

Automated Testing and Deployment of DAGs

Austin Bennett

Automation and Testing Foundations

For Python/AirFlow

Best Practices: Foundations

There are some “Best Practices”:

<https://airflow.apache.org/docs/apache-airflow/stable/best-practices.html>

But, let's talk about the FOUNDATIONS/BASICS

P.S. I take for Git use as a given; though, I think much of the talk will be applicable even if not(?)

Code Hygiene

- Might not seem fun, or even a distraction ...
- BUT ...



StaleBot

<https://github.com/actions/stale>

```
1  name: Mark and close stale pull requests
2
3  on:
4    schedule:
5      - cron: '00 01 * * *'
6
7  env:
8    PR_STALE_DAYS: 14
9    PR_CLOSE_DAYS: 14
10
11 jobs:
12   stale:
13     runs-on: ubuntu-latest
14     permissions:
15       #contents: write -- ADD this if/once allowing branches to be deleted [ also need to add
16       issues: read
17       pull-requests: write
18     steps:
19       - uses: actions/stale@v8
20         with:
21           repo-token: ${ secrets.GITHUB_TOKEN }
22           stale-pr-message: "This pull request has been marked as stale. It will be closed in ${
23           close-pr-message: "This pull request has been closed due to lack of activity."
24           days-before-pr-stale: ${ env.PR_STALE_DAYS }
25           days-before-pr-close: ${ env.PR_CLOSE_DAYS }
```

Consistency, Commonalities



Where to begin?

- We are talking airflow, so esp. Python code in:
 - `$AIRFLOW_HOME/plugins`
 - `$AIRFLOW_HOME/dags`

Some Tools

- CONSISTENT ENVIRONMENTS
 - Poetry
 - Nix
 - Devcontainers
- pre-commit
- GH Actions
- PYTHON
 - Black
 - Ruff / Flake8
 - Type-checking(?!)
- SQL
 - SQLFluff
-

CONSISTENCY Across Team

Define and Pin Versions/Dependencies to save debug headaches!!



<https://python-poetry.org/>

At least a requirements.txt

Poetry is great if sticking to ONLY Python

BUT, for environments not tied to python, see:

- NIX
- DevContainers

Wider than Python Ecosystem: NIX

<https://nixos.org/>

Pros/Cons.

Cool Design.

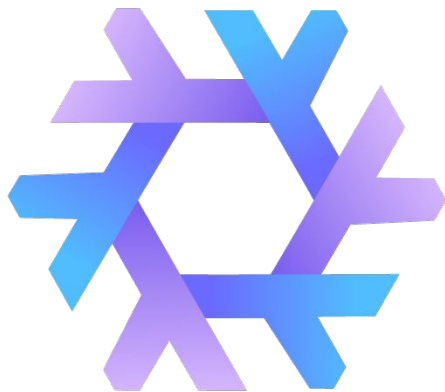
Works.

Not Common.

Plays well with poetry:

<https://github.com/nix-community/poetry2nix>

Tutorial: <https://determinate.systems/posts/zero-to-nix/>



NixOS

DevContainers

<https://containers.dev/>

Can run locally while developing, and beyond

Well supported by some major IDEs

Can be 'just' Dockerfiles



Pre-Commit

Run things locally “before” the commit

Official supported hooks:

- <https://pre-commit.com/hooks.html>

Also, can just be a ‘hook’/trigger based on the action of attempting to commit

Airflow Repo uses Pre-Commit:

- <https://github.com/apache/airflow/blob/main/.pre-commit-config.yaml>

^^ is currently 1346 lines on master



Pre-Commit

<https://pre-commit.com/>

Pre-Commit: Starter

```
1 ▼ repos:
2 ▼   - repo: https://github.com/pre-commit/pre-commit-hooks
3     rev: v4.4.0
4 ▼   hooks:
5     - id: check-yaml
6     - id: end-of-file-fixer
7     - id: trailing-whitespace
8     - id: check-toml
9     - id: check-json
10
```

A good starting place for Python



- Black is the uncompromising Python code formatter. By using it, you agree to cede control over minutiae of hand-formatting. In return, Black gives you speed, determinism, and freedom from pycodestyle nagging about formatting. You will save time and mental energy for more important matters.

Pre-Commit: Black ["custom"]

```
1  repos:
2  |   - repo: local
3  |     hooks:
4  |       - id: black
5  |         name: black
6  |         language: system
7  |         entry: poetry run black --check
8  |         types: [python]
```

GitHub Action: Black [off-the-shelf]

```
1  name: Lint
2
3  on: [push, pull_request]
4
5  jobs:
6    lint:
7      runs-on: ubuntu-latest
8      steps:
9        - uses: actions/checkout@v4
10       - uses: psf/black@stable
11
```


GitHub Action: Black [“custom”]

This can even be improved

Ex: find changed files

- jitterbit/get-changed-files

Or custom code ...

Then only check the ‘new’/updated

Again, pros/cons

```
1  name: Black
2
3  on: [pull_request]
4
5  jobs:
6    black:
7      runs-on: ubuntu-latest
8      container:
9        image: ${{ vars.CONTAINER }}
10       options: --platform linux/amd64
11       credentials:
12         username: ${{ github.actor }}
13         password: ${{ secrets.github_token }}
14
15     steps:
16       - uses: actions/checkout@v4
17
18         # included to aid debuggin
19       - name: Echo value for container build
20         run: echo "${{ vars.CONTAINER }}"
21
22       - name: Run Black
23         shell: 'bash'
24         run: poetry run black --check .
25
```

AIRFLOW RULES

- Ex: <https://github.com/BasPH/pylint-airflow> [needs updated]
- Also: <https://github.com/feluella/airflint> [says not production ready]

The current codes are:

Code	Symbol	Description
C8300	different-operator-varname-taskid	For consistency assign the same variable name and task_id to operators.
C8301	match-callable-taskid	For consistency name the callable function '_[task_id]', e.g. PythonOperator(task_id='mytask', python_callable=_mytask).
C8302	mixed-dependency-directions	For consistency don't mix directions in a single statement, instead split over multiple statements.
C8303	task-no-dependencies	Sometimes a task without any dependency is desired, however often it is the result of a forgotten dependency.
C8304	task-context-argname	Indicate you expect Airflow task context variables in the **kwargs argument by renaming to **context.
C8305	task-context-separate-arg	To avoid unpacking kwargs from the Airflow task context in a function, you can set the needed variables as arguments in the function.
C8306	match-dagid-filename	For consistency match the DAG filename with the dag_id.
R8300	unused-xcom	Return values from a python_callable function or execute() method are automatically pushed as XCom.
W8300	basehook-top-level	Airflow executes DAG scripts periodically and anything at the top level of a script is executed. Therefore, move BaseHook calls into functions/hooks/operators.
E8300	duplicate-dag-name	DAG name should be unique.
E8301	duplicate-task-name	Task name within a DAG should be unique.
E8302	duplicate-dependency	Task dependencies can be defined only once.
E8303	dag-with-cycles	A DAG is acyclic and cannot contain cycles.
E8304	task-no-dag	A task must know a DAG instance to run.

AIRFLOW RULES: Ruff!

See: <https://github.com/astral-sh/ruff/issues/4421>

Currently 'just' one rule:

- “task variable name should be same as task_id”
 - <https://github.com/astral-sh/ruff/pull/4687>

Room for more. Community is open for contributions!

RUFF

Has been in use for awhile

Charlie Marsh @charliermarsh

Apache Airflow adopts the Ruff formatter! 🚀🚀🚀

Switch from Black to Ruff formatter #35287

kanal merged 10 commits into apache/airflow from julianlaneve/ruff-formatter 34 minutes ago

Conversation 26 · Comments 10 · Checks 96 · Files changed 82

Juliane commented yesterday · edited · Contributor

This PR switches the formatter we use from Black to Ruff, now that Ruff's introduced a `formatter`. This PR also upgrades the version of Ruff to the latest.

Running locally:

```
> time black --executable"/usr/bin/black" --target-version["py310"] --config"/dev/null" --diff
All done!
All files up to date.
black - 0.85s user 1.65s system 9% cpu 6.82s total

> time ruff format --executable"/usr/bin/ruff" --target-version["py310"] --config"/dev/null" --diff
2889 files left unchanged
ruff format - 0.85s user 0.18s system 89% cpu 0.13s total
```

This represents means Ruff is ~40x faster than Black (when there are no changes - when there are changes, I've seen it even higher).

Note that this doesn't swap the formatter used in inline Python code in the docs as I haven't seen an easy way of getting the Ruff formatter working with it. Because of this, a lot of the Black code hasn't been removed (i.e. the code to upgrade black, etc).

* Add meaningful description above

Julian LaNeve @JulianLaneve · Oct 31, 2023

Just migrated Apache Airflow's ~1m lines of code from Black to Ruff's new formatter... it's about 40x faster and took less than an hour. Incredible work that @charliermarsh & the team at @astral_sh are doing!

<https://x.com/charliermarsh/status/1719496146815422536>

Ruff: Pre-Commit

```
- repo: https://github.com/astral-sh/ruff-pre-commit
# Ruff version.
rev: v0.6.4
hooks:
  # Run the linter.
  - id: ruff
    args: [ --fix ]
  # Run the formatter.
  - id: ruff-format
```

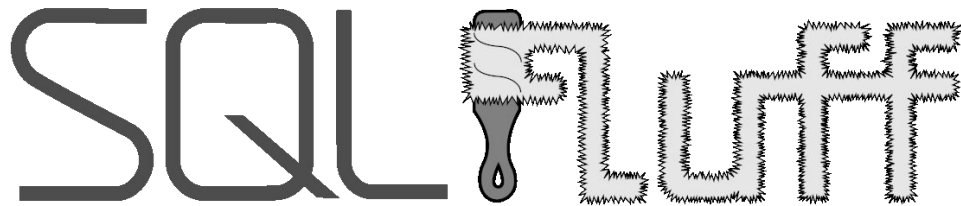
<https://github.com/astral-sh/ruff-pre-commit>

Ruff GitHub Action

```
name: Ruff
on: [push, pull_request]
jobs:
  ruff:
    runs-on: ubuntu-latest
    steps:
      - uses: actions/checkout@v4
      - uses: chartboost/ruff-action@v1
```

SQLFluff

DAGs orchestrate ALOT of SQL



 The SQL Linter for humans.

Lint

SQLFluff finds issues with your SQL code and reports them back to you (and your team) automatically so that your code reviews can be more about function and less about form.

Fix

SQLFluff saves time by fixing linting issues found in your code to save you time, and make it easy to have consistent and legible SQL.

Parse

SQLFluff parses your SQL to catch a range of syntax issues without needing access to the database, so you can catch mistakes earlier in your development process.

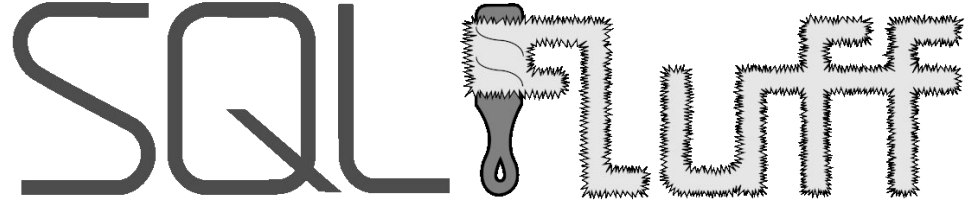
Configure

SQLFluff is configurable to work with a range of SQL dialects and style choices. It has opinionated defaults, so you can get going easily, but a range of flexible configuration options to fit your local style.

SQLFluff

DIALECTS:

ANSI
Athena
BigQuery
ClickHouse
Databricks
Db2
DuckDB
Exasol
Greenplum
Hive
Materialize
MariaDB
MySQL
Oracle
PostgreSQL
Redshift
Snowflake
SOQL
SparkSQL
SQLite
T-SQL
Teradata
Trino
Vertica



The SQL Linter for humans.

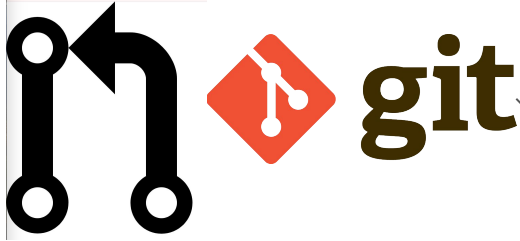
See Contribution guide if need more than this →

<https://github.com/sqlfluff/sqlfluff/wiki/Contributing-Dialect-Changes>

Templating [env/project] variables

<https://docs.sqlfluff.com/en/2.1.3/developingplugins.html>

Manual Deployment



Manual LOG IN
To deploy !?!

`$AIRFLOW_HOME/dags`

Auto-Deploy

With Safeguards, once the code is OK, deploy it...

BUT ... “PROTECT THE REPO”

DEPLOY

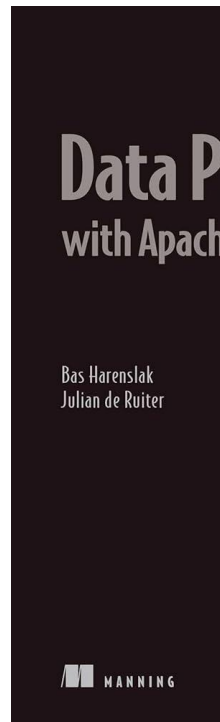
```
1 name: Deploy Airflow DAGs
2
3 on:
4   push:
5     branches:
6       - master
7
8   permissions:
9     id-token: write
10    contents: read
11
12  jobs:
13    airflow-dags:
14      runs-on: ubuntu-latest
15      steps:
16        - uses: actions/checkout@v3
17        - id: 'auth'
18          uses: 'google-github-actions/auth@v2'
19          with:
20            workload_identity_provider: ${vars.PROD_WIP}
21            service_account: ${vars.PROD_SA}
22        - name: 'Set up Cloud SDK'
23          uses: 'google-github-actions/setup-gcloud@v2'
24        - name: 'sync DAG files'
25          run: 'gsutil -m rsync -d -r $REPO_PATH/dags $AF_ENV/dags'
26
```

Airflow Rules!

9 Testing

This chapter covers

- Testing Airflow tasks in a CI/CD pipeline
- Structuring a project for testing with pytest
- Mimicking a DAG run to test tasks that apply templating
- Faking external system events with mocking
- Testing behavior in external systems with containers



MOR

- ActionLint :-p
 - Lint yor GH Actions
- Types ... MyPy, PyRight, PyType
 - !!
- PyTest and more
- Commit/extend RUFF!
- Don't forget to read/implement the `Data Pipelines with Airflow` book
- With all of these you'll be in solid shape