list of class(es) derived from BaseOperator
    operators = []
    # A list of class(es) derived from BaseHook
    hooks = []
    # A list of class(es) derived from BaseExecutor
    executors = []
    # A list of references to inject into the macros namespace
    macros = []
    # A list of objects created from a class derived
    # from flask_admin.BaseView
    admin_views = []
    # A list of Blueprint object created from flask.Blueprint
    flask_blueprints = []
    # A list of menu links (flask_admin.base.MenuLink)
    menu_links = []
```

3.12.4 Example

The code below defines a plugin that injects a set of dummy object definitions in Airflow.

```

# This is the class you derive to create a plugin
from airflow.plugins_manager import AirflowPlugin

from flask import Blueprint
from flask_admin import BaseView, expose
from flask_admin.base import MenuLink

# Importing base classes that we need to derive
from airflow.hooks.base_hook import BaseHook
from airflow.models import BaseOperator
from airflow.executors.base_executor import BaseExecutor

# Will show up under airflow.hooks.test_plugin.PluginHook
class PluginHook(BaseHook):
    pass

# Will show up under airflow.operators.test_plugin.PluginOperator
class PluginOperator(BaseOperator):
    pass

# Will show up under airflow.executors.test_plugin.PluginExecutor
class PluginExecutor(BaseExecutor):
    pass

# Will show up under airflow.macros.test_plugin.plugin_macro
def plugin_macro():
    pass

# Creating a flask admin BaseView
class TestView(BaseView):
    @expose('/')
    def test(self):
        # in this example, put your test_plugin/test.html template at airflow/plugins/
        ↪ templates/test_plugin/test.html
        return self.render("test_plugin/test.html", content="Hello galaxy!")
v = TestView(category="Test Plugin", name="Test View")

# Creating a flask blueprint to intergrate the templates and static folder
bp = Blueprint(
    "test_plugin", __name__,
    template_folder='templates', # registers airflow/plugins/templates as a Jinja_
    ↪ template folder
    static_folder='static',
    static_url_path='/static/test_plugin')

ml = MenuLink(
    category='Test Plugin',
    name='Test Menu Link',
    url='http://pythonhosted.org/airflow/')

# Defining the plugin class
class AirflowTestPlugin(AirflowPlugin):
    name = "test_plugin"
    operators = [PluginOperator]
    hooks = [PluginHook]
    executors = [PluginExecutor]
    macros = [plugin_macro]
    admin_views = [v]

```



```
flask_blueprints = [bp]
menu_links = [ml]
```

3.13 Security

By default, all gates are opened. An easy way to restrict access to the web application is to do it at the network level, or by using SSH tunnels.

It is however possible to switch on authentication by either using one of the supplied backends or creating your own.

3.13.1 Web Authentication

3.13.1.1 Password

One of the simplest mechanisms for authentication is requiring users to specify a password before logging in. Password authentication requires the use of the `password` subpackage in your requirements file. Password hashing uses `bcrypt` before storing passwords.

```
[webserver]
authenticate = True
auth_backend = airflow.contrib.auth.backends.password_auth
```

When password auth is enabled, an initial user credential will need to be created before anyone can login. An initial user was not created in the migrations for this authentication backend to prevent default Airflow installations from attack. Creating a new user has to be done via a Python REPL on the same machine Airflow is installed.

```
# navigate to the airflow installation directory
$ cd ~/airflow
$ python
Python 2.7.9 (default, Feb 10 2015, 03:28:08)
Type "help", "copyright", "credits" or "license" for more information.
>>> import airflow
>>> from airflow import models, settings
>>> from airflow.contrib.auth.backends.password_auth import PasswordUser
>>> user = PasswordUser(models.User())
>>> user.username = 'new_user_name'
>>> user.email = 'new_user_email@example.com'
>>> user.password = 'set_the_password'
>>> session = settings.Session()
>>> session.add(user)
>>> session.commit()
>>> session.close()
>>> exit()
```

3.13.1.2 LDAP

To turn on LDAP authentication configure your `airflow.cfg` as follows. Please note that the example uses an encrypted connection to the ldap server as you probably do not want passwords be readable on the network level. It is however possible to configure without encryption if you really want to.

Additionally, if you are using Active Directory, and are not explicitly specifying an OU that your users are in, you will need to change `search_scope` to “SUBTREE”.

Valid search_scope options can be found in the [ldap3 Documentation](#)

```
[webserver]
authenticate = True
auth_backend = airflow.contrib.auth.backends.ldap_auth

[ldap]
# set a connection without encryption: uri = ldap://<your.ldap.server>:<port>
uri = ldaps://<your.ldap.server>:<port>
user_filter = objectClass=*
# in case of Active Directory you would use: user_name_attr = sAMAccountName
user_name_attr = uid
# group_member_attr should be set accordingly with *_filter
# eg :
#     group_member_attr = groupMembership
#     superuser_filter = groupMembership=CN=airflow-super-users...
group_member_attr = memberOf
superuser_filter = memberOf=CN=airflow-super-users,OU=Groups,OU=RWC,OU=US,OU=NORAM,
↳DC=example,DC=com
data_profiler_filter = memberOf=CN=airflow-data-profilers,OU=Groups,OU=RWC,OU=US,
↳OU=NORAM,DC=example,DC=com
bind_user = cn=Manager,dc=example,dc=com
bind_password = insecure
basedn = dc=example,dc=com
cacert = /etc/ca/ldap_ca.crt
# Set search_scope to one of them: BASE, LEVEL, SUBTREE
# Set search_scope to SUBTREE if using Active Directory, and not specifying an_
↳Organizational Unit
search_scope = LEVEL
```

The superuser_filter and data_profiler_filter are optional. If defined, these configurations allow you to specify LDAP groups that users must belong to in order to have superuser (admin) and data-profiler permissions. If undefined, all users will be superusers and data profilers.

3.13.1.3 Roll your own

Airflow uses flask_login and exposes a set of hooks in the airflow.default_login module. You can alter the content and make it part of the PYTHONPATH and configure it as a backend in airflow.cfg.

```
[webserver]
authenticate = True
auth_backend = mypackage.auth
```

3.13.2 Multi-tenancy

You can filter the list of dags in webserver by owner name when authentication is turned on by setting webserver:filter_by_owner in your config. With this, a user will see only the dags which it is owner of, unless it is a superuser.

```
[webserver]
filter_by_owner = True
```

3.13.3 Kerberos

Airflow has initial support for Kerberos. This means that airflow can renew kerberos tickets for itself and store it in the ticket cache. The hooks and dags can make use of ticket to authenticate against kerberized services.

3.13.3.1 Limitations

Please note that at this time, not all hooks have been adjusted to make use of this functionality. Also it does not integrate kerberos into the web interface and you will have to rely on network level security for now to make sure your service remains secure.

Celery integration has not been tried and tested yet. However, if you generate a key tab for every host and launch a ticket renewer next to every worker it will most likely work.

3.13.3.2 Enabling kerberos

Airflow

To enable kerberos you will need to generate a (service) key tab.

```
# in the kadmin.local or kadmin shell, create the airflow principal
kadmin: addprinc -randkey airflow/fully.qualified.domain.name@YOUR-REALM.COM

# Create the airflow keytab file that will contain the airflow principal
kadmin: xst -norandkey -k airflow.keytab airflow/fully.qualified.domain.name
```

Now store this file in a location where the airflow user can read it (chmod 600). And then add the following to your `airflow.cfg`

```
[core]
security = kerberos

[kerberos]
keytab = /etc/airflow/airflow.keytab
reinit_frequency = 3600
principal = airflow
```

Launch the ticket renewer by

```
# run ticket renewer
airflow kerberos
```

Hadoop

If want to use impersonation this needs to be enabled in `core-site.xml` of your hadoop config.

```
<property>
  <name>hadoop.proxyuser.airflow.groups</name>
  <value>*</value>
</property>

<property>
  <name>hadoop.proxyuser.airflow.users</name>
  <value>*</value>
```

```
</property>

<property>
  <name>hadoop.proxyuser.airflow.hosts</name>
  <value>*</value>
</property>
```

Of course if you need to tighten your security replace the asterisk with something more appropriate.

3.13.3.3 Using kerberos authentication

The hive hook has been updated to take advantage of kerberos authentication. To allow your DAGs to use it, simply update the connection details with, for example:

```
{ "use_beeline": true, "principal": "hive/_HOST@EXAMPLE.COM" }
```

Adjust the principal to your settings. The `_HOST` part will be replaced by the fully qualified domain name of the server.

You can specify if you would like to use the dag owner as the user for the connection or the user specified in the login section of the connection. For the login user, specify the following as extra:

```
{ "use_beeline": true, "principal": "hive/_HOST@EXAMPLE.COM", "proxy_user": "login" }
```

For the DAG owner use:

```
{ "use_beeline": true, "principal": "hive/_HOST@EXAMPLE.COM", "proxy_user": "owner" }
```

and in your DAG, when initializing the `HiveOperator`, specify:

```
run_as_owner=True
```

3.13.4 OAuth Authentication

3.13.4.1 GitHub Enterprise (GHE) Authentication

The GitHub Enterprise authentication backend can be used to authenticate users against an installation of GitHub Enterprise using OAuth2. You can optionally specify a team whitelist (composed of slug cased team names) to restrict login to only members of those teams.

```
[webserver]
authenticate = True
auth_backend = airflow.contrib.auth.backends.github_enterprise_auth

[github_enterprise]
host = github.example.com
client_id = oauth_key_from_github_enterprise
client_secret = oauth_secret_from_github_enterprise
oauth_callback_route = /example/ghe_oauth/callback
allowed_teams = 1, 345, 23
```

Note: If you do not specify a team whitelist, anyone with a valid account on your GHE installation will be able to login to Airflow.

Setting up GHE Authentication

An application must be setup in GHE before you can use the GHE authentication backend. In order to setup an application:

1. Navigate to your GHE profile
2. Select 'Applications' from the left hand nav
3. Select the 'Developer Applications' tab
4. Click 'Register new application'
5. Fill in the required information (the 'Authorization callback URL' must be fully qualified e.g. http://airflow.example.com/example/ghe_oauth/callback)
6. Click 'Register application'
7. Copy 'Client ID', 'Client Secret', and your callback route to your airflow.cfg according to the above example

3.13.4.2 Google Authentication

The Google authentication backend can be used to authenticate users against Google using OAuth2. You must specify a domain to restrict login to only members of that domain.

```
[webserver]
authenticate = True
auth_backend = airflow.contrib.auth.backends.google_auth

[google]
client_id = google_client_id
client_secret = google_client_secret
oauth_callback_route = /oauth2callback
domain = example.com
```

Setting up Google Authentication

An application must be setup in the Google API Console before you can use the Google authentication backend. In order to setup an application:

1. Navigate to <https://console.developers.google.com/apis/>
2. Select 'Credentials' from the left hand nav
2. Select 'Credentials' from the left hand nav
3. Click 'Create credentials' and choose 'OAuth client ID'
4. Choose 'Web application'
5. Fill in the required information (the 'Authorized redirect URIs' must be fully qualified e.g. <http://airflow.example.com/oauth2callback>)
6. Click 'Create'
7. Copy 'Client ID', 'Client Secret', and your redirect URI to your airflow.cfg according to the above example

3.13.5 SSL

SSL can be enabled by providing a certificate and key. Once enabled, be sure to use "https://" in your browser.

```
[webserver]
web_server_ssl_cert = <path to cert>
web_server_ssl_key = <path to key>
```

Enabling SSL will not automatically change the web server port. If you want to use the standard port 443, you'll need to configure that too. Be aware that super user privileges (or `cap_net_bind_service` on Linux) are required to listen on port 443.

```
# Optionally, set the server to listen on the standard SSL port.
web_server_port = 443
base_url = http://<hostname or IP>:443
```

3.13.6 Impersonation

Airflow has the ability to impersonate a unix user while running task instances based on the task's `run_as_user` parameter, which takes a user's name.

NOTE: For impersonations to work, Airflow must be run with `sudo` as subtasks are run with `sudo -u` and permissions of files are changed. Furthermore, the unix user needs to exist on the worker. Here is what a simple sudoers file entry could look like to achieve this, assuming as airflow is running as the *airflow* user. Note that this means that the airflow user must be trusted and treated the same way as the root user.

```
airflow ALL=(ALL) NOPASSWD: ALL
```

Subtasks with impersonation will still log to the same folder, except that the files they log to will have permissions changed such that only the unix user can write to it.

3.13.6.1 Default Impersonation

To prevent tasks that don't use impersonation to be run with `sudo` privileges, you can set the `core:default_impersonation` config which sets a default user impersonate if `run_as_user` is not set.

```
[core]
default_impersonation = airflow
```

3.14 Experimental Rest API

Airflow exposes an experimental Rest API. It is available through the webserver. Endpoints are available at `/api/experimental/`. Please note that we expect the endpoint definitions to change.

3.14.1 Endpoints

This is a place holder until the swagger definitions are active

- `/api/experimental/dags/<DAG_ID>/tasks/<TASK_ID>` returns info for a task (GET).
- `/api/experimental/dags/<DAG_ID>/dag_runs` creates a dag_run for a given dag id (POST).

3.14.2 CLI

For some functions the cli can use the API. To configure the CLI to use the API when available configure as follows:

```
[cli]
api_client = airflow.api.client.json_client
endpoint_url = http://<WEBSERVER>:<PORT>
```

3.14.3 Authentication

Only Kerberos authentication is currently supported for the API. To enable this set the following in the configuration:

```
[api]
auth_backend = airflow.api.auth.backend.default

[kerberos]
keytab = <KEYTAB>
```

The Kerberos service is configured as *airflow/fully.qualified.domainname@REALM*. Make sure this principal exists in the keytab file.

3.15 Integration

- *Azure: Microsoft Azure*
- *AWS: Amazon Webservices*
- *GCP: Google Cloud Platform*

3.15.1 Azure: Microsoft Azure

Airflow has limited support for Microsoft Azure: interfaces exist only for Azure Blob Storage. Note that the Hook, Sensor and Operator are in the contrib section.

3.15.1.1 Azure Blob Storage

All classes communicate via the Windows Azure Storage Blob protocol. Make sure that a Airflow connection of type *wasb* exists. Authorization can be done by supplying a login (=Storage account name) and password (=KEY), or login and SAS token in the extra field (see connection *wasb_default* for an example).

- *WasbBlobSensor*: Checks if a blob is present on Azure Blob storage.
- *WasbPrefixSensor*: Checks if blobs matching a prefix are present on Azure Blob storage.
- *FileToWasbOperator*: Uploads a local file to a container as a blob.
- *WasbHook*: Interface with Azure Blob Storage.

WasbBlobSensor

WasbPrefixSensor

FileToWasbOperator

WasbHook

3.15.2 AWS: Amazon Webservices

3.15.3 Databricks

Databricks has contributed an Airflow operator which enables submitting runs to the Databricks platform. Internally the operator talks to the `api/2.0/jobs/runs/submit` endpoint.

3.15.3.1 DatabricksSubmitRunOperator

```
class airflow.contrib.operators.databricks_operator.DatabricksSubmitRunOperator (json=None,
                                                                                   spark_jar_task=None,
                                                                                   notebook_task=None,
                                                                                   new_cluster=None,
                                                                                   existing_cluster_id=None,
                                                                                   libraries=None,
                                                                                   run_name=None,
                                                                                   timeout_seconds=None,
                                                                                   databricks_conn_id=None,
                                                                                   polling_period_seconds=None,
                                                                                   databricks_retry_limit=None,
                                                                                   **kwargs)
```

Submits an Spark job run to Databricks using the `api/2.0/jobs/runs/submit` API endpoint.

There are two ways to instantiate this operator.

In the first way, you can take the JSON payload that you typically use to call the `api/2.0/jobs/runs/submit` endpoint and pass it directly to our `DatabricksSubmitRunOperator` through the `json` parameter. For example

```
json = {
    'new_cluster': {
        'spark_version': '2.1.0-db3-scala2.11',
        'num_workers': 2
    },
    'notebook_task': {
        'notebook_path': '/Users/airflow@example.com/PrepareData',
    },
}
notebook_run = DatabricksSubmitRunOperator(task_id='notebook_run', json=json)
```

Another way to accomplish the same thing is to use the named parameters of the `DatabricksSubmitRunOperator` directly. Note that there is exactly one named parameter for each top level parameter in the `runs/submit` endpoint. In this method, your code would look like this:

```
new_cluster = {
    'spark_version': '2.1.0-db3-scala2.11',
    'num_workers': 2
}
notebook_task = {
    'notebook_path': '/Users/airflow@example.com/PrepareData',
}
notebook_run = DatabricksSubmitRunOperator(
    task_id='notebook_run',
```



```
new_cluster=new_cluster,
notebook_task=notebook_task)
```

In the case where both the json parameter **AND** the named parameters are provided, they will be merged together. If there are conflicts during the merge, the named parameters will take precedence and override the top level json keys.

Currently the named parameters that `DatabricksSubmitRunOperator` supports are

- `spark_jar_task`
- `notebook_task`
- `new_cluster`
- `existing_cluster_id`
- `libraries`
- `run_name`
- `timeout_seconds`

Parameters

- **json** (*dict*) – A JSON object containing API parameters which will be passed directly to the `api/2.0/jobs/runs/submit` endpoint. The other named parameters (i.e. `spark_jar_task`, `notebook_task`..) to this operator will be merged with this json dictionary if they are provided. If there are conflicts during the merge, the named parameters will take precedence and override the top level json keys. This field will be templated.

See also:

For more information about templating see *Jinja Templating*. <https://docs.databricks.com/api/latest/jobs.html#runs-submit>

- **spark_jar_task** (*dict*) – The main class and parameters for the JAR task. Note that the actual JAR is specified in the `libraries`. *EITHER* `spark_jar_task` *OR* `notebook_task` should be specified. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/jobs.html#jobssparkjartask>

- **notebook_task** (*dict*) – The notebook path and parameters for the notebook task. *EITHER* `spark_jar_task` *OR* `notebook_task` should be specified. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/jobs.html#jobsnotebooktask>

- **new_cluster** (*dict*) – Specs for a new cluster on which this task will be run. *EITHER* `new_cluster` *OR* `existing_cluster_id` should be specified. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/jobs.html#jobsclusterspecnewcluster>

- **existing_cluster_id** (*string*) – ID for existing cluster on which to run this task. *EITHER* `new_cluster` *OR* `existing_cluster_id` should be specified. This field will be templated.

- **libraries** (*list of dicts*) – Libraries which this run will use. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/libraries.html#managedlibrarieslibrary>

- **run_name** (*string*) – The run name used for this task. By default this will be set to the Airflow `task_id`. This `task_id` is a required parameter of the superclass `BaseOperator`. This field will be templated.
- **timeout_seconds** (*int32*) – The timeout for this run. By default a value of 0 is used which means to have no timeout. This field will be templated.
- **databricks_conn_id** (*string*) – The name of the Airflow connection to use. By default and in the common case this will be `databricks_default`.
- **polling_period_seconds** (*int*) – Controls the rate which we poll for the result of this run. By default the operator will poll every 30 seconds.
- **databricks_retry_limit** (*int*) – Amount of times retry if the Databricks backend is unreachable. Its value must be greater than or equal to 1.

3.15.4 GCP: Google Cloud Platform

Airflow has extensive support for the Google Cloud Platform. But note that most Hooks and Operators are in the contrib section. Meaning that they have a *beta* status, meaning that they can have breaking changes between minor releases.

3.15.4.1 BigQuery

- *BigQueryCheckOperator* : Performs checks against a SQL query that will return a single row with different values.
- *BigQueryValueCheckOperator* : Performs a simple value check using SQL code.
- *BigQueryIntervalCheckOperator* : Checks that the values of metrics given as SQL expressions are within a certain tolerance of the ones from `days_back` before.
- *BigQueryOperator* : Executes BigQuery SQL queries in a specific BigQuery database.
- *BigQueryToBigQueryOperator* : Copy a BigQuery table to another BigQuery table.
- *BigQueryToCloudStorageOperator* : Transfers a BigQuery table to a Google Cloud Storage bucket

BigQueryCheckOperator

```
class airflow.contrib.operators.bigquery_check_operator.BigQueryCheckOperator (sql,
                                                                              big-
                                                                              query_conn_id='bigq
                                                                              *args,
                                                                              **kwargs)
```

Performs checks against Presto. The `BigQueryCheckOperator` expects a `sql` query that will return a single row. Each value on that first row is evaluated using `python bool` casting. If any of the values return `False` the check is failed and errors out.

Note that Python `bool` casting evals the following as `False`:

- `False`

- 0
- Empty string ("")
- Empty list ([])
- Empty dictionary or set ({})

Given a query like `SELECT COUNT(*) FROM foo`, it will fail only if the count == 0. You can craft much more complex query that could, for instance, check that the table has the same number of rows as the source table upstream, or that the count of today's partition is greater than yesterday's partition, or that a set of metrics are less than 3 standard deviation for the 7 day average.

This operator can be used as a data quality check in your pipeline, and depending on where you put it in your DAG, you have the choice to stop the critical path, preventing from publishing dubious data, or on the side and receive email alerts without stopping the progress of the DAG.

Parameters

- **sql** (*string*) – the sql to be executed
- **bigquery_conn_id** – reference to the BigQuery database

BigQueryValueCheckOperator

```
class airflow.contrib.operators.bigquery_check_operator.BigQueryValueCheckOperator (sql,
                                                                                       pass_value,
                                                                                       tol-
                                                                                       er-
                                                                                       ance=None,
                                                                                       big-
                                                                                       query_conn_id,
                                                                                       *args,
                                                                                       **kwargs)
```

Performs a simple value check using sql code.

Parameters **sql** (*string*) – the sql to be executed

BigQueryIntervalCheckOperator

```
class airflow.contrib.operators.bigquery_check_operator.BigQueryIntervalCheckOperator (table,
                                                                                       met-
                                                                                       rics_thres-
                                                                                       date_filter
                                                                                       days_back
                                                                                       7,
                                                                                       big-
                                                                                       query_con
                                                                                       *args,
                                                                                       **kwargs)
```

Checks that the values of metrics given as SQL expressions are within a certain tolerance of the ones from `days_back` before.

This method constructs a query like so:

```
SELECT {metrics_threshold_dict_key} FROM {table} WHERE {date_filter_column}=<date>
```

Parameters

- **table** (*str*) – the table name
- **days_back** (*int*) – number of days between ds and the ds we want to check against. Defaults to 7 days
- **metrics_threshold** (*dict*) – a dictionary of ratios indexed by metrics, for example 'COUNT(*)': 1.5 would require a 50 percent or less difference between the current day, and the prior days_back.

BigQueryOperator

```
class airflow.contrib.operators.bigquery_operator.BigQueryOperator(bql, destination_dataset_table=False, write_disposition='WRITE_EMPTY', allow_large_results=False, bigquery_conn_id='bigquery_default', delegate_to=None, udf_config=False, use_legacy_sql=True, *args, **kwargs)
```

Executes BigQuery SQL queries in a specific BigQuery database

Parameters

- **bql** (Can receive a *str* representing a *sql* statement, a list of *str* (*sql* statements), or reference to a template file. Template reference are recognized by *str* ending in *'.sql'*) – the *sql* code to be executed
- **destination_dataset_table** (*string*) – A dotted (<project>.<project>:<dataset>.<table>) that, if set, will store the results of the query.
- **bigquery_conn_id** (*string*) – reference to a specific BigQuery hook.
- **delegate_to** (*string*) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.
- **udf_config** (*list*) – The User Defined Function configuration for the query. See <https://cloud.google.com/bigquery/user-defined-functions> for details.
- **use_legacy_sql** (*boolean*) – Whether to use legacy SQL (true) or standard SQL (false).

BigQueryToBigQueryOperator

```
class airflow.contrib.operators.bigquery_to_bigquery.BigQueryToBigQueryOperator (source_project_data-
des-
ti-
na-
tion_project_data=
write_disposition=
cre-
ate_disposition='C
big-
query_conn_id='b
del-
e-
gate_to=None,
*args,
**kwargs)
```

Copies data from one BigQuery table to another. See here:

<https://cloud.google.com/bigquery/docs/reference/v2/jobs#configuration.copy>

For more details about these parameters.

Parameters

- **source_project_dataset_tables** (*list/string*) – One or more dotted (project:**|**project.)<dataset>.<table> BigQuery tables to use as the source data. If <project> is not included, project will be the project defined in the connection json. Use a list if there are multiple source tables.
- **destination_project_dataset_table** (*string*) – The destination BigQuery table. Format is: (project:**|**project.)<dataset>.<table>
- **write_disposition** (*string*) – The write disposition if the table already exists.
- **create_disposition** (*string*) – The create disposition if the table doesn't exist.
- **bigquery_conn_id** (*string*) – reference to a specific BigQuery hook.
- **delegate_to** (*string*) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.

BigQueryToCloudStorageOperator

```
class airflow.contrib.operators.bigquery_to_gcs.BigQueryToCloudStorageOperator (source_project_data
des-
ti-
na-
tion_cloud_storage_
com-
pres-
sion='NONE',
ex-
port_format='CSV',
field_delimiter=',
',
print_header=True,
big-
query_conn_id='big
del-
e-
gate_to=None,
*args,
**kwargs)
```

Transfers a BigQuery table to a Google Cloud Storage bucket.

See here:

<https://cloud.google.com/bigquery/docs/reference/v2/jobs>

For more details about these parameters.

Parameters

- **source_project_dataset_table** (*string*) – The dotted (<project>.<project>:<dataset>.<table> BigQuery table to use as the source data. If <project> is not included, project will be the project defined in the connection json.
- **destination_cloud_storage_uris** (*list*) – The destination Google Cloud Storage URI (e.g. gs://some-bucket/some-file.txt). Follows convention defined here: <https://cloud.google.com/bigquery/exporting-data-from-bigquery#exportingmultiple>
- **compression** (*string*) – Type of compression to use.
- **export_format** – File format to export.
- **field_delimiter** (*string*) – The delimiter to use when extracting to a CSV.
- **print_header** (*boolean*) – Whether to print a header for a CSV file extract.
- **bigquery_conn_id** (*string*) – reference to a specific BigQuery hook.
- **delegate_to** (*string*) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.

BigQueryHook

```
class airflow.contrib.hooks.bigquery_hook.BigQueryHook (bigquery_conn_id='bigquery_default',
delegate_to=None)
```

Interact with BigQuery. This hook uses the Google Cloud Platform connection.

get_conn()

Returns a BigQuery PEP 249 connection object.

get_pandas_df (*bql*, *parameters=None*, *dialect='legacy'*)

Returns a Pandas DataFrame for the results produced by a BigQuery query. The DbApiHook method must be overridden because Pandas doesn't support PEP 249 connections, except for SQLite. See:

<https://github.com/pydata/pandas/blob/master/pandas/io/sql.py#L447> <https://github.com/pydata/pandas/issues/6900>

Parameters

- **bql** (*string*) – The BigQuery SQL to execute.
- **parameters** (*mapping or iterable*) – The parameters to render the SQL query with (not used, leave to override superclass method)
- **dialect** (*string in {'legacy', 'standard'}, default 'legacy'*) – Dialect of BigQuery SQL – legacy SQL or standard SQL

get_service()

Returns a BigQuery service object.

insert_rows (*table*, *rows*, *target_fields=None*, *commit_every=1000*)

Insertion is currently unsupported. Theoretically, you could use BigQuery's streaming API to insert rows into a table, but this hasn't been implemented.

table_exists (*project_id*, *dataset_id*, *table_id*)

Checks for the existence of a table in Google BigQuery.

Parameters **project_id** – The Google cloud project in which to look for the table. The connection supplied to the hook

must provide access to the specified project. :type project_id: string :param dataset_id: The name of the dataset in which to look for the table.

storage bucket.

Parameters **table_id** (*string*) – The name of the table to check the existence of.

3.15.4.2 Cloud DataFlow

- *DataFlowJavaOperator* :

DataFlowJavaOperator

```
class airflow.contrib.operators.dataflow_operator.DataFlowJavaOperator(jar,
                                                                    dataflow_default_options=None,
                                                                    op-
                                                                    tions=None,
                                                                    gcp_conn_id='google_cloud_de
                                                                    dele-
                                                                    gate_to=None,
                                                                    *args,
                                                                    **kwargs)
```

Start a Java Cloud DataFlow batch job. The parameters of the operation will be passed to the job.

It's a good practice to define dataflow_* parameters in the default_args of the dag like the project, zone and staging location.

```

““ default_args = {
    'dataflow_default_options': { 'project': 'my-gcp-project', 'zone': 'europe-west1-d', 'stagingLo-
        cation': 'gs://my-staging-bucket/staging/'
    }
}

```

You need to pass the path to your dataflow as a file reference with the `jar` parameter, the jar needs to be a self executing jar. Use `options` to pass on options to your job.

```

““ t1 = DataFlowOperation(
    task_id='datapflow_example',    jar='{{ var.value.gcp_dataflow_base }}pipeline/build/libs/pipeline-
example-1.0.jar', options={
        'autoscalingAlgorithm': 'BASIC', 'maxNumWorkers': '50', 'start': '{{ds}}', 'partition-
        Type': 'DAY'
    }, dag=my-dag)
““

```

Both `jar` and `options` are templated so you can use variables in them.

```

default_args = {
    'owner': 'airflow',
    'depends_on_past': False,
    'start_date':
        (2016, 8, 1),
    'email': ['alex@vanboxel.be'],
    'email_on_failure': False,
    'email_on_retry': False,
    'retries': 1,
    'retry_delay': timedelta(minutes=30),
    'dataflow_default_options': {
        'project': 'my-gcp-project',
        'zone': 'us-centrall-f',
        'stagingLocation': 'gs://bucket/tmp/dataflow/staging/',
    }
}

dag = DAG('test-dag', default_args=default_args)

task = DataFlowJavaOperator(
    gcp_conn_id='gcp_default',
    task_id='normalize-cal',
    jar='{{var.value.gcp_dataflow_base}}pipeline-ingress-cal-normalize-1.0.jar',
    options={
        'autoscalingAlgorithm': 'BASIC',
        'maxNumWorkers': '50',
        'start': '{{ds}}',
        'partitionType': 'DAY'
    },
    dag=dag)

```

DataFlowHook

```

class airflow.contrib.hooks.gcp_dataflow_hook.DataFlowHook(gcp_conn_id='google_cloud_default',
    delegate_to=None)

```



```
get_conn()
```

Returns a Google Cloud Storage service object.

3.15.4.3 Cloud DataProc

- *DataProcPigOperator* : Start a Pig query Job on a Cloud DataProc cluster.
- *DataProcHiveOperator* : Start a Hive query Job on a Cloud DataProc cluster.
- *DataProcSparkSqlOperator* : Start a Spark SQL query Job on a Cloud DataProc cluster.
- *DataProcSparkOperator* : Start a Spark Job on a Cloud DataProc cluster.
- *DataProcHadoopOperator* : Start a Hadoop Job on a Cloud DataProc cluster.
- *DataProcPySparkOperator* : Start a PySpark Job on a Cloud DataProc cluster.

DataProcPigOperator

```
class airflow.contrib.operators.dataproc_operator.DataProcPigOperator(query=None,
                                                                       query_uri=None,
                                                                       variables=None,
                                                                       job_name='{{task.task_id}}_{{ds}}',
                                                                       dataproc_cluster='cluster-1',
                                                                       dataproc_pig_properties=None,
                                                                       dataproc_pig_jars=None,
                                                                       gcp_conn_id='google_cloud_default',
                                                                       delete_gateway=None,
                                                                       *args,
                                                                       **kwargs)
```

Start a Pig query Job on a Cloud DataProc cluster. The parameters of the operation will be passed to the cluster.

It's a good practice to define `dataproc_*` parameters in the `default_args` of the dag like the cluster name and UDFs.

```
““ default_args = {
    'dataproc_cluster': 'cluster-1', 'dataproc_pig_jars': [
        'gs://example/udf/jar/datafu/1.2.0/datafu.jar', 'gs://example/udf/jar/gpig/1.2/gpig.jar'
    ]
}
```

You can pass a pig script as string or file reference. Use variables to pass on variables for the pig script to be resolved on the cluster or use the parameters to be resolved in the script as template parameters.

```
““ t1 = DataProcPigOperator(
    task_id='dataproc_pig', query='a_pig_script.pig', variables={'out': 'gs://example/output/{{ds}}'},
    dag=dag) ““
```

DataProcHiveOperator

```
class airflow.contrib.operators.dataproc_operator.DataProcHiveOperator(query,
                                                                    vari-
                                                                    ables=None,
                                                                    job_name='{{task.task_id}}_{{a
                                                                    dataproc_cluster='cluster-
                                                                    l', dat-
                                                                    aproc_hive_properties=None,
                                                                    dat-
                                                                    aproc_hive_jars=None,
                                                                    gcp_conn_id='google_cloud_de
                                                                    dele-
                                                                    gate_to=None,
                                                                    *args,
                                                                    **kwargs)
```

Start a Hive query Job on a Cloud DataProc cluster.

DataProcSparkSqlOperator

```
class airflow.contrib.operators.dataproc_operator.DataProcSparkSqlOperator(query,
                                                                    vari-
                                                                    ables=None,
                                                                    job_name='{{task.task_id
                                                                    dataproc_cluster='cluster
                                                                    l',
                                                                    dat-
                                                                    aproc_spark_properties=N
                                                                    dat-
                                                                    aproc_spark_jars=None,
                                                                    gcp_conn_id='google_clo
                                                                    del-
                                                                    e-
                                                                    gate_to=None,
                                                                    *args,
                                                                    **kwargs)
```

Start a Spark SQL query Job on a Cloud DataProc cluster.

DataProcSparkOperator

```
class airflow.contrib.operators.dataproc_operator.DataProcSparkOperator (main_jar=None,
                                                                    main_class=None,
                                                                    arguments=None,
                                                                    archives=None,
                                                                    files=None,
                                                                    job_name='{{task.task_id}}_{{task_id}}',
                                                                    dataproc_cluster='cluster-1',
                                                                    dataproc_spark_properties=None,
                                                                    dataproc_spark_jars=None,
                                                                    gcp_conn_id='google_cloud_default',
                                                                    gate_to=None,
                                                                    *args,
                                                                    **kwargs)
```

Start a Spark Job on a Cloud DataProc cluster.

DataProcHadoopOperator

```
class airflow.contrib.operators.dataproc_operator.DataProcHadoopOperator (main_jar=None,
                                                                    main_class=None,
                                                                    arguments=None,
                                                                    archives=None,
                                                                    files=None,
                                                                    job_name='{{task.task_id}}_{{task_id}}',
                                                                    dataproc_cluster='cluster-1',
                                                                    dataproc_hadoop_properties=None,
                                                                    dataproc_hadoop_jars=None,
                                                                    gcp_conn_id='google_cloud_default',
                                                                    gate_to=None,
                                                                    *args,
                                                                    **kwargs)
```

Start a Hadoop Job on a Cloud DataProc cluster.

DataProcPySparkOperator

```
class airflow.contrib.operators.dataproc_operator.DataProcPySparkOperator (main,
                                                                           ar-
                                                                           gu-
                                                                           ments=None,
                                                                           archives=None,
                                                                           py-
                                                                           files=None,
                                                                           files=None,
                                                                           job_name='{{task.task_id}}',
                                                                           dataproc_cluster='cluster-
                                                                           1',
                                                                           dat-
                                                                           aproc_pyspark_properties=
                                                                           dat-
                                                                           aproc_pyspark_jars=None,
                                                                           gcp_conn_id='google_cloud_
                                                                           del-
                                                                           e-
                                                                           gate_to=None,
                                                                           *args,
                                                                           **kwargs)
```

Start a PySpark Job on a Cloud DataProc cluster.

3.15.4.4 Cloud Datastore

```
class airflow.contrib.hooks.datastore_hook.DatastoreHook (datastore_conn_id='google_cloud_datastore_default',
                                                           delegate_to=None)
```

Interact with Google Cloud Datastore. This hook uses the Google Cloud Platform connection.

This object is not threads safe. If you want to make multiple requests simultaneously, you will need to create a hook per thread.

allocate_ids (*partialKeys*)

Allocate IDs for incomplete keys. see <https://cloud.google.com/datastore/docs/apis/v1beta2/datasets/allocateIds>

Parameters *partialKeys* – a list of partial keys

Returns a list of full keys.

begin_transaction ()

Get a new transaction handle see <https://cloud.google.com/datastore/docs/apis/v1beta2/datasets/beginTransaction>

Returns a transaction handle

commit (*body*)

Commit a transaction, optionally creating, deleting or modifying some entities. see <https://cloud.google.com/datastore/docs/apis/v1beta2/datasets/commit>

Parameters *body* – the body of the commit request

Returns the response body of the commit request

get_conn ()

Returns a Google Cloud Storage service object.

lookup (*keys, read_consistency=None, transaction=None*)

Lookup some entities by key see <https://cloud.google.com/datastore/docs/apis/v1beta2/datasets/lookup>
:param keys: the keys to lookup :param read_consistency: the read consistency to use. default, strong or eventual.

Cannot be used with a transaction.

Parameters **transaction** – the transaction to use, if any.

Returns the response body of the lookup request.

rollback (*transaction*)

Roll back a transaction see <https://cloud.google.com/datastore/docs/apis/v1beta2/datasets/rollback> :param transaction: the transaction to roll back

run_query (*body*)

Run a query for entities. see <https://cloud.google.com/datastore/docs/apis/v1beta2/datasets/runQuery>
:param body: the body of the query request :return: the batch of query results.

3.15.4.5 Cloud Storage

- *GoogleCloudStorageDownloadOperator* : Downloads a file from Google Cloud Storage.
- *GoogleCloudStorageToBigQueryOperator* : Loads files from Google cloud storage into BigQuery.

GoogleCloudStorageDownloadOperator

```
class airflow.contrib.operators.gcs_download_operator.GoogleCloudStorageDownloadOperator (bucket, object, filename, store_to_xcom_key, google_credentials, delete_on_completion, *args, **kwargs)
```

Downloads a file from Google Cloud Storage.

Parameters

- **bucket** (*string*) – The Google cloud storage bucket where the object is.
- **object** (*string*) – The name of the object to download in the Google cloud storage bucket.
- **filename** (*string*) – The file path on the local file system (where the operator is being executed) that the file should be downloaded to. If false, the downloaded data will not be stored on the local file system.
- **store_to_xcom_key** (*string*) – If this param is set, the operator will push the contents of the downloaded file to XCom with the key set in this parameter. If false, the downloaded data will not be pushed to XCom.

- **google_cloud_storage_conn_id**(*string*) – The connection ID to use when connecting to Google cloud storage.
- **delegate_to**(*string*) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.

GoogleCloudStorageToBigQueryOperator

```
class airflow.contrib.operators.gcs_to_bq.GoogleCloudStorageToBigQueryOperator (bucket,
                                         source_objects,
                                         dest_project_dataset,
                                         schema_fields=None,
                                         schema_object=None,
                                         source_format='CSV',
                                         create_disposition='CREATE_IF_EXISTS',
                                         skip_leading_rows=0,
                                         write_disposition='WRITE_TRUNCATE',
                                         field_delimiter=',',
                                         max_bad_records=0,
                                         max_id_key=None,
                                         big_query_conn_id='big_query_default',
                                         google_cloud_storage_conn_id=None,
                                         delegate_to=None,
                                         schema_update_option=None,
                                         *args,
                                         **kwargs)
```

Loads files from Google cloud storage into BigQuery.

GoogleCloudStorageHook

```
class airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook (google_cloud_storage_conn_id='google_cloud_storage_default',
                                                              delegate_to=None)
```

Interact with Google Cloud Storage. This hook uses the Google Cloud Platform connection.

delete (*bucket, object, generation=None*)

Delete an object if versioning is not enabled for the bucket, or if generation parameter is used. :param bucket: name of the bucket, where the object resides :type bucket: string :param object: name of the object to delete :type object: string :param generation: if present, permanently delete the object of this generation :type generation: string :return: True if succeeded

download (*bucket, object, filename=False*)

Get a file from Google Cloud Storage.

Parameters

- **bucket** (*string*) – The bucket to fetch from.

- **object** (*string*) – The object to fetch.
- **filename** (*string*) – If set, a local file path where the file should be written to.

exists (*bucket, object*)

Checks for the existence of a file in Google Cloud Storage.

Parameters

- **bucket** (*string*) – The Google cloud storage bucket where the object is.
- **object** (*string*) – The name of the object to check in the Google cloud storage bucket.

get_conn ()

Returns a Google Cloud Storage service object.

is_updated_after (*bucket, object, ts*)

Checks if an object is updated in Google Cloud Storage.

Parameters

- **bucket** (*string*) – The Google cloud storage bucket where the object is.
- **object** (*string*) – The name of the object to check in the Google cloud storage bucket.
- **ts** (*datetime*) – The timestamp to check against.

list (*bucket, versions=None, maxResults=None, prefix=None*)

List all objects from the bucket with the give string prefix in name :param bucket: bucket name :type bucket: string :param versions: if true, list all versions of the objects :type versions: boolean :param maxResults: max count of items to return in a single page of responses :type maxResults: integer :param prefix: prefix string which filters objects whose name begin with this prefix :type prefix: string :return: a stream of object names matching the filtering criteria

upload (*bucket, object, filename, mime_type='application/octet-stream'*)

Uploads a local file to Google Cloud Storage.

Parameters

- **bucket** (*string*) – The bucket to upload to.
- **object** (*string*) – The object name to set when uploading the local file.
- **filename** (*string*) – The local file path to the file to be uploaded.
- **mime_type** (*string*) – The MIME type to set when uploading the file.

3.16 FAQ

3.16.1 Why isn't my task getting scheduled?

There are very many reasons why your task might not be getting scheduled. Here are some of the common causes:

- Does your script “compile”, can the Airflow engine parse it and find your DAG object. To test this, you can run `airflow list_dags` and confirm that your DAG shows up in the list. You can also run `airflow list_tasks foo_dag_id --tree` and confirm that your task shows up in the list as expected. If you use the CeleryExecutor, you may way to confirm that this works both where the scheduler runs as well as where the worker runs.
- Is your `start_date` set properly? The Airflow scheduler triggers the task soon after the `start_date + scheduler_interval` is passed.

- Is your `schedule_interval` set properly? The default `schedule_interval` is one day (`datetime.timedelta(1)`). You must specify a different `schedule_interval` directly to the DAG object you instantiate, not as a `default_param`, as task instances do not override their parent DAG's `schedule_interval`.
- Is your `start_date` beyond where you can see it in the UI? If you set your it to some time say 3 months ago, you won't be able to see it in the main view in the UI, but you should be able to see it in the Menu -> Browse -> Task Instances.
- Are the dependencies for the task met. The task instances directly upstream from the task need to be in a success state. Also, if you have set `depends_on_past=True`, the previous task instance needs to have succeeded (except if it is the first run for that task). Also, if `wait_for_downstream=True`, make sure you understand what it means. You can view how these properties are set from the Task Instance Details page for your task.
- Are the DagRuns you need created and active? A DagRun represents a specific execution of an entire DAG and has a state (running, success, failed, ...). The scheduler creates new DagRun as it moves forward, but never goes back in time to create new ones. The scheduler only evaluates running DagRuns to see what task instances it can trigger. Note that clearing tasks instances (from the UI or CLI) does set the state of a DagRun back to running. You can bulk view the list of DagRuns and alter states by clicking on the schedule tag for a DAG.
- Is the `concurrency` parameter of your DAG reached? `concurrency` defines how many running task instances a DAG is allowed to have, beyond which point things get queued.
- Is the `max_active_runs` parameter of your DAG reached? `max_active_runs` defines how many running concurrent instances of a DAG there are allowed to be.

You may also want to read the Scheduler section of the docs and make sure you fully understand how it proceeds.

3.16.2 How do I trigger tasks based on another task's failure?

Check out the `Trigger Rule` section in the Concepts section of the documentation

3.16.3 Why are connection passwords still not encrypted in the metadata db after I installed airflow[crypto]?

Check out the `Connections` section in the Configuration section of the documentation

3.16.4 What's the deal with `start_date`?

`start_date` is partly legacy from the pre-DagRun era, but it is still relevant in many ways. When creating a new DAG, you probably want to set a global `start_date` for your tasks using `default_args`. The first DagRun to be created will be based on the `min(start_date)` for all your task. From that point on, the scheduler creates new DagRuns based on your `schedule_interval` and the corresponding task instances run as your dependencies are met. When introducing new tasks to your DAG, you need to pay special attention to `start_date`, and may want to reactivate inactive DagRuns to get the new task to get onboarded properly.

We recommend against using dynamic values as `start_date`, especially `datetime.now()` as it can be quite confusing. The task is triggered once the period closes, and in theory an `@hourly` DAG would never get to an hour after now as `now()` moves along.

Previously we also recommended using rounded `start_date` in relation to your `schedule_interval`. This meant an `@hourly` would be at 00:00 minutes:seconds, a `@daily` job at midnight, a `@monthly` job on the first of the month. This is no longer required. Airflow will now auto align the `start_date` and the `schedule_interval`, by using the `start_date` as the moment to start looking.

You can use any sensor or a `TimeDeltaSensor` to delay the execution of tasks within the schedule interval. While `schedule_interval` does allow specifying a `datetime.timedelta` object, we recommend using the macros or cron expressions instead, as it enforces this idea of rounded schedules.

When using `depends_on_past=True` it's important to pay special attention to `start_date` as the past dependency is not enforced only on the specific schedule of the `start_date` specified for the task. It's also important to watch `DagRun` activity status in time when introducing new `depends_on_past=True`, unless you are planning on running a backfill for the new task(s).

Also important to note is that the tasks `start_date`, in the context of a backfill CLI command, get overridden by the backfill's command `start_date`. This allows for a backfill on tasks that have `depends_on_past=True` to actually start, if it wasn't the case, the backfill just wouldn't start.

3.16.5 How can I create DAGs dynamically?

Airflow looks in your `DAGS_FOLDER` for modules that contain DAG objects in their global namespace, and adds the objects it finds in the `DagBag`. Knowing this all we need is a way to dynamically assign variable in the global namespace, which is easily done in python using the `globals()` function for the standard library which behaves like a simple dictionary.

```
for i in range(10):
    dag_id = 'foo_{}'.format(i)
    globals()[dag_id] = DAG(dag_id)
    # or better, call a function that returns a DAG object!
```

3.16.6 What are all the airflow run commands in my process list?

There are many layers of `airflow run` commands, meaning it can call itself.

- Basic `airflow run`: fires up an executor, and tell it to run an `airflow run --local` command. if using Celery, this means it puts a command in the queue for it to run remote, on the worker. If using `LocalExecutor`, that translates into running it in a subprocess pool.
- Local `airflow run --local`: starts an `airflow run --raw` command (described below) as a subprocess and is in charge of emitting heartbeats, listening for external kill signals and ensures some cleanup takes place if the subprocess fails
- Raw `airflow run --raw` runs the actual operator's `execute` method and performs the actual work

3.17 API Reference

3.17.1 Operators

Operators allow for generation of certain types of tasks that become nodes in the DAG when instantiated. All operators derive from `BaseOperator` and inherit many attributes and methods that way. Refer to the `BaseOperator` documentation for more details.

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, or a specific time of the day. Sensors are derived from `BaseSensorOperator` and run a `poke` method at a specified `poke_interval` until it returns `True`.

3.17.1.1 BaseOperator

All operators are derived from `BaseOperator` and acquire much functionality through inheritance. Since this is the core of the engine, it's worth taking the time to understand the parameters of `BaseOperator` to understand the primitive features that can be leveraged in your DAGs.

```
class airflow.models.BaseOperator(task_id, owner='Airflow', email=None,
                                  email_on_retry=True, email_on_failure=True, re-
                                  tries=0, retry_delay=datetime.timedelta(0, 300),
                                  retry_exponential_backoff=False, max_retry_delay=None,
                                  start_date=None, end_date=None, schedule_interval=None,
                                  depends_on_past=False, wait_for_downstream=False,
                                  dag=None, params=None, default_args=None, adhoc=False,
                                  priority_weight=1, queue='default', pool=None, sla=None,
                                  execution_timeout=None, on_failure_callback=None,
                                  on_success_callback=None, on_retry_callback=None, trig-
                                  ger_rule='all_success', resources=None, run_as_user=None,
                                  *args, **kwargs)
```

Abstract base class for all operators. Since operators create objects that become node in the dag, `BaseOperator` contains many recursive methods for dag crawling behavior. To derive this class, you are expected to override the constructor as well as the 'execute' method.

Operators derived from this class should perform or trigger certain tasks synchronously (wait for completion). Example of operators could be an operator the runs a Pig job (`PigOperator`), a sensor operator that waits for a partition to land in Hive (`HiveSensorOperator`), or one that moves data from Hive to MySQL (`Hive2MySQLOperator`). Instances of these operators (tasks) target specific operations, running specific scripts, functions or data transfers.

This class is abstract and shouldn't be instantiated. Instantiating a class derived from this one results in the creation of a task object, which ultimately becomes a node in DAG objects. Task dependencies should be set by using the `set_upstream` and/or `set_downstream` methods.

Note that this class is derived from SQLAlchemy's Base class, which allows us to push metadata regarding tasks to the database. Deriving this classes needs to implement the polymorphic specificities documented in SQLAlchemy. This should become clear while reading the code for other operators.

Parameters

- **task_id** (*string*) – a unique, meaningful id for the task
- **owner** (*string*) – the owner of the task, using the unix username is recommended
- **retries** (*int*) – the number of retries that should be performed before failing the task
- **retry_delay** (*timedelta*) – delay between retries
- **retry_exponential_backoff** (*bool*) – allow progressive longer waits between retries by using exponential backoff algorithm on retry delay (delay will be converted into seconds)
- **max_retry_delay** (*timedelta*) – maximum delay interval between retries
- **start_date** (*datetime*) – The `start_date` for the task, determines the `execution_date` for the first task instance. The best practice is to have the `start_date`

rounded to your DAG's `schedule_interval`. Daily jobs have their `start_date` some day at 00:00:00, hourly jobs have their `start_date` at 00:00 of a specific hour. Note that Airflow simply looks at the latest `execution_date` and adds the `schedule_interval` to determine the next `execution_date`. It is also very important to note that different tasks' dependencies need to line up in time. If task A depends on task B and their `start_date` are offset in a way that their `execution_date` don't line up, A's dependencies will never be met. If you are looking to delay a task, for example running a daily task at 2AM, look into the `TimeSensor` and `TimeDeltaSensor`. We advise against using dynamic `start_date` and recommend using fixed ones. Read the FAQ entry about `start_date` for more information.

- **end_date** (*datetime*) – if specified, the scheduler won't go beyond this date
- **depends_on_past** (*bool*) – when set to true, task instances will run sequentially while relying on the previous task's schedule to succeed. The task instance for the `start_date` is allowed to run.
- **wait_for_downstream** (*bool*) – when set to true, an instance of task X will wait for tasks immediately downstream of the previous instance of task X to finish successfully before it runs. This is useful if the different instances of a task X alter the same asset, and this asset is used by tasks downstream of task X. Note that `depends_on_past` is forced to True wherever `wait_for_downstream` is used.
- **queue** (*str*) – which queue to target when running this job. Not all executors implement queue management, the CeleryExecutor does support targeting specific queues.
- **dag** (*DAG*) – a reference to the dag the task is attached to (if any)
- **priority_weight** (*int*) – priority weight of this task against other task. This allows the executor to trigger higher priority tasks before others when things get backed up.
- **pool** (*str*) – the slot pool this task should run in, slot pools are a way to limit concurrency for certain tasks
- **sla** (*datetime.timedelta*) – time by which the job is expected to succeed. Note that this represents the `timedelta` after the period is closed. For example if you set an SLA of 1 hour, the scheduler would send an email soon after 1:00AM on the 2016-01-02 if the 2016-01-01 instance has not succeeded yet. The scheduler pays special attention for jobs with an SLA and sends alert emails for SLA misses. SLA misses are also recorded in the database for future reference. All tasks that share the same SLA time get bundled in a single email, sent soon after that time. SLA notification are sent once and only once for each task instance.
- **execution_timeout** (*datetime.timedelta*) – max time allowed for the execution of this task instance, if it goes beyond it will raise and fail.
- **on_failure_callback** (*callable*) – a function to be called when a task instance of this task fails. a context dictionary is passed as a single parameter to this function. Context contains references to related objects to the task instance and is documented under the macros section of the API.
- **on_retry_callback** – much like the `on_failure_callback` except that it is executed when retries occur.
- **on_success_callback** (*callable*) – much like the `on_failure_callback` except that it is executed when the task succeeds.
- **trigger_rule** (*str*) – defines the rule by which dependencies are applied for the task to get triggered. Options are: { `all_success` | `all_failed` | `all_done` | `one_success` | `one_failed` | `dummy` } default is `all_success`. Options

can be set as string or using the constants defined in the static class `airflow.utils.TriggerRule`

- **resources** (*dict*) – A map of resource parameter names (the argument names of the Resources constructor) to their values.
- **run_as_user** (*str*) – unix username to impersonate while running the task

3.17.1.2 BaseSensorOperator

All sensors are derived from `BaseSensorOperator`. All sensors inherit the `timeout` and `poke_interval` on top of the `BaseOperator` attributes.

```
class airflow.operators.sensors.BaseSensorOperator (poke_interval=60, timeout=604800,
                                                    soft_fail=False, *args, **kwargs)
```

Sensor operators are derived from this class and inherit these attributes.

Sensor operators keep executing at a time interval and succeed when a criteria is met and fail if and when they time out.

Parameters

- **soft_fail** (*bool*) – Set to true to mark the task as SKIPPED on failure
- **poke_interval** (*int*) – Time in seconds that the job should wait in between each tries
- **timeout** (*int*) – Time, in seconds before the task times out and fails.

3.17.1.3 Operator API

Importer that dynamically loads a class and module from its parent. This allows Airflow to support `from airflow.operators import BashOperator` even though `BashOperator` is actually in `airflow.operators.bash_operator`.

The importer also takes over for the parent_module by wrapping it. This is required to support attribute-based usage:

```
from airflow import operators
operators.BashOperator(...)
```

```
class airflow.operators.BashOperator (bash_command, xcom_push=False, env=None,
                                       output_encoding='utf-8', *args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Execute a Bash script, command or set of commands.

Parameters

- **bash_command** (*string*) – The command, set of commands or reference to a bash script (must be `‘.sh’`) to be executed.
- **xcom_push** (*bool*) – If `xcom_push` is True, the last line written to stdout will also be pushed to an XCom when the bash command completes.
- **env** (*dict*) – If `env` is not None, it must be a mapping that defines the environment variables for the new process; these are used instead of inheriting the current process environment, which is the default behavior. (templated)

execute (*context*)

Execute the bash command in a temporary directory which will be cleaned afterwards

```
class airflow.operators.BranchPythonOperator (python_callable,          op_args=None,
                                              op_kwargs=None,    provide_context=False,
                                              templates_dict=None, templates_exts=None,
                                              *args, **kwargs)
```

Bases: `python_operator.PythonOperator`

Allows a workflow to “branch” or follow a single path following the execution of this task.

It derives the `PythonOperator` and expects a Python function that returns the `task_id` to follow. The `task_id` returned should point to a task directly downstream from `{self}`. All other “branches” or directly downstream tasks are marked with a state of `skipped` so that these paths can’t move forward. The `skipped` states are propagated downstream to allow for the DAG state to fill up and the DAG run’s state to be inferred.

Note that using tasks with `depends_on_past=True` downstream from `BranchPythonOperator` is logically unsound as `skipped` status will invariably lead to block tasks that depend on their past successes. `skipped` states propagate where all directly upstream tasks are `skipped`.

```
class airflow.operators.TriggerDagRunOperator (trigger_dag_id,  python_callable,  *args,
                                              **kwargs)
```

Bases: `airflow.models.BaseOperator`

Triggers a DAG run for a specified `dag_id` if a criteria is met

Parameters

- **trigger_dag_id** (*str*) – the `dag_id` to trigger
- **python_callable** (*python callable*) – a reference to a python function that will be called while passing it the `context` object and a placeholder object `obj` for your callable to fill and return if you want a `DagRun` created. This `obj` object contains a `run_id` and `payload` attribute that you can modify in your function. The `run_id` should be a unique identifier for that DAG run, and the `payload` has to be a picklable object that will be made available to your tasks while executing that DAG run. Your function header should look like `def foo(context, dag_run_obj):`

```
class airflow.operators.DummyOperator (*args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Operator that does literally nothing. It can be used to group tasks in a DAG.

```
class airflow.operators.EmailOperator (to,  subject,  html_content,  files=None,  cc=None,
                                       bcc=None, mime_subtype='mixed', *args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Sends an email.

Parameters

- **to** (*list or string (comma or semicolon delimited)*) – list of emails to send the email to
- **subject** (*string*) – subject line for the email (templated)
- **html_content** (*string*) – content of the email (templated), html markup is allowed
- **files** (*list*) – file names to attach in email
- **cc** (*list or string (comma or semicolon delimited)*) – list of recipients to be added in CC field
- **bcc** (*list or string (comma or semicolon delimited)*) – list of recipients to be added in BCC field

```
class airflow.operators.ExternalTaskSensor (external_dag_id, external_task_id, allowed_states=None, execution_delta=None, execution_date_fn=None, *args, **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Waits for a task to complete in a different DAG

Parameters

- **external_dag_id** (*string*) – The dag_id that contains the task you want to wait for
- **external_task_id** (*string*) – The task_id that contains the task you want to wait for
- **allowed_states** (*list*) – list of allowed states, default is ['success']
- **execution_delta** (*datetime.timedelta*) – time difference with the previous execution to look at, the default is the same execution_date as the current task. For yesterday, use [positive!] `datetime.timedelta(days=1)`. Either `execution_delta` or `execution_date_fn` can be passed to `ExternalTaskSensor`, but not both.
- **execution_date_fn** (*callable*) – function that receives the current execution date and returns the desired execution date to query. Either `execution_delta` or `execution_date_fn` can be passed to `ExternalTaskSensor`, but not both.

```
class airflow.operators.GenericTransfer (sql, destination_table, source_conn_id, destination_conn_id, preoperator=None, *args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Moves data from a connection to another, assuming that they both provide the required methods in their respective hooks. The source hook needs to expose a `get_records` method, and the destination a `insert_rows` method.

This is meant to be used on small-ish datasets that fit in memory.

Parameters

- **sql** (*str*) – SQL query to execute against the source database
- **destination_table** (*str*) – target table
- **source_conn_id** (*str*) – source connection
- **destination_conn_id** (*str*) – source connection
- **preoperator** (*str or list of str*) – sql statement or list of statements to be executed prior to loading the data

```
class airflow.operators.HdfsSensor (filepath, hdfs_conn_id='hdfs_default', ignored_ext=['_COPYING_'], ignore_copying=True, file_size=None, hook=<class 'airflow.hooks.hdfs_hook.HDFSHook'>, *args, **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Waits for a file or folder to land in HDFS

static filter_for_filesize (*result, size=None*)

Will test the filepath result and test if its size is at least self.filesize :param result: a list of dicts returned by Snakebite ls :param size: the file size in MB a file should be at least to trigger True :return: (bool) depending on the matching criteria

static filter_for_ignored_ext (*result, ignored_ext, ignore_copying*)

Will filter if instructed to do so the result to remove matching criteria :param result: (list) of dicts returned by Snakebite ls :param ignored_ext: (list) of ignored extensions :param ignore_copying: (bool) shall we ignore ? :return:

```
class airflow.operators.HivePartitionSensor(table, partition="ds={{ ds }}", metastore_conn_id='metastore_default',
                                             schema='default', poke_interval=180, *args,
                                             **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Waits for a partition to show up in Hive.

Note: Because `partition` supports general logical operators, it can be inefficient. Consider using `NamedHivePartitionSensor` instead if you don't need the full flexibility of `HivePartitionSensor`.

Parameters

- **table** (*string*) – The name of the table to wait for, supports the dot notation (`my_database.my_table`)
- **partition** (*string*) – The partition clause to wait for. This is passed as is to the metastore Thrift client `get_partitions_by_filter` method, and apparently supports SQL like notation as in `ds='2015-01-01' AND type='value'` and comparison operators as in `"ds>=2015-01-01"`
- **metastore_conn_id** (*str*) – reference to the metastore thrift service connection id

```
class airflow.operators.SimpleHttpOperator(endpoint, method='POST', data=None,
                                             headers=None, response_check=None,
                                             extra_options=None, xcom_push=False,
                                             http_conn_id='http_default', *args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Calls an endpoint on an HTTP system to execute an action

Parameters

- **http_conn_id** (*string*) – The connection to run the sensor against
- **endpoint** (*string*) – The relative part of the full url
- **method** (*string*) – The HTTP method to use, default = "POST"
- **data** (*For POST/PUT, depends on the content-type parameter, for GET a dictionary of key/value string pairs*) – The data to pass. POST-data in POST/PUT and params in the URL for a GET request.
- **headers** (*a dictionary of string key/value pairs*) – The HTTP headers to be added to the GET request
- **response_check** (*A lambda or defined function.*) – A check against the 'requests' response object. Returns True for 'pass' and False otherwise.
- **extra_options** (*A dictionary of options, where key is string and value depends on the option that's being modified.*) – Extra options for the 'requests' library, see the 'requests' documentation (options to modify timeout, ssl, etc.)

```
class airflow.operators.HttpSensor(endpoint, http_conn_id='http_default', method='GET',
                                    params=None, headers=None, response_check=None,
                                    extra_options=None, *args, **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Executes a HTTP get statement and returns False on failure: 404 not found or `response_check` function returned False

Parameters

- **http_conn_id** (*string*) – The connection to run the sensor against
- **method** (*string*) – The HTTP request method to use
- **endpoint** (*string*) – The relative part of the full url
- **params** (*a dictionary of string key/value pairs*) – The parameters to be added to the GET url
- **headers** (*a dictionary of string key/value pairs*) – The HTTP headers to be added to the GET request
- **response_check** (*A lambda or defined function.*) – A check against the ‘requests’ response object. Returns True for ‘pass’ and False otherwise.
- **extra_options** (*A dictionary of options, where key is string and value depends on the option that’s being modified.*) – Extra options for the ‘requests’ library, see the ‘requests’ documentation (options to modify timeout, ssl, etc.)

```
class airflow.operators.MetastorePartitionSensor (table, partition_name, schema='default',
                                                  mysql_conn_id='metastore_mysql',
                                                  *args, **kwargs)
```

Bases: `sensors.SqlSensor`

An alternative to the `HivePartitionSensor` that talk directly to the MySQL db. This was created as a result of observing sub optimal queries generated by the Metastore thrift service when hitting subpartitioned tables. The Thrift service’s queries were written in a way that wouldn’t leverage the indexes.

Parameters

- **schema** (*str*) – the schema
- **table** (*str*) – the table
- **partition_name** (*str*) – the partition name, as defined in the PARTITIONS table of the Metastore. Order of the fields does matter. Examples: `ds=2016-01-01` or `ds=2016-01-01/sub=foo` for a sub partitioned table
- **mysql_conn_id** (*str*) – a reference to the MySQL conn_id for the metastore

```
class airflow.operators.NamedHivePartitionSensor (partition_names,
                                                  metastore_conn_id='metastore_default',
                                                  poke_interval=180, *args, **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Waits for a set of partitions to show up in Hive.

Parameters

- **partition_names** (*list of strings*) – List of fully qualified names of the partitions to wait for. A fully qualified name is of the form `schema.table/pk1=pv1/pk2=pv2`, for example, `default.users/ds=2016-01-01`. This is passed as is to the metastore Thrift client `get_partitions_by_name` method. Note that you cannot use logical or comparison operators as in `HivePartitionSensor`.
- **metastore_conn_id** (*str*) – reference to the metastore thrift service connection id

```
class airflow.operators.PythonOperator (python_callable, op_args=None, op_kwargs=None,
                                       provide_context=False, templates_dict=None, templates_exts=None, *args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Executes a Python callable

Parameters

- **python_callable** (*python callable*) – A reference to an object that is callable
- **op_kwargs** (*dict*) – a dictionary of keyword arguments that will get unpacked in your function
- **op_args** (*list*) – a list of positional arguments that will get unpacked when calling your callable
- **provide_context** (*bool*) – if set to true, Airflow will pass a set of keyword arguments that can be used in your function. This set of kwargs correspond exactly to what you can use in your jinja templates. For this to work, you need to define ***kwargs* in your function header.
- **templates_dict** (*dict of str*) – a dictionary where the values are templates that will get templated by the Airflow engine sometime between `__init__` and `execute` takes place and are made available in your callable's context after the template has been applied
- **templates_exts** – a list of file extensions to resolve while processing templated fields, for examples `['.sql', '.hql']`

```
class airflow.operators.S3KeySensor(bucket_key, bucket_name=None, wildcard_match=False,
                                    s3_conn_id='s3_default', *args, **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Waits for a key (a file-like instance on S3) to be present in a S3 bucket. S3 being a key/value it does not support folders. The path is just a key a resource.

Parameters

- **bucket_key** (*str*) – The key being waited on. Supports full `s3://` style url or relative path from root level.
- **bucket_name** (*str*) – Name of the S3 bucket
- **wildcard_match** (*bool*) – whether the bucket_key should be interpreted as a Unix wildcard pattern
- **s3_conn_id** (*str*) – a reference to the s3 connection

```
class airflow.operators.ShortCircuitOperator(python_callable, op_args=None,
                                              op_kwargs=None, provide_context=False,
                                              templates_dict=None, templates_exts=None,
                                              *args, **kwargs)
```

Bases: `python_operator.PythonOperator`

Allows a workflow to continue only if a condition is met. Otherwise, the workflow “short-circuits” and downstream tasks are skipped.

The ShortCircuitOperator is derived from the PythonOperator. It evaluates a condition and short-circuits the workflow if the condition is False. Any downstream tasks are marked with a state of “skipped”. If the condition is True, downstream tasks proceed as normal.

The condition is determined by the result of *python_callable*.

```
class airflow.operators.SqlSensor(conn_id, sql, *args, **kwargs)
```

Bases: `sensors.BaseSensorOperator`

Runs a sql statement until a criteria is met. It will keep trying until sql returns no row, or if the first cell in (0, '0', '').

Parameters

- **conn_id** (*string*) – The connection to run the sensor against
- **sql** – The sql to run. To pass, it needs to return at least one cell that contains a non-zero / empty string value.

class airflow.operators.**TimeSensor** (*target_time*, *args, **kwargs)

Bases: *sensors.BaseSensorOperator*

Waits until the specified time of the day.

Parameters **target_time** (*datetime.time*) – time after which the job succeeds

class airflow.operators.**WebHdfsSensor** (*filepath*, *webhdfs_conn_id*=*'webhdfs_default'*, *args, **kwargs)

Bases: *sensors.BaseSensorOperator*

Waits for a file or folder to land in HDFS

class airflow.operators.docker_operator.**DockerOperator** (*image*, *api_version*=None, *command*=None, *cpus*=1.0, *docker_url*=*'unix://var/run/docker.sock'*, *environment*=None, *force_pull*=False, *mem_limit*=None, *network_mode*=None, *tls_ca_cert*=None, *tls_client_cert*=None, *tls_client_key*=None, *tls_hostname*=None, *tls_ssl_version*=None, *tmp_dir*=*'/tmp/airflow'*, *user*=None, *volumes*=None, *xcom_push*=False, *xcom_all*=False, *args, **kwargs)

Execute a command inside a docker container.

A temporary directory is created on the host and mounted into a container to allow storing files that together exceed the default disk size of 10GB in a container. The path to the mounted directory can be accessed via the environment variable `AIRFLOW_TMP_DIR`.

Parameters

- **image** (*str*) – Docker image from which to create the container.
- **api_version** (*str*) – Remote API version.
- **command** (*str* or *list*) – Command to be run in the container.
- **cpus** (*float*) – Number of CPUs to assign to the container. This value gets multiplied with 1024. See <https://docs.docker.com/engine/reference/run/#cpu-share-constraint>
- **docker_url** (*str*) – URL of the host running the docker daemon.
- **environment** (*dict*) – Environment variables to set in the container.
- **force_pull** (*bool*) – Pull the docker image on every run.
- **mem_limit** (*float* or *str*) – Maximum amount of memory the container can use. Either a float value, which represents the limit in bytes, or a string like 128m or 1g.
- **network_mode** (*str*) – Network mode for the container.

- **tls_ca_cert** (*str*) – Path to a PEM-encoded certificate authority to secure the docker connection.
- **tls_client_cert** (*str*) – Path to the PEM-encoded certificate used to authenticate docker client.
- **tls_client_key** (*str*) – Path to the PEM-encoded key used to authenticate docker client.
- **tls_hostname** (*str or bool*) – Hostname to match against the docker server certificate or False to disable the check.
- **tls_ssl_version** (*str*) – Version of SSL to use when communicating with docker daemon.
- **tmp_dir** (*str*) – Mount point inside the container to a temporary directory created on the host by the operator. The path is also made available via the environment variable `AIRFLOW_TMP_DIR` inside the container.
- **user** (*int or str*) – Default user inside the docker container.
- **volumes** – List of volumes to mount into the container, e.g. `['/host/path:/container/path', '/host/path2:/container/path2:ro']`.
- **xcom_push** (*bool*) – Does the stdout will be pushed to the next step using XCom. The default is False.
- **xcom_all** (*bool*) – Push all the stdout or just the last line. The default is False (last line).

3.17.1.4 Community-contributed Operators

Importer that dynamically loads a class and module from its parent. This allows Airflow to support `from airflow.operators import BashOperator` even though `BashOperator` is actually in `airflow.operators.bash_operator`.

The importer also takes over for the `parent_module` by wrapping it. This is required to support attribute-based usage:

```
from airflow import operators
operators.BashOperator(...)
```

```
class airflow.contrib.operators.SSHExecuteOperator(ss_hook,          bash_command,
                                                    xcom_push=False,    env=None,
                                                    *args, **kwargs)
```

Bases: `airflow.models.BaseOperator`

Execute a Bash script, command or set of commands at remote host.

Parameters

- **ssh_hook** (*string*) – A SSHHook that indicates the remote host you want to run the script
- **bash_command** (*string*) – The command, set of commands or reference to a bash script (must be `.sh`) to be executed.
- **env** (*dict*) – If `env` is not `None`, it must be a mapping that defines the environment variables for the new process; these are used instead of inheriting the current process environment, which is the default behavior.

```
class airflow.contrib.operators.bigquery_operator.BigQueryOperator (bql, destination_dataset_table=False,
                                                                    write_disposition='WRITE_EMPTY',
                                                                    allow_large_results=False,
                                                                    bigquery_conn_id='bigquery_default',
                                                                    delegate_to=None,
                                                                    udf_config=False,
                                                                    use_legacy_sql=True,
                                                                    *args,
                                                                    **kwargs)
```

Executes BigQuery SQL queries in a specific BigQuery database

Parameters

- **bql** (Can receive a str representing a sql statement, a list of str (sql statements), or reference to a template file. Template reference are recognized by str ending in '.sql') – the sql code to be executed
- **destination_dataset_table** (string) – A dotted (<project>.<project>:<dataset>.<table> that, if set, will store the results of the query.
- **bigquery_conn_id** (string) – reference to a specific BigQuery hook.
- **delegate_to** (string) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.
- **udf_config** (list) – The User Defined Function configuration for the query. See <https://cloud.google.com/bigquery/user-defined-functions> for details.
- **use_legacy_sql** (boolean) – Whether to use legacy SQL (true) or standard SQL (false).

```
class airflow.contrib.operators.bigquery_to_gcs.BigQueryToCloudStorageOperator (source_project_data,
                                                                    destination_cloud_storage_com-
                                                                    pression='NONE',
                                                                    export_format='CSV',
                                                                    field_delimiter=',',
                                                                    print_header=True,
                                                                    bigquery_conn_id='big-
                                                                    del-
                                                                    e-
                                                                    gate_to=None,
                                                                    *args,
                                                                    **kwargs)
```

Transfers a BigQuery table to a Google Cloud Storage bucket.

See here:

<https://cloud.google.com/bigquery/docs/reference/v2/jobs>

For more details about these parameters.

Parameters

- **source_project_dataset_table** (*string*) – The dotted (<project>.<project>:<dataset>.<table> BigQuery table to use as the source data. If <project> is not included, project will be the project defined in the connection json.
- **destination_cloud_storage_uris** (*list*) – The destination Google Cloud Storage URI (e.g. gs://some-bucket/some-file.txt). Follows convention defined here: <https://cloud.google.com/bigquery/exporting-data-from-bigquery#exportingmultiple>
- **compression** (*string*) – Type of compression to use.
- **export_format** – File format to export.
- **field_delimiter** (*string*) – The delimiter to use when extracting to a CSV.
- **print_header** (*boolean*) – Whether to print a header for a CSV file extract.
- **bigquery_conn_id** (*string*) – reference to a specific BigQuery hook.
- **delegate_to** (*string*) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.

```
class airflow.contrib.operators.databricks_operator.DatabricksSubmitRunOperator (json=None,
                                                                              spark_jar_task=None,
                                                                              notebook_task=None,
                                                                              new_cluster=None,
                                                                              existing_cluster_id=None,
                                                                              libraries=None,
                                                                              run_name=None,
                                                                              timeout_seconds=None,
                                                                              databricks_conn_id=None,
                                                                              polling_period_seconds=None,
                                                                              databricks_retry_limit=None,
                                                                              **kwargs)
```

Submits an Spark job run to Databricks using the [api/2.0/jobs/runs/submit](#) API endpoint.

There are two ways to instantiate this operator.

In the first way, you can take the JSON payload that you typically use to call the `api/2.0/jobs/runs/submit` endpoint and pass it directly to our `DatabricksSubmitRunOperator` through the `json` parameter. For example

```
json = {
    'new_cluster': {
        'spark_version': '2.1.0-db3-scala2.11',
        'num_workers': 2
    },
    'notebook_task': {
        'notebook_path': '/Users/airflow@example.com/PrepareData',
    },
}
```

```
}
notebook_run = DatabricksSubmitRunOperator(task_id='notebook_run', json=json)
```

Another way to accomplish the same thing is to use the named parameters of the `DatabricksSubmitRunOperator` directly. Note that there is exactly one named parameter for each top level parameter in the runs/submit endpoint. In this method, your code would look like this:

```
new_cluster = {
    'spark_version': '2.1.0-db3-scala2.11',
    'num_workers': 2
}
notebook_task = {
    'notebook_path': '/Users/airflow@example.com/PrepareData',
}
notebook_run = DatabricksSubmitRunOperator(
    task_id='notebook_run',
    new_cluster=new_cluster,
    notebook_task=notebook_task)
```

In the case where both the json parameter **AND** the named parameters are provided, they will be merged together. If there are conflicts during the merge, the named parameters will take precedence and override the top level json keys.

Currently the named parameters that `DatabricksSubmitRunOperator` supports are

- `spark_jar_task`
- `notebook_task`
- `new_cluster`
- `existing_cluster_id`
- `libraries`
- `run_name`
- `timeout_seconds`

Parameters

- **json** (*dict*) – A JSON object containing API parameters which will be passed directly to the `api/2.0/jobs/runs/submit` endpoint. The other named parameters (i.e. `spark_jar_task`, `notebook_task`..) to this operator will be merged with this json dictionary if they are provided. If there are conflicts during the merge, the named parameters will take precedence and override the top level json keys. This field will be templated.

See also:

For more information about templating see *Jinja Templating*. <https://docs.databricks.com/api/latest/jobs.html#runs-submit>

- **spark_jar_task** (*dict*) – The main class and parameters for the JAR task. Note that the actual JAR is specified in the `libraries`. *EITHER* `spark_jar_task` *OR* `notebook_task` should be specified. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/jobs.html#jobssparkjartask>

- **notebook_task** (*dict*) – The notebook path and parameters for the notebook task. *EITHER* `spark_jar_task` *OR* `notebook_task` should be specified. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/jobs.html#jobsnotebooktask>

- **new_cluster** (*dict*) – Specs for a new cluster on which this task will be run. *EITHER* `new_cluster` *OR* `existing_cluster_id` should be specified. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/jobs.html#jobsclusterspecnewcluster>

- **existing_cluster_id** (*string*) – ID for existing cluster on which to run this task. *EITHER* `new_cluster` *OR* `existing_cluster_id` should be specified. This field will be templated.
- **libraries** (*list of dicts*) – Libraries which this run will use. This field will be templated.

See also:

<https://docs.databricks.com/api/latest/libraries.html#managedlibrarieslibrary>

- **run_name** (*string*) – The run name used for this task. By default this will be set to the Airflow `task_id`. This `task_id` is a required parameter of the superclass `BaseOperator`. This field will be templated.
- **timeout_seconds** (*int32*) – The timeout for this run. By default a value of 0 is used which means to have no timeout. This field will be templated.
- **databricks_conn_id** (*string*) – The name of the Airflow connection to use. By default and in the common case this will be `databricks_default`.
- **polling_period_seconds** (*int*) – Controls the rate which we poll for the result of this run. By default the operator will poll every 30 seconds.
- **databricks_retry_limit** (*int*) – Amount of times retry if the Databricks backend is unreachable. Its value must be greater than or equal to 1.

```
class airflow.contrib.operators.ecs_operator.ECSOperator(task_definition, cluster, overrides, aws_conn_id=None, region_name=None, **kwargs)
```

Execute a task on AWS EC2 Container Service

Parameters

- **task_definition** (*str*) – the task definition name on EC2 Container Service
- **cluster** (*str*) – the cluster name on EC2 Container Service
- **aws_conn_id** (*str*) – connection id of AWS credentials / region name. If None, credential boto3 strategy will be used (<http://boto3.readthedocs.io/en/latest/guide/configuration.html>).
- **region_name** – region name to use in AWS Hook. Override the `region_name` in connection (if provided)

Param overrides: the same parameter that boto3 will receive: http://boto3.readthedocs.org/en/latest/reference/services/ecs.html#ECS.Client.run_task

Type overrides: dict

```
class airflow.contrib.operators.gcs_download_operator.GoogleCloudStorageDownloadOperator (bucket, object, filename, store_to_xcom_key, google_cloud_storage_conn_id, delegate_to, *args, **kwargs)
```

Downloads a file from Google Cloud Storage.

Parameters

- **bucket** (*string*) – The Google cloud storage bucket where the object is.
- **object** (*string*) – The name of the object to download in the Google cloud storage bucket.
- **filename** (*string*) – The file path on the local file system (where the operator is being executed) that the file should be downloaded to. If false, the downloaded data will not be stored on the local file system.
- **store_to_xcom_key** (*string*) – If this param is set, the operator will push the contents of the downloaded file to XCom with the key set in this parameter. If false, the downloaded data will not be pushed to XCom.
- **google_cloud_storage_conn_id** (*string*) – The connection ID to use when connecting to Google cloud storage.
- **delegate_to** (*string*) – The account to impersonate, if any. For this to work, the service account making the request must have domain-wide delegation enabled.

```
class airflow.contrib.operators.hipchat_operator.HipChatAPIOperator (token, base_url='https://api.hipchat.com/v2', *args, **kwargs)
```

Base HipChat Operator. All derived HipChat operators reference from HipChat's official REST API documentation at <https://www.hipchat.com/docs/apiv2>. Before using any HipChat API operators you need to get an authentication token at <https://www.hipchat.com/docs/apiv2/auth>. In the future additional HipChat operators will be derived from this class as well.

Parameters

- **token** (*str*) – HipChat REST API authentication token
- **base_url** (*str*) – HipChat REST API base url.

```
class airflow.contrib.operators.hipchat_operator.HipChatAPISendRoomNotificationOperator (room_id, message, *args, **kwargs)
```

Send notification to a specific HipChat room. More info: https://www.hipchat.com/docs/apiv2/method/send_room_notification

Parameters

- **room_id** (*str*) – Room in which to send notification on HipChat
- **message** (*str*) – The message body
- **frm** (*str*) – Label to be shown in addition to sender’s name
- **message_format** (*str*) – How the notification is rendered: html or text
- **color** (*str*) – Background color of the msg: yellow, green, red, purple, gray, or random
- **attach_to** (*str*) – The message id to attach this notification to
- **notify** (*bool*) – Whether this message should trigger a user notification
- **card** (*dict*) – HipChat-defined card object

3.17.2 Macros

Here’s a list of variables and macros that can be used in templates

3.17.2.1 Default Variables

The Airflow engine passes a few variables by default that are accessible in all templates

Variable	Description
{{ ds }}	the execution date as YYYY-MM-DD
{{ ds_nodash }}	the execution date as YYYYMMDD
{{ yesterday_ds }}	yesterday's date as YYYY-MM-DD
{{ yesterday_ds_nodash }}	yesterday's date as YYYYMMDD
{{ tomorrow_ds }}	tomorrow's date as YYYY-MM-DD
{{ tomorrow_ds_nodash }}	tomorrow's date as YYYYMMDD
{{ ts }}	same as <code>execution_date.isoformat()</code>
{{ ts_nodash }}	same as <code>ts</code> without <code>-</code> and <code>:</code>
{{ execution_date }}	the <code>execution_date</code> , (<code>datetime.datetime</code>)
{{ prev_execution_date }}	the previous execution date (if available) (<code>datetime.datetime</code>)
{{ next_execution_date }}	the next execution date (<code>datetime.datetime</code>)
{{ dag }}	the DAG object
{{ task }}	the Task object
{{ macros }}	a reference to the macros package, described below
{{ task_instance }}	the <code>task_instance</code> object
{{ end_date }}	same as {{ ds }}
{{ latest_date }}	same as {{ ds }}
{{ ti }}	same as {{ task_instance }}
{{ params }}	a reference to the user-defined params dictionary
{{ var.value.my_var }}	global defined variables represented as a dictionary
{{ var.json.my_var.path }}	global defined variables represented as a dictionary with deserialized JSON object, append the path to the key within the JSON object
{{ task_instance_key_str }}	a unique, human-readable key to the task instance formatted as <code>{dag_id}_{task_id}_{ds}</code>
conf	the full configuration object located at <code>airflow.configuration.conf</code> which represents the content of your <code>airflow.cfg</code>
run_id	the <code>run_id</code> of the current DAG run
dag_run	a reference to the <code>DagRun</code> object
test_mode	whether the task instance was called using the CLI's test subcommand

Note that you can access the object's attributes and methods with simple dot notation. Here are some examples of what is possible: `{{ task.owner }}`, `{{ task.task_id }}`, `{{ ti.hostname }}`, ... Refer to the models documentation for more information on the objects' attributes and methods.

The `var` template variable allows you to access variables defined in Airflow's UI. You can access them as either plain-text or JSON. If you use JSON, you are also able to walk nested structures, such as dictionaries like: `{{ var.json.my_dict_var.key1 }}`

3.17.2.2 Macros

Macros are a way to expose objects to your templates and live under the `macros` namespace in your templates.

A few commonly used libraries and methods are made available.

Variable	Description
<code>macros.datetime</code>	The standard lib's <code>datetime.datetime</code>
<code>macros.timedelta</code>	The standard lib's <code>datetime.timedelta</code>
<code>macros.dateutil</code>	A reference to the <code>dateutil</code> package
<code>macros.time</code>	The standard lib's <code>time</code>
<code>macros.uuid</code>	The standard lib's <code>uuid</code>
<code>macros.random</code>	The standard lib's <code>random</code>

Some airflow specific macros are also defined:

`airflow.macros.ds_add(ds, days)`

Add or subtract days from a YYYY-MM-DD

Parameters

- **ds** (*str*) – anchor date in YYYY-MM-DD format to add to
- **days** (*int*) – number of days to add to the ds, you can use negative values

```
>>> ds_add('2015-01-01', 5)
'2015-01-06'
>>> ds_add('2015-01-06', -5)
'2015-01-01'
```

`airflow.macros.ds_format(ds, input_format, output_format)`

Takes an input string and outputs another string as specified in the output format

Parameters

- **ds** (*str*) – input string which contains a date
- **input_format** (*str*) – input string format. E.g. `%Y-%m-%d`
- **output_format** (*str*) – output string format E.g. `%Y-%m-%d`

```
>>> ds_format('2015-01-01', "%Y-%m-%d", "%m-%d-%Y")
'01-01-15'
>>> ds_format('1/5/2015', "%m/%d/%Y", "%Y-%m-%d")
'2015-01-05'
```

`airflow.macros.random()` → x in the interval [0, 1).

`airflow.macros.hive.closest_ds_partition(table, ds, before=True, schema='default', metastore_conn_id='metastore_default')`

This function finds the date in a list closest to the target date. An optional parameter can be given to get the closest before or after.

Parameters

- **table** (*str*) – A hive table name
- **ds** (*datetime.date list*) – A timestamp `%Y-%m-%d` e.g. `yyyy-mm-dd`
- **before** (*bool or None*) – closest before (True), after (False) or either side of ds

Returns The closest date

Return type str or None

```
>>> tbl = 'airflow.static_babynames_partitioned'
>>> closest_ds_partition(tbl, '2015-01-02')
'2015-01-01'
```

`airflow.macros.hive.max_partition` (*table*, *schema*='default', *field*=None, *filter*=None, *metastore_conn_id*='metastore_default')

Gets the max partition for a table.

Parameters

- **schema** (*string*) – The hive schema the table lives in
- **table** (*string*) – The hive table you are interested in, supports the dot notation as in “my_database.my_table”, if a dot is found, the schema param is disregarded
- **hive_conn_id** (*string*) – The hive connection you are interested in. If your default is set you don’t need to use this parameter.
- **filter** (*string*) – filter on a subset of partition as in *sub_part*='specific_value'
- **field** – the field to get the max value from. If there’s only one partition field, this will be inferred

```
>>> max_partition('airflow.static_babynames_partitioned')
'2015-01-01'
```

3.17.3 Models

Models are built on top of the SQLAlchemy ORM Base class, and instances are persisted in the database.

```
class airflow.models.DAG(dag_id, description=u'', schedule_interval=datetime.timedelta(1),
                        start_date=None, end_date=None, full_filepath=None,
                        template_searchpath=None, user_defined_macros=None,
                        user_defined_filters=None, default_args=None, concurrency=16,
                        max_active_runs=16, dagrun_timeout=None, sla_miss_callback=None,
                        default_view=u'tree', orientation='LR', catchup=True, params=None)
```

Bases: `airflow.dag.base_dag.BaseDag`, `airflow.utils.logging.LoggingMixin`

A dag (directed acyclic graph) is a collection of tasks with directional dependencies. A dag also has a schedule, a start end an end date (optional). For each schedule, (say daily or hourly), the DAG needs to run each individual tasks as their dependencies are met. Certain tasks have the property of depending on their own past, meaning that they can’t run until their previous schedule (and upstream tasks) are completed.

DAGs essentially act as namespaces for tasks. A task_id can only be added once to a DAG.

Parameters

- **dag_id** (*string*) – The id of the DAG
- **description** (*string*) – The description for the DAG to e.g. be shown on the web-server
- **schedule_interval** (*datetime.timedelta or dateutil.relativedelta.relativedelta or str that acts as a cron expression*) – Defines how often that DAG runs, this timedelta object gets added to your latest task instance’s execution_date to figure out the next schedule
- **start_date** (*datetime.datetime*) – The timestamp from which the scheduler will attempt to backfill

- **end_date** (*datetime.datetime*) – A date beyond which your DAG won’t run, leave to None for open ended scheduling
- **template_searchpath** (*string or list of strings*) – This list of folders (non relative) defines where jinja will look for your templates. Order matters. Note that jinja/airflow includes the path of your DAG file by default
- **user_defined_macros** (*dict*) – a dictionary of macros that will be exposed in your jinja templates. For example, passing `dict(foo='bar')` to this argument allows you to `{{ foo }}` in all jinja templates related to this DAG. Note that you can pass any type of object here.
- **user_defined_filters** (*dict*) – a dictionary of filters that will be exposed in your jinja templates. For example, passing `dict(hello=lambda name: 'Hello %s' % name)` to this argument allows you to `{{ 'world' | hello }}` in all jinja templates related to this DAG.
- **default_args** (*dict*) – A dictionary of default parameters to be used as constructor keyword parameters when initialising operators. Note that operators have the same hook, and precede those defined here, meaning that if your dict contains `'depends_on_past': True` here and `'depends_on_past': False` in the operator’s call `default_args`, the actual value will be `False`.
- **params** (*dict*) – a dictionary of DAG level parameters that are made accessible in templates, namespaced under *params*. These params can be overridden at the task level.
- **concurrency** (*int*) – the number of task instances allowed to run concurrently
- **max_active_runs** (*int*) – maximum number of active DAG runs, beyond this number of DAG runs in a running state, the scheduler won’t create new active DAG runs
- **dagrun_timeout** (*datetime.timedelta*) – specify how long a DagRun should be up before timing out / failing, so that new DagRuns can be created
- **sla_miss_callback** (*types.FunctionType*) – specify a function to call when reporting SLA timeouts.
- **default_view** (*string*) – Specify DAG default view (tree, graph, duration, gantt, landing_times)
- **orientation** (*string*) – Specify DAG orientation in graph view (LR, TB, RL, BT)
- **catchup** – Perform scheduler catchup (or only run latest)? Defaults to True

“type catchup: bool”

add_task (*task*)

Add a task to the DAG

Parameters *task* (*task*) – the task you want to add

add_tasks (*tasks*)

Add a list of tasks to the DAG

Parameters *tasks* (*list of tasks*) – a list of tasks you want to add

clear (*start_date=None, end_date=None, only_failed=False, only_running=False, confirm_prompt=False, include_subdags=True, reset_dag_runs=True, dry_run=False*)

Clears a set of task instances associated with the current dag for a specified date range.

cli ()

Exposes a CLI specific to this DAG

concurrency_reached

Returns a boolean indicating whether the concurrency limit for this DAG has been reached

create_dagrun (*args, **kwargs)

Creates a dag run from this dag including the tasks associated with this dag. Returns the dag run.

Parameters

- **run_id** (*string*) – defines the the run id for this dag run
- **execution_date** (*datetime*) – the execution date of this dag run
- **state** (*State*) – the state of the dag run
- **start_date** (*datetime*) – the date this dag run should be evaluated
- **external_trigger** (*bool*) – whether this dag run is externally triggered
- **session** (*Session*) – database session

static deactivate_stale_dags (*args, **kwargs)

Deactivate any DAGs that were last touched by the scheduler before the expiration date. These DAGs were likely deleted.

Parameters **expiration_date** – set inactive DAGs that were touched before this

time :type expiration_date: datetime :return: None

static deactivate_unknown_dags (*args, **kwargs)

Given a list of known DAGs, deactivate any other DAGs that are marked as active in the ORM

Parameters **active_dag_ids** (*list [unicode]*) – list of DAG IDs that are active

Returns None

filepath

File location of where the dag object is instantiated

folder

Folder location of where the dag object is instantiated

get_active_runs (*args, **kwargs)

Returns a list of “running” tasks :param session: :return: List of execution dates

get_dagrun (*args, **kwargs)

Returns the dag run for a given execution date if it exists, otherwise none. :param execution_date: The execution date of the DagRun to find. :param session: :return: The DagRun if found, otherwise None.

get_last_dagrun (*args, **kwargs)

Returns the last dag run for this dag, None if there was none. Last dag run can be any type of run eg. scheduled or backfilled. Overridden DagRuns are ignored

static get_num_task_instances (*args, **kwargs)

Returns the number of task instances in the given DAG.

Parameters

- **session** – ORM session
- **dag_id** (*unicode*) – ID of the DAG to get the task concurrency of
- **task_ids** (*list [unicode]*) – A list of valid task IDs for the given DAG
- **states** (*list [state]*) – A list of states to filter by if supplied

Returns The number of running tasks

Return type int

get_template_env()

Returns a jinja2 Environment while taking into account the DAGs template_searchpath, user_defined_macros and user_defined_filters

is_paused

Returns a boolean indicating whether this DAG is paused

latest_execution_date

Returns the latest date for which at least one dag run exists

normalize_schedule (dtm)

Returns dtm + interval unless dtm is first interval then it returns dtm

run (start_date=None, end_date=None, mark_success=False, include_adhoc=False, local=False, executor=None, do_not_pickle=False, ignore_task_deps=False, ignore_first_depends_on_past=False, pool=None)
Runs the DAG.

set_dependency (upstream_task_id, downstream_task_id)

Simple utility method to set dependency between two tasks that already have been added to the DAG using add_task()

sub_dag (task_regex, include_downstream=False, include_upstream=True)

Returns a subset of the current dag as a deep copy of the current dag based on a regex that should match one or many tasks, and includes upstream and downstream neighbours based on the flag passed.

subdags

Returns a list of the subdag objects associated to this DAG

static sync_to_db (*args, **kwargs)

Save attributes about this DAG to the DB. Note that this method can be called for both DAGs and SubDAGs. A SubDag is actually a SubDagOperator.

Parameters dag (DAG) – the DAG object to save to the DB

:own :param sync_time: The time that the DAG should be marked as sync'ed :type sync_time: datetime
:return: None

topological_sort ()

Sorts tasks in topographical order, such that a task comes after any of its upstream dependencies.

Heavily inspired by: <http://blog.jupo.org/2012/04/06/topological-sorting-acyclic-directed-graphs/> :return: list of tasks in topological order

tree_view ()

Shows an ascii tree representation of the DAG

```
class airflow.models.BaseOperator(task_id, owner='Airflow', email=None,
                                  email_on_retry=True, email_on_failure=True, re-
                                  tries=0, retry_delay=datetime.timedelta(0, 300),
                                  retry_exponential_backoff=False, max_retry_delay=None,
                                  start_date=None, end_date=None, schedule_interval=None,
                                  depends_on_past=False, wait_for_downstream=False,
                                  dag=None, params=None, default_args=None, adhoc=False,
                                  priority_weight=1, queue='default', pool=None, sla=None,
                                  execution_timeout=None, on_failure_callback=None,
                                  on_success_callback=None, on_retry_callback=None, trig-
                                  ger_rule=u'all_success', resources=None, run_as_user=None,
                                  *args, **kwargs)
```

Bases: future.types.newobject.newobject

Abstract base class for all operators. Since operators create objects that become node in the dag, BaseOperator contains many recursive methods for dag crawling behavior. To derive this class, you are expected to override the constructor as well as the 'execute' method.

Operators derived from this class should perform or trigger certain tasks synchronously (wait for completion). Example of operators could be an operator the runs a Pig job (PigOperator), a sensor operator that waits for a partition to land in Hive (HiveSensorOperator), or one that moves data from Hive to MySQL (Hive2MySQLOperator). Instances of these operators (tasks) target specific operations, running specific scripts, functions or data transfers.

This class is abstract and shouldn't be instantiated. Instantiating a class derived from this one results in the creation of a task object, which ultimately becomes a node in DAG objects. Task dependencies should be set by using the `set_upstream` and/or `set_downstream` methods.

Note that this class is derived from SQLAlchemy's Base class, which allows us to push metadata regarding tasks to the database. Deriving this classes needs to implement the polymorphic specificities documented in SQLAlchemy. This should become clear while reading the code for other operators.

Parameters

- **task_id** (*string*) – a unique, meaningful id for the task
- **owner** (*string*) – the owner of the task, using the unix username is recommended
- **retries** (*int*) – the number of retries that should be performed before failing the task
- **retry_delay** (*timedelta*) – delay between retries
- **retry_exponential_backoff** (*bool*) – allow progressive longer waits between retries by using exponential backoff algorithm on retry delay (delay will be converted into seconds)
- **max_retry_delay** (*timedelta*) – maximum delay interval between retries
- **start_date** (*datetime*) – The `start_date` for the task, determines the `execution_date` for the first task instance. The best practice is to have the `start_date` rounded to your DAG's `schedule_interval`. Daily jobs have their `start_date` some day at 00:00:00, hourly jobs have their `start_date` at 00:00 of a specific hour. Note that Airflow simply looks at the latest `execution_date` and adds the `schedule_interval` to determine the next `execution_date`. It is also very important to note that different tasks' dependencies need to line up in time. If task A depends on task B and their `start_date` are offset in a way that their `execution_date` don't line up, A's dependencies will never be met. If you are looking to delay a task, for example running a daily task at 2AM, look into the `TimeSensor` and `TimeDeltaSensor`. We advise against using dynamic `start_date` and recommend using fixed ones. Read the FAQ entry about `start_date` for more information.
- **end_date** (*datetime*) – if specified, the scheduler won't go beyond this date
- **depends_on_past** (*bool*) – when set to true, task instances will run sequentially while relying on the previous task's schedule to succeed. The task instance for the `start_date` is allowed to run.
- **wait_for_downstream** (*bool*) – when set to true, an instance of task X will wait for tasks immediately downstream of the previous instance of task X to finish successfully before it runs. This is useful if the different instances of a task X alter the same asset, and this asset is used by tasks downstream of task X. Note that `depends_on_past` is forced to True wherever `wait_for_downstream` is used.
- **queue** (*str*) – which queue to target when running this job. Not all executors implement queue management, the CeleryExecutor does support targeting specific queues.

- **dag** (*DAG*) – a reference to the dag the task is attached to (if any)
- **priority_weight** (*int*) – priority weight of this task against other task. This allows the executor to trigger higher priority tasks before others when things get backed up.
- **pool** (*str*) – the slot pool this task should run in, slot pools are a way to limit concurrency for certain tasks
- **sla** (*datetime.timedelta*) – time by which the job is expected to succeed. Note that this represents the *timedelta* after the period is closed. For example if you set an SLA of 1 hour, the scheduler would send an email soon after 1:00AM on the 2016-01-02 if the 2016-01-01 instance has not succeeded yet. The scheduler pays special attention for jobs with an SLA and sends alert emails for sla misses. SLA misses are also recorded in the database for future reference. All tasks that share the same SLA time get bundled in a single email, sent soon after that time. SLA notification are sent once and only once for each task instance.
- **execution_timeout** (*datetime.timedelta*) – max time allowed for the execution of this task instance, if it goes beyond it will raise and fail.
- **on_failure_callback** (*callable*) – a function to be called when a task instance of this task fails. a context dictionary is passed as a single parameter to this function. Context contains references to related objects to the task instance and is documented under the macros section of the API.
- **on_retry_callback** – much like the *on_failure_callback* except that it is executed when retries occur.
- **on_success_callback** (*callable*) – much like the *on_failure_callback* except that it is executed when the task succeeds.
- **trigger_rule** (*str*) – defines the rule by which dependencies are applied for the task to get triggered. Options are: { *all_success* | *all_failed* | *all_done* | *one_success* | *one_failed* | *dummy* } default is *all_success*. Options can be set as string or using the constants defined in the static class *airflow.utils.TriggerRule*
- **resources** (*dict*) – A map of resource parameter names (the argument names of the Resources constructor) to their values.
- **run_as_user** (*str*) – unix username to impersonate while running the task

clear (*start_date=None, end_date=None, upstream=False, downstream=False*)

Clears the state of task instances associated with the task, following the parameters specified.

dag

Returns the Operator's DAG if set, otherwise raises an error

deps

Returns the list of dependencies for the operator. These differ from execution context dependencies in that they are specific to tasks and can be extended/overridden by subclasses.

detect_downstream_cycle (*task=None*)

When invoked, this routine will raise an exception if a cycle is detected downstream from self. It is invoked when tasks are added to the DAG to detect cycles.

downstream_list

@property: list of tasks directly downstream

execute (*context*)

This is the main method to derive when creating an operator. Context is the same dictionary used as when rendering jinja templates.

Refer to `get_template_context` for more context.

get_direct_relatives (*upstream=False*)

Get the direct relatives to the current task, upstream or downstream.

get_flat_relatives (*upstream=False, l=None*)

Get a flat list of relatives, either upstream or downstream.

get_task_instances (*session, start_date=None, end_date=None*)

Get a set of task instance related to this task for a specific date range.

has_dag ()

Returns True if the Operator has been assigned to a DAG.

on_kill ()

Override this method to cleanup subprocesses when a task instance gets killed. Any use of the threading, subprocess or multiprocessing module within an operator needs to be cleaned up or it will leave ghost processes behind.

post_execute (*context, result=None*)

This hook is triggered right after `self.execute()` is called. It is passed the execution context and any results returned by the operator.

pre_execute (*context*)

This hook is triggered right before `self.execute()` is called.

prepare_template ()

Hook that is triggered after the templated fields get replaced by their content. If you need your operator to alter the content of the file before the template is rendered, it should override this method to do so.

render_template (*attr, content, context*)

Renders a template either from a file or directly in a field, and returns the rendered result.

render_template_from_field (*attr, content, context, jinja_env*)

Renders a template from a field. If the field is a string, it will simply render the string and return the result. If it is a collection or nested set of collections, it will traverse the structure and render all strings in it.

run (*start_date=None, end_date=None, ignore_first_depends_on_past=False, ignore_ti_state=False, mark_success=False*)

Run a set of task instances for a date range.

schedule_interval

The schedule interval of the DAG always wins over individual tasks so that tasks within a DAG always line up. The task still needs a `schedule_interval` as it may not be attached to a DAG.

set_downstream (*task_or_task_list*)

Set a task, or a task task to be directly downstream from the current task.

set_upstream (*task_or_task_list*)

Set a task, or a task task to be directly upstream from the current task.

upstream_list

@property: list of tasks directly upstream

xcom_pull (*context, task_ids, dag_id=None, key='return_value', include_prior_dates=None*)

See `TaskInstance.xcom_pull()`

xcom_push (*context, key, value, execution_date=None*)

See `TaskInstance.xcom_push()`

class `airflow.models.TaskInstance` (*task, execution_date, state=None*)

Bases: `sqlalchemy.ext.declarative.api.Base`

Task instances store the state of a task instance. This table is the authority and single source of truth around what tasks have run and the state they are in.

The SQLAlchemy model doesn't have a SQLAlchemy foreign key to the task or dag model deliberately to have more control over transactions.

Database transactions on this table should insure double triggers and any confusion around what task instances are or aren't ready to run even while multiple schedulers may be firing task instances.

are_dependencies_met (*args, **kwargs)

Returns whether or not all the conditions are met for this task instance to be run given the context for the dependencies (e.g. a task instance being force run from the UI will ignore some dependencies).

Parameters

- **dep_context** (*DepContext*) – The execution context that determines the dependencies that should be evaluated.
- **session** (*Session*) – database session
- **verbose** (*boolean*) – whether or not to print details on failed dependencies

are_dependents_done (*args, **kwargs)

Checks whether the dependents of this task instance have all succeeded. This is meant to be used by wait_for_downstream.

This is useful when you do not want to start processing the next schedule of a task until the dependents are done. For instance, if the task DROPS and recreates a table.

clear_xcom_data (*args, **kwargs)

Clears all XCom data from the database for the task instance

command (mark_success=False, ignore_all_deps=False, ignore_depends_on_past=False, ignore_task_deps=False, ignore_ti_state=False, local=False, pickle_id=None, raw=False, job_id=None, pool=None, cfg_path=None)

Returns a command that can be executed anywhere where airflow is installed. This command is part of the message sent to executors by the orchestrator.

command_as_list (mark_success=False, ignore_all_deps=False, ignore_task_deps=False, ignore_depends_on_past=False, ignore_ti_state=False, local=False, pickle_id=None, raw=False, job_id=None, pool=None, cfg_path=None)

Returns a command that can be executed anywhere where airflow is installed. This command is part of the message sent to executors by the orchestrator.

current_state (*args, **kwargs)

Get the very latest state from the database, if a session is passed, we use and looking up the state becomes part of the session, otherwise a new session is used.

error (*args, **kwargs)

Forces the task instance's state to FAILED in the database.

static generate_command (dag_id, task_id, execution_date, mark_success=False, ignore_all_deps=False, ignore_depends_on_past=False, ignore_task_deps=False, ignore_ti_state=False, local=False, pickle_id=None, file_path=None, raw=False, job_id=None, pool=None, cfg_path=None)

Generates the shell command required to execute this task instance.

Parameters

- **dag_id** (*unicode*) – DAG ID
- **task_id** (*unicode*) – Task ID

- **execution_date** (*datetime*) – Execution date for the task
- **mark_success** (*bool*) – Whether to mark the task as successful
- **ignore_all_deps** (*boolean*) – Ignore all ignorable dependencies. Overrides the other `ignore_*` parameters.
- **ignore_depends_on_past** (*boolean*) – Ignore `depends_on_past` parameter of DAGs (e.g. for Backfills)
- **ignore_task_deps** (*boolean*) – Ignore task-specific dependencies such as `depends_on_past` and trigger rule
- **ignore_ti_state** (*boolean*) – Ignore the task instance’s previous failure/success
- **local** (*bool*) – Whether to run the task locally
- **pickle_id** – If the DAG was serialized to the DB, the ID

associated with the pickled DAG :type pickle_id: unicode :param file_path: path to the file containing the DAG definition :param raw: raw mode (needs more details) :param job_id: job ID (needs more details) :param pool: the Airflow pool that the task should run in :type pool: unicode :return: shell command that can be used to run the task instance

get_dagrun (*args, **kwargs)

Returns the DagRun for this TaskInstance :param session: :return: DagRun

init_on_load ()

Initialize the attributes that aren’t stored in the DB.

is_premature

Returns whether a task is in UP_FOR_RETRY state and its retry interval has elapsed.

key

Returns a tuple that identifies the task instance uniquely

next_retry_datetime ()

Get datetime of the next retry if the task instance fails. For exponential backoff, `retry_delay` is used as base and will be converted to seconds.

pool_full (*args, **kwargs)

Returns a boolean as to whether the slot pool has room for this task to run

previous_ti

The task instance for the task that ran before this task instance

ready_for_retry ()

Checks on whether the task instance is in the right state and timeframe to be retried.

refresh_from_db (*args, **kwargs)

Refreshes the task instance from the database based on the primary key

Parameters lock_for_update – if True, indicates that the database should lock the TaskInstance (issuing a FOR UPDATE clause) until the session is committed.

run (*args, **kwargs)

Runs the task instance.

Parameters

- **verbose** (*boolean*) – whether to turn on more verbose logging
- **ignore_all_deps** (*boolean*) – Ignore all of the non-critical dependencies, just runs
- **ignore_depends_on_past** (*boolean*) – Ignore `depends_on_past` DAG attribute

- **ignore_task_deps** (*boolean*) – Don't check the dependencies of this TI's task
- **ignore_ti_state** (*boolean*) – Disregards previous task instance state
- **mark_success** (*boolean*) – Don't run the task, mark its state as success
- **test_mode** (*boolean*) – Doesn't record success or failure in the DB
- **pool** (*str*) – specifies the pool to use to run the task instance

xcom_pull (*task_ids*, *dag_id=None*, *key=u'return_value'*, *include_prior_dates=False*)

Pull XComs that optionally meet certain criteria.

The default value for *key* limits the search to XComs that were returned by other tasks (as opposed to those that were pushed manually). To remove this filter, pass *key=None* (or any desired value).

If a single *task_id* string is provided, the result is the value of the most recent matching XCom from that *task_id*. If multiple *task_ids* are provided, a tuple of matching values is returned. *None* is returned whenever no matches are found.

Parameters

- **key** (*string*) – A key for the XCom. If provided, only XComs with matching keys will be returned. The default key is 'return_value', also available as a constant `XCOM_RETURN_KEY`. This key is automatically given to XComs returned by tasks (as opposed to being pushed manually). To remove the filter, pass *key=None*.
- **task_ids** (*string or iterable of strings (representing task_ids)*) – Only XComs from tasks with matching ids will be pulled. Can pass *None* to remove the filter.
- **dag_id** (*string*) – If provided, only pulls XComs from this DAG. If *None* (default), the DAG of the calling task is used.
- **include_prior_dates** (*bool*) – If *False*, only XComs from the current execution_date are returned. If *True*, XComs from previous dates are returned as well.

xcom_push (*key*, *value*, *execution_date=None*)

Make an XCom available for tasks to pull.

Parameters

- **key** (*string*) – A key for the XCom
- **value** (*any pickleable object*) – A value for the XCom. The value is pickled and stored in the database.
- **execution_date** (*datetime*) – if provided, the XCom will not be visible until this date. This can be used, for example, to send a message to a task on a future date without it being immediately visible.

class `airflow.models.DagBag` (*dag_folder=None*, *executor=None*, *include_examples=True*)

Bases: `airflow.dag.base_dag.BaseDagBag`, `airflow.utils.logging.LoggingMixin`

A dagbag is a collection of dags, parsed out of a folder tree and has high level configuration settings, like what database to use as a backend and what executor to use to fire off tasks. This makes it easier to run distinct environments for say production and development, tests, or for different teams or security profiles. What would have been system level settings are now dagbag level so that one system can run multiple, independent settings sets.

Parameters

- **dag_folder** (*unicode*) – the folder to scan to find DAGs
- **executor** – the executor to use when executing task instances in this DagBag

- **include_examples** (*bool*) – whether to include the examples that ship with airflow or not
- **sync_to_db** (*bool*) – whether to sync the properties of the DAGs to the metadata DB while finding them, typically should be done by the scheduler job only

bag_dag (*dag, parent_dag, root_dag*)

Adds the DAG into the bag, recurses into sub dags.

collect_dags (*dag_folder=None, only_if_updated=True*)

Given a file path or a folder, this method looks for python modules, imports them and adds them to the dagbag collection.

Note that if a .airflowignore file is found while processing, the directory, it will behaves much like a .gitignore does, ignoring files that match any of the regex patterns specified in the file.

dagbag_report ()

Prints a report around DagBag loading stats

get_dag (*dag_id*)

Gets the DAG out of the dictionary, and refreshes it if expired

kill_zombies (**args, **kwargs*)

Fails tasks that haven't had a heartbeat in too long

process_file (*filepath, only_if_updated=True, safe_mode=True*)

Given a path to a python module or zip file, this method imports the module and look for dag objects within it.

size ()

Returns the amount of dags contained in this dagbag

class airflow.models.**Connection** (*conn_id=None, conn_type=None, host=None, login=None, password=None, schema=None, port=None, extra=None, uri=None*)

Bases: sqlalchemy.ext.declarative.api.Base

Placeholder to store information about different database instances connection information. The idea here is that scripts use references to database instances (*conn_id*) instead of hard coding hostname, logins and passwords when using operators or hooks.

extra_dejson

Returns the extra property by deserializing json.

3.17.4 Hooks

Importer that dynamically loads a class and module from its parent. This allows Airflow to support `from airflow.operators import BashOperator` even though `BashOperator` is actually in `airflow.operators.bash_operator`.

The importer also takes over for the `parent_module` by wrapping it. This is required to support attribute-based usage:

```
from airflow import operators
operators.BashOperator(...)
```

class airflow.hooks.**DbApiHook** (**args, **kwargs*)

Bases: airflow.hooks.base_hook.BaseHook

Abstract base class for sql hooks.

bulk_dump (*table, tmp_file*)

Dumps a database table into a tab-delimited file

Parameters

- **table** (*str*) – The name of the source table
- **tmp_file** (*str*) – The path of the target file

bulk_load (*table, tmp_file*)

Loads a tab-delimited file into a database table

Parameters

- **table** (*str*) – The name of the target table
- **tmp_file** (*str*) – The path of the file to load into the table

get_conn ()

Returns a connection object

get_cursor ()

Returns a cursor

get_first (*sql, parameters=None*)

Executes the sql and returns the first resulting row.

Parameters

- **sql** (*str or list*) – the sql statement to be executed (*str*) or a list of sql statements to execute
- **parameters** (*mapping or iterable*) – The parameters to render the SQL query with.

get_pandas_df (*sql, parameters=None*)

Executes the sql and returns a pandas dataframe

Parameters

- **sql** (*str or list*) – the sql statement to be executed (*str*) or a list of sql statements to execute
- **parameters** (*mapping or iterable*) – The parameters to render the SQL query with.

get_records (*sql, parameters=None*)

Executes the sql and returns a set of records.

Parameters

- **sql** (*str or list*) – the sql statement to be executed (*str*) or a list of sql statements to execute
- **parameters** (*mapping or iterable*) – The parameters to render the SQL query with.

insert_rows (*table, rows, target_fields=None, commit_every=1000*)

A generic way to insert a set of tuples into a table, a new transaction is created every *commit_every* rows

Parameters

- **table** (*str*) – Name of the target table
- **rows** (*iterable of tuples*) – The rows to insert into the table
- **target_fields** (*iterable of strings*) – The names of the columns to fill in the table

- **commit_every** (*int*) – The maximum number of rows to insert in one transaction. Set to 0 to insert all rows in one transaction.

run (*sql*, *autocommit=False*, *parameters=None*)

Runs a command or a list of commands. Pass a list of sql statements to the *sql* parameter to get them to execute sequentially

Parameters

- **sql** (*str* or *list*) – the sql statement to be executed (*str*) or a list of sql statements to execute
- **autocommit** (*bool*) – What to set the connection's autocommit setting to before executing the query.
- **parameters** (*mapping* or *iterable*) – The parameters to render the SQL query with.

class `airflow.hooks.HttpHook` (*method='POST'*, *http_conn_id='http_default'*)

Bases: `airflow.hooks.base_hook.BaseHook`

Interact with HTTP servers.

get_conn (*headers*)

Returns http session for use with requests

run (*endpoint*, *data=None*, *headers=None*, *extra_options=None*)

Performs the request

run_and_check (*session*, *prepped_request*, *extra_options*)

Grabs extra options like timeout and actually runs the request, checking for the result

class `airflow.hooks.SqliteHook` (**args*, ***kwargs*)

Bases: `airflow.hooks.dbapi_hook.DbApiHook`

Interact with SQLite.

get_conn ()

Returns a sqlite connection object

3.17.4.1 Community contributed hooks

Importer that dynamically loads a class and module from its parent. This allows Airflow to support `from airflow.operators import BashOperator` even though `BashOperator` is actually in `airflow.operators.bash_operator`.

The importer also takes over for the `parent_module` by wrapping it. This is required to support attribute-based usage:

```
from airflow import operators
operators.BashOperator(...)
```

class `airflow.contrib.hooks.BigQueryHook` (*bigquery_conn_id='bigquery_default'*, *delete_to=None*)

Bases: `airflow.contrib.hooks.gcp_api_base_hook.GoogleCloudBaseHook`, `airflow.hooks.dbapi_hook.DbApiHook`

Interact with BigQuery. This hook uses the Google Cloud Platform connection.

get_conn ()

Returns a BigQuery PEP 249 connection object.

get_pandas_df (*bql*, *parameters=None*, *dialect='legacy'*)

Returns a Pandas DataFrame for the results produced by a BigQuery query. The DbApiHook method must be overridden because Pandas doesn't support PEP 249 connections, except for SQLite. See:

<https://github.com/pydata/pandas/blob/master/pandas/io/sql.py#L447> <https://github.com/pydata/pandas/issues/6900>

Parameters

- **bql** (*string*) – The BigQuery SQL to execute.
- **parameters** (*mapping or iterable*) – The parameters to render the SQL query with (not used, leave to override superclass method)
- **dialect** (*string in {'legacy', 'standard'}, default 'legacy'*) – Dialect of BigQuery SQL – legacy SQL or standard SQL

get_service ()

Returns a BigQuery service object.

insert_rows (*table*, *rows*, *target_fields=None*, *commit_every=1000*)

Insertion is currently unsupported. Theoretically, you could use BigQuery's streaming API to insert rows into a table, but this hasn't been implemented.

table_exists (*project_id*, *dataset_id*, *table_id*)

Checks for the existence of a table in Google BigQuery.

Parameters **project_id** – The Google cloud project in which to look for the table. The connection supplied to the hook

must provide access to the specified project. :type project_id: string :param dataset_id: The name of the dataset in which to look for the table.

storage bucket.

Parameters **table_id** (*string*) – The name of the table to check the existence of.

class airflow.contrib.hooks.**GoogleCloudStorageHook** (*google_cloud_storage_conn_id='google_cloud_storage_default'*, *delegate_to=None*)

Bases: airflow.contrib.hooks.gcp_api_base_hook.GoogleCloudBaseHook

Interact with Google Cloud Storage. This hook uses the Google Cloud Platform connection.

delete (*bucket*, *object*, *generation=None*)

Delete an object if versioning is not enabled for the bucket, or if generation parameter is used. :param bucket: name of the bucket, where the object resides :type bucket: string :param object: name of the object to delete :type object: string :param generation: if present, permanently delete the object of this generation :type generation: string :return: True if succeeded

download (*bucket*, *object*, *filename=False*)

Get a file from Google Cloud Storage.

Parameters

- **bucket** (*string*) – The bucket to fetch from.
- **object** (*string*) – The object to fetch.
- **filename** (*string*) – If set, a local file path where the file should be written to.

exists (*bucket*, *object*)

Checks for the existence of a file in Google Cloud Storage.

Parameters

- **bucket** (*string*) – The Google cloud storage bucket where the object is.
- **object** (*string*) – The name of the object to check in the Google cloud storage bucket.

get_conn ()

Returns a Google Cloud Storage service object.

is_updated_after (*bucket, object, ts*)

Checks if an object is updated in Google Cloud Storage.

Parameters

- **bucket** (*string*) – The Google cloud storage bucket where the object is.
- **object** (*string*) – The name of the object to check in the Google cloud storage bucket.
- **ts** (*datetime*) – The timestamp to check against.

list (*bucket, versions=None, maxResults=None, prefix=None*)

List all objects from the bucket with the give string prefix in name :param bucket: bucket name :type bucket: string :param versions: if true, list all versions of the objects :type versions: boolean :param maxResults: max count of items to return in a single page of responses :type maxResults: integer :param prefix: prefix string which filters objects whose name begin with this prefix :type prefix: string :return: a stream of object names matching the filtering criteria

upload (*bucket, object, filename, mime_type='application/octet-stream'*)

Uploads a local file to Google Cloud Storage.

Parameters

- **bucket** (*string*) – The bucket to upload to.
- **object** (*string*) – The object name to set when uploading the local file.
- **filename** (*string*) – The local file path to the file to be uploaded.
- **mime_type** (*string*) – The MIME type to set when uploading the file.

class airflow.contrib.hooks.**FTPHook** (*ftp_conn_id='ftp_default'*)

Bases: airflow.hooks.base_hook.BaseHook

Interact with FTP.

Errors that may occur throughout but should be handled downstream.

close_conn ()

Closes the connection. An error will occur if the connection wasn't ever opened.

create_directory (*path*)

Creates a directory on the remote system.

Parameters **path** (*str*) – full path to the remote directory to create

delete_directory (*path*)

Deletes a directory on the remote system.

Parameters **path** (*str*) – full path to the remote directory to delete

delete_file (*path*)

Removes a file on the FTP Server.

Parameters **path** (*str*) – full path to the remote file

describe_directory (*path*)

Returns a dictionary of {filename: {attributes}} for all files on the remote system (where the MLSD command is supported).

Parameters `path` (*str*) – full path to the remote directory

get_conn ()

Returns a FTP connection object

list_directory (*path*, *nlst=False*)

Returns a list of files on the remote system.

Parameters `path` (*str*) – full path to the remote directory to list

rename (*from_name*, *to_name*)

Rename a file.

Parameters

- **from_name** – rename file from name
- **to_name** – rename file to name

retrieve_file (*remote_full_path*, *local_full_path_or_buffer*)

Transfers the remote file to a local location.

If `local_full_path_or_buffer` is a string path, the file will be put at that location; if it is a file-like buffer, the file will be written to the buffer but not closed.

Parameters

- **remote_full_path** (*str*) – full path to the remote file
- **local_full_path_or_buffer** – full path to the local file or a file-like buffer

store_file (*remote_full_path*, *local_full_path_or_buffer*)

Transfers a local file to the remote location.

If `local_full_path_or_buffer` is a string path, the file will be read from that location; if it is a file-like buffer, the file will be read from the buffer but not closed.

Parameters

- **remote_full_path** (*str*) – full path to the remote file
- **local_full_path_or_buffer** (*str or file-like buffer*) – full path to the local file or a file-like buffer

class `airflow.contrib.hooks.SSHHook` (*conn_id='ssh_default'*)

Bases: `airflow.hooks.base_hook.BaseHook`

Light-weight remote execution library and utilities.

Using this hook (which is just a convenience wrapper for subprocess), is created to let you stream data from a remotely stored file.

As a bonus, `SSHHook` also provides a really cool feature that let's you set up ssh tunnels super easily using a python context manager (there is an example in the integration part of unittests).

Parameters

- **key_file** (*str*) – Typically the SSHHook uses the keys that are used by the user airflow is running under. This sets the behavior to use another file instead.
- **connect_timeout** (*int*) – sets the connection timeout for this connection.
- **no_host_key_check** (*bool*) – whether to check to host key. If True host keys will not be checked, but are also not stored in the current users's known_hosts file.
- **tty** (*bool*) – allocate a tty.

- **sshpas** (*bool*) – Use to non-interactively perform password authentication by using sshpass.

Popen (*cmd*, ***kwargs*)

Remote Popen

Parameters

- **cmd** – command to remotely execute
- **kwargs** – extra arguments to Popen (see subprocess.Popen)

Returns handle to subprocess

check_output (*cmd*)

Executes a remote command and returns the stdout a remote process. Simplified version of Popen when you only want the output as a string and detect any errors.

Parameters **cmd** – command to remotely execute

Returns stdout

tunnel (**args*, ***kws*)

Creates a tunnel between two hosts. Like ssh -L <LOCAL_PORT>:host:<REMOTE_PORT>. Remember to close() the returned “tunnel” object in order to clean up after yourself when you are done with the tunnel.

Parameters

- **local_port** (*int*) –
- **remote_port** (*int*) –
- **remote_host** (*str*) –

Returns

class airflow.contrib.hooks.gcs_hook.**GoogleCloudStorageHook** (*google_cloud_storage_conn_id='google_cloud_*
delegate_to=None)

Interact with Google Cloud Storage. This hook uses the Google Cloud Platform connection.

3.17.5 Executors

Executors are the mechanism by which task instances get run.

class airflow.executors.**LocalExecutor** (*parallelism=32*)

Bases: airflow.executors.base_executor.BaseExecutor

LocalExecutor executes tasks locally in parallel. It uses the multiprocessing Python library and queues to parallelize the execution of tasks.

class airflow.executors.**SequentialExecutor**

Bases: airflow.executors.base_executor.BaseExecutor

This executor will only run one task instance at a time, can be used for debugging. It is also the only executor that can be used with sqlite since sqlite doesn’t support multiple connections.

Since we want airflow to work out of the box, it defaults to this SequentialExecutor alongside sqlite as you first install it.

3.17.5.1 Community-contributed executors

class `airflow.contrib.executors.mesos_executor.MesosExecutor` (*parallelism=32*)

MesosExecutor allows distributing the execution of task instances to multiple mesos workers.

Apache Mesos is a distributed systems kernel which abstracts CPU, memory, storage, and other compute resources away from machines (physical or virtual), enabling fault-tolerant and elastic distributed systems to easily be built and run effectively. See <http://mesos.apache.org/>

a

- `airflow.contrib.hooks`, [122](#)
- `airflow.contrib.operators`, [101](#)
- `airflow.executors`, [126](#)
- `airflow.hooks`, [120](#)
- `airflow.macros`, [109](#)
- `airflow.macros.hive`, [109](#)
- `airflow.models`, [110](#)
- `airflow.operators`, [94](#)

A

[add_task\(\)](#) (airflow.models.DAG method), 111
[add_tasks\(\)](#) (airflow.models.DAG method), 111
[airflow.contrib.hooks](#) (module), 122
[airflow.contrib.operators](#) (module), 101
[airflow.executors](#) (module), 126
[airflow.hooks](#) (module), 120
[airflow.macros](#) (module), 109
[airflow.macros.hive](#) (module), 109
[airflow.models](#) (module), 110
[airflow.operators](#) (module), 94
[allocate_ids\(\)](#) (airflow.contrib.hooks.datastore_hook.DatastoreHook method), 86
[are_dependencies_met\(\)](#) (airflow.models.TaskInstance method), 117
[are_dependents_done\(\)](#) (airflow.models.TaskInstance method), 117

B

[bag_dag\(\)](#) (airflow.models.DagBag method), 120
[BaseOperator](#) (class in airflow.models), 92, 113
[BaseSensorOperator](#) (class in airflow.operators.sensors), 94
[BashOperator](#) (class in airflow.operators), 94
[begin_transaction\(\)](#) (airflow.contrib.hooks.datastore_hook.DatastoreHook method), 86
[BigQueryCheckOperator](#) (class in airflow.contrib.operators.bigquery_check_operator), 76
[BigQueryHook](#) (class in airflow.contrib.hooks), 122
[BigQueryHook](#) (class in airflow.contrib.hooks.bigquery_hook), 80
[BigQueryIntervalCheckOperator](#) (class in airflow.contrib.operators.bigquery_check_operator), 77
[BigQueryOperator](#) (class in airflow.contrib.operators.bigquery_operator), 78, 101

[BigQueryToBigQueryOperator](#) (class in airflow.contrib.operators.bigquery_to_bigquery), 79
[BigQueryToCloudStorageOperator](#) (class in airflow.contrib.operators.bigquery_to_gcs), 80, 102
[BigQueryValueCheckOperator](#) (class in airflow.contrib.operators.bigquery_check_operator), 77
[BranchPythonOperator](#) (class in airflow.operators), 94
[bulk_dump\(\)](#) (airflow.hooks.DbApiHook method), 120
[bulk_load\(\)](#) (airflow.hooks.DbApiHook method), 121

C

[check_output\(\)](#) (airflow.contrib.hooks.SSHHook method), 126
[clear\(\)](#) (airflow.models.BaseOperator method), 115
[clear\(\)](#) (airflow.models.DAG method), 111
[clear_xcom_data\(\)](#) (airflow.models.TaskInstance method), 117
[cli\(\)](#) (airflow.models.DAG method), 111
[close_conn\(\)](#) (airflow.contrib.hooks.FTPHook method), 124
[closest_ds_partition\(\)](#) (in module airflow.macros.hive), 109
[collect_dags\(\)](#) (airflow.models.DagBag method), 120
[command\(\)](#) (airflow.models.TaskInstance method), 117
[command_as_list\(\)](#) (airflow.models.TaskInstance method), 117
[commit\(\)](#) (airflow.contrib.hooks.datastore_hook.DatastoreHook method), 86
[concurrency_reached](#) (airflow.models.DAG attribute), 111
[Connection](#) (class in airflow.models), 120
[create_dagrun\(\)](#) (airflow.models.DAG method), 112
[create_directory\(\)](#) (airflow.contrib.hooks.FTPHook method), 124
[current_state\(\)](#) (airflow.models.TaskInstance method), 117

D

dag (airflow.models.BaseOperator attribute), 115

DAG (class in airflow.models), 110

DagBag (class in airflow.models), 119

dagbag_report() (airflow.models.DagBag method), 120

DatabricksSubmitRunOperator (class in airflow.contrib.operators.databricks_operator), 74, 103

DataFlowHook (class in airflow.contrib.hooks.gcp_dataflow_hook), 82

DataFlowJavaOperator (class in airflow.contrib.operators.dataflow_operator), 81

DataProcHadoopOperator (class in airflow.contrib.operators.dataproc_operator), 85

DataProcHiveOperator (class in airflow.contrib.operators.dataproc_operator), 84

DataProcPigOperator (class in airflow.contrib.operators.dataproc_operator), 83

DataProcPySparkOperator (class in airflow.contrib.operators.dataproc_operator), 86

DataProcSparkOperator (class in airflow.contrib.operators.dataproc_operator), 85

DataProcSparkSqlOperator (class in airflow.contrib.operators.dataproc_operator), 84

DatastoreHook (class in airflow.contrib.hooks.datastore_hook), 86

DbApiHook (class in airflow.hooks), 120

deactivate_stale_dags() (airflow.models.DAG static method), 112

deactivate_unknown_dags() (airflow.models.DAG static method), 112

delete() (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), 88

delete() (airflow.contrib.hooks.GoogleCloudStorageHook method), 123

delete_directory() (airflow.contrib.hooks.FTPHook method), 124

delete_file() (airflow.contrib.hooks.FTPHook method), 124

deps (airflow.models.BaseOperator attribute), 115

describe_directory() (airflow.contrib.hooks.FTPHook method), 124

detect_downstream_cycle() (airflow.models.BaseOperator method), 115

DockerOperator (class in airflow.operators.docker_operator), 100

download() (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), 88

download() (airflow.contrib.hooks.GoogleCloudStorageHook method), 123

downstream_list (airflow.models.BaseOperator attribute), 115

ds_add() (in module airflow.macros), 109

ds_format() (in module airflow.macros), 109

DummyOperator (class in airflow.operators), 95

E

ECSOperator (class in airflow.contrib.operators.ecs_operator), 105

EmailOperator (class in airflow.operators), 95

error() (airflow.models.TaskInstance method), 117

execute() (airflow.models.BaseOperator method), 115

execute() (airflow.operators.BashOperator method), 94

exists() (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), 89

exists() (airflow.contrib.hooks.GoogleCloudStorageHook method), 123

ExternalTaskSensor (class in airflow.operators), 95

extra_dejson (airflow.models.Connection attribute), 120

F

filepath (airflow.models.DAG attribute), 112

filter_for_filesize() (airflow.operators.HdfsSensor static method), 96

filter_for_ignored_ext() (airflow.operators.HdfsSensor static method), 96

folder (airflow.models.DAG attribute), 112

FTPHook (class in airflow.contrib.hooks), 124

G

generate_command() (airflow.models.TaskInstance static method), 117

GenericTransfer (class in airflow.operators), 96

get_active_runs() (airflow.models.DAG method), 112

get_conn() (airflow.contrib.hooks.bigquery_hook.BigQueryHook method), 80

get_conn() (airflow.contrib.hooks.BigQueryHook method), 122

get_conn() (airflow.contrib.hooks.datastore_hook.DatastoreHook method), 86

get_conn() (airflow.contrib.hooks.FTPHook method), 125

get_conn() (airflow.contrib.hooks.gcp_dataflow_hook.DataFlowHook method), 82

get_conn() (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), 89

get_conn() (airflow.contrib.hooks.GoogleCloudStorageHook method), 124

get_conn() (airflow.hooks.DbApiHook method), 121

get_conn() (airflow.hooks.HttpHook method), 122

get_conn() (airflow.hooks.SQLiteHook method), 122

[get_cursor\(\)](#) (airflow.hooks.DbApiHook method), [121](#)
[get_dag\(\)](#) (airflow.models.DagBag method), [120](#)
[get_dagrun\(\)](#) (airflow.models.DAG method), [112](#)
[get_dagrun\(\)](#) (airflow.models.TaskInstance method), [118](#)
[get_direct_relatives\(\)](#) (airflow.models.BaseOperator method), [116](#)
[get_first\(\)](#) (airflow.hooks.DbApiHook method), [121](#)
[get_flat_relatives\(\)](#) (airflow.models.BaseOperator method), [116](#)
[get_last_dagrun\(\)](#) (airflow.models.DAG method), [112](#)
[get_num_task_instances\(\)](#) (airflow.models.DAG static method), [112](#)
[get_pandas_df\(\)](#) (airflow.contrib.hooks.bigquery_hook.BigQueryHook method), [81](#)
[get_pandas_df\(\)](#) (airflow.contrib.hooks.BigQueryHook method), [122](#)
[get_pandas_df\(\)](#) (airflow.hooks.DbApiHook method), [121](#)
[get_records\(\)](#) (airflow.hooks.DbApiHook method), [121](#)
[get_service\(\)](#) (airflow.contrib.hooks.bigquery_hook.BigQueryHook method), [81](#)
[get_service\(\)](#) (airflow.contrib.hooks.BigQueryHook method), [123](#)
[get_task_instances\(\)](#) (airflow.models.BaseOperator method), [116](#)
[get_template_env\(\)](#) (airflow.models.DAG method), [113](#)
[GoogleCloudStorageDownloadOperator](#) (class in airflow.contrib.operators.gcs_download_operator), [87](#), [106](#)
[GoogleCloudStorageHook](#) (class in airflow.contrib.hooks), [123](#)
[GoogleCloudStorageHook](#) (class in airflow.contrib.hooks.gcs_hook), [88](#), [126](#)
[GoogleCloudStorageToBigQueryOperator](#) (class in airflow.contrib.operators.gcs_to_bq), [88](#)

H

[has_dag\(\)](#) (airflow.models.BaseOperator method), [116](#)
[HdfsSensor](#) (class in airflow.operators), [96](#)
[HipChatAPIOperator](#) (class in airflow.contrib.operators.hipchat_operator), [106](#)
[HipChatAPISendRoomNotificationOperator](#) (class in airflow.contrib.operators.hipchat_operator), [106](#)
[HivePartitionSensor](#) (class in airflow.operators), [96](#)
[HttpHook](#) (class in airflow.hooks), [122](#)
[HttpSensor](#) (class in airflow.operators), [97](#)

I

[init_on_load\(\)](#) (airflow.models.TaskInstance method), [118](#)
[insert_rows\(\)](#) (airflow.contrib.hooks.bigquery_hook.BigQueryHook method), [81](#)

[insert_rows\(\)](#) (airflow.contrib.hooks.BigQueryHook method), [123](#)
[insert_rows\(\)](#) (airflow.hooks.DbApiHook method), [121](#)
[is_paused](#) (airflow.models.DAG attribute), [113](#)
[is_premature](#) (airflow.models.TaskInstance attribute), [118](#)
[is_updated_after\(\)](#) (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), [89](#)
[is_updated_after\(\)](#) (airflow.contrib.hooks.GoogleCloudStorageHook method), [124](#)

K

[key](#) (airflow.models.TaskInstance attribute), [118](#)
[keyframes\(\)](#) (airflow.models.DagBag method), [120](#)

L

[latest_execution_date](#) (airflow.models.DAG attribute), [113](#)
[list\(\)](#) (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), [89](#)
[list\(\)](#) (airflow.contrib.hooks.GoogleCloudStorageHook method), [124](#)
[list_directory\(\)](#) (airflow.contrib.hooks.FTPHook method), [125](#)
[LocalExecutor](#) (class in airflow.executors), [126](#)
[lookup\(\)](#) (airflow.contrib.hooks.datastore_hook.DatastoreHook method), [86](#)

M

[max_partition\(\)](#) (in module airflow.macros.hive), [110](#)
[MesosExecutor](#) (class in airflow.contrib.executors.mesos_executor), [127](#)
[MetastorePartitionSensor](#) (class in airflow.operators), [98](#)

N

[NamedHivePartitionSensor](#) (class in airflow.operators), [98](#)
[next_retry_datetime\(\)](#) (airflow.models.TaskInstance method), [118](#)
[normalize_schedule\(\)](#) (airflow.models.DAG method), [113](#)

O

[on_kill\(\)](#) (airflow.models.BaseOperator method), [116](#)

P

[pool_full\(\)](#) (airflow.models.TaskInstance method), [118](#)
[Popen\(\)](#) (airflow.contrib.hooks.SSHHook method), [126](#)
[post_execute\(\)](#) (airflow.models.BaseOperator method), [116](#)
[pre_execute\(\)](#) (airflow.models.BaseOperator method), [116](#)
[prepare_template\(\)](#) (airflow.models.BaseOperator method), [116](#)

previous_ti (airflow.models.TaskInstance attribute), 118
 process_file() (airflow.models.DagBag method), 120
 PythonOperator (class in airflow.operators), 98

R

random() (in module airflow.macros), 109
 ready_for_retry() (airflow.models.TaskInstance method), 118
 refresh_from_db() (airflow.models.TaskInstance method), 118
 rename() (airflow.contrib.hooks.FTPHook method), 125
 render_template() (airflow.models.BaseOperator method), 116
 render_template_from_field() (airflow.models.BaseOperator method), 116
 retrieve_file() (airflow.contrib.hooks.FTPHook method), 125
 rollback() (airflow.contrib.hooks.datastore_hook.DatastoreHook method), 87
 run() (airflow.hooks.DbApiHook method), 122
 run() (airflow.hooks.HttpHook method), 122
 run() (airflow.models.BaseOperator method), 116
 run() (airflow.models.DAG method), 113
 run() (airflow.models.TaskInstance method), 118
 run_and_check() (airflow.hooks.HttpHook method), 122
 run_query() (airflow.contrib.hooks.datastore_hook.DatastoreHook method), 87

S

S3KeySensor (class in airflow.operators), 99
 schedule_interval (airflow.models.BaseOperator attribute), 116
 SequentialExecutor (class in airflow.executors), 126
 set_dependency() (airflow.models.DAG method), 113
 set_downstream() (airflow.models.BaseOperator method), 116
 set_upstream() (airflow.models.BaseOperator method), 116
 ShortCircuitOperator (class in airflow.operators), 99
 SimpleHttpOperator (class in airflow.operators), 97
 size() (airflow.models.DagBag method), 120
 SqliteHook (class in airflow.hooks), 122
 SqlSensor (class in airflow.operators), 99
 SSHExecuteOperator (class in airflow.contrib.operators), 101
 SSHHook (class in airflow.contrib.hooks), 125
 store_file() (airflow.contrib.hooks.FTPHook method), 125
 sub_dag() (airflow.models.DAG method), 113
 subdags (airflow.models.DAG attribute), 113
 sync_to_db() (airflow.models.DAG static method), 113

T

table_exists() (airflow.contrib.hooks.bigquery_hook.BigQueryHook

method), 81
 table_exists() (airflow.contrib.hooks.BigQueryHook method), 123
 TaskInstance (class in airflow.models), 116
 TimeSensor (class in airflow.operators), 100
 topological_sort() (airflow.models.DAG method), 113
 tree_view() (airflow.models.DAG method), 113
 TriggerDagRunOperator (class in airflow.operators), 95
 tunnel() (airflow.contrib.hooks.SSHHook method), 126

U

upload() (airflow.contrib.hooks.gcs_hook.GoogleCloudStorageHook method), 89
 upload() (airflow.contrib.hooks.GoogleCloudStorageHook method), 124
 upstream_list (airflow.models.BaseOperator attribute), 116

W

WebHdfsSensor (class in airflow.operators), 100

X

xcom_pull() (airflow.models.BaseOperator method), 116
 xcom_pull() (airflow.models.TaskInstance method), 119
 xcom_push() (airflow.models.BaseOperator method), 116
 xcom_push() (airflow.models.TaskInstance method), 119