

# Projekt 1 – orkiestracja – Apache Airflow

Apache Airflow to platforma open source do opracowywania, harmonogramowania i monitorowania wsadowych przepływów pracy. Funkcjonalność, którą ona oferuje, często określa się mianem orkiestracja.

Niektórzy uważają, że jest to jedno z najlepszych narzędzi w swojej klasie.

Apache Airflow jest oparte na Pythonie. Prawie w każdym elemencie systemu wykorzystywany jest właśnie ten język programowania. Począwszy od sposobu definiowania przepływów pracy przez użytkownika, a skończywszy na implementacji wtyczek. Dzięki wykorzystywaniu tak uniwersalnego języka, Apache Airflow pozwala na tworzenie przepływów, które mogą wykorzystywać praktycznie dowolną technologię.

Istotnym elementem systemu jest jego interfejs sieciowy. Pozwala on nie tylko podglądać zdefiniowane przepływy, oraz je uruchamiać, ale także w pełnym zakresie zarządzać stanem przepływów.

Airflow można wdrożyć na wiele sposobów, od pojedynczego procesu na laptopie po konfigurację rozproszoną, która obsługuje potężne przepływy. My wykorzystamy bardzo prostą konfigurację uruchomioną na maszynie master klastra *Dataproc*.

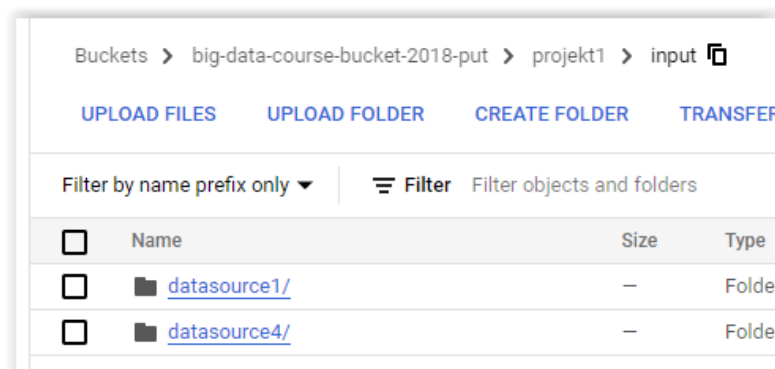
Naszym celem ostatecznym będzie przygotowanie przepływu, który będzie zarządzał uruchamianiem (orkiestrował) elementów składowych projektu 1.

## Elementy wyjściowe

W założeniu wszystkie elementy składowe dla projektu 1 są już gotowe. W przypadku niniejszego warsztatu będzie to wariant: *Hadoop Streaming*, *Pig*, *Zestaw danych nr 3*

Powyższe oznacza, że są już gotowe:

- Program MapReduce:
  - skrypt mappera: `mapper.py`
  - skrypt reduktora: `reducer.py`
  - skrypt combinera: `combiner.py`
- Skrypt Piga
  - skrypt `transform5.pig`
- Dane wejściowe
  - znajdują się one w zasobniku w katalogu `gs://nazwa_zasobnika/projekt1/input` (w podkatalogach `datasource1` i `datasource4`)



Ponadto jest gotowa pobrana źródłowa definicja przepływu, plik projekt1.py

Oczywiście każdy projekt może dotyczyć innego wariantu, a zatem elementy składowe projektu mogą być inne. Należy to uwzględnić realizując ten warsztat. W razie potrzeby szukaj wsparcia u prowadzącego.

Wszystkie elementy składowe powinny być przetestowane i działające. Polecenia uruchamiające powinny być opracowane i znane. Skrypt Pig/Hive powinien być parametryzowany:

- input\_dir3 – katalogiem wejściowym dla przetworzonego pierwszego zbioru danych (wynik przetwarzania MapReduce)
- input\_dir4 – katalogiem wejściowym dla drugiego zbioru danych
- output\_dir6 – katalogiem wyjściowym, który będzie zawierał ostateczny wynik całości przetwarzania

W razie potrzeby wprowadź odpowiednie poprawki oraz przetestuj jeszcze raz wszystkie składowe projektu 1

## Uruchomienie klastra Dataproc

*Apache Airflow*, nie należy do komponentów dostępnych w ramach klastra *Dataproc*, o czym możemy się przekonać: <https://cloud.google.com/dataproc/docs/concepts/versioning/dataproc-release-2.0>

Oznacza to, że będziemy go musieli sobie zainstalować niezależnie na uruchomionym uprzednio klastrze.

Tak na marginesie, w ramach usług dostarczanych przez platformę GCP istnieje narzędzie *Cloud Composer*. Jest to usługa pozwalająca na orkiestrację przetwarzania (danych). Opiera się ona całkowicie na *Apache Airflow*. Nie będziemy jednak z niej korzystali w ramach tego warsztatu.

1. Utwórz klaster *Dataproc* za pomocą poniższego polecenia

```
gcloud beta dataproc clusters create ${CLUSTER_NAME} \
--enable-component-gateway --bucket ${BUCKET_NAME} \
--region ${REGION} --subnet default --zone ${ZONE} \
--master-machine-type n1-standard-4 --master-boot-disk-size 50 \
--num-workers 2 \
--worker-machine-type n1-standard-2 --worker-boot-disk-size 50 \
--image-version 2.0-debian10 \
--optional-components DOCKER \
--project ${PROJECT_ID} --max-age=3h
```

2. Uruchom terminal SSH do maszyny master naszego klastra. Sprawdź za pomocą poniższego polecenia wersję zainstalowanego Pythona

```
jankiewicz_krzysztof@cluster-0999-m:~$ python --version
Python 3.8.13
```

```
export AIRFLOW_HOME=~/.airflow
```

```
pip install apache-airflow
```

```
export PATH=$PATH:~/local/bin
```

```
airflow db init
```

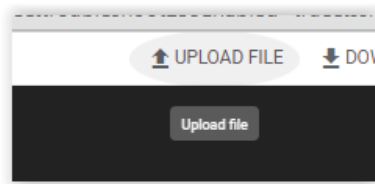
---

airflow standalone

```
triggerer | triggerer  
triggerer | _____|_____|_____  
triggerer | ____/_||_/___/_/_____\__/_/\_/  
triggerer | ____/_|\_/_/___/_/_____\__/_/\_/  
triggerer | _/_/_\___/_/___/_/_____\__/_/\_/  
triggerer | [2022-09-20 10:01:47,268] {triggerer_job.py:101} INFO - Starting the triggerer  
webserver | [2022-09-20 10:01:47 +0000] [8970] [INFO] Starting gunicorn 20.1.0  
webserver | [2022-09-20 10:01:48 +0000] [8970] [INFO] Listening at: http://0.0.0.0:8080 (8970)  
webserver | [2022-09-20 10:01:48 +0000] [8970] [INFO] Using worker: sync  
webserver | [2022-09-20 10:01:48 +0000] [8999] [INFO] Booting worker with pid: 8999  
webserver | [2022-09-20 10:01:48 +0000] [9000] [INFO] Booting worker with pid: 9000  
webserver | [2022-09-20 10:01:48 +0000] [9001] [INFO] Booting worker with pid: 9001  
webserver | [2022-09-20 10:01:48 +0000] [9002] [INFO] Booting worker with pid: 9002  
standalone |  
standalone | Airflow is ready  
standalone | Login with username: admin password: F5G5xSdZyuNnenWp  
standalone | Airflow Standalone is for development purposes only. Do not use this in production!
```

## Instalacja wzorcowej wersji przepływu

- Otwórz nowy terminal SSH do maszyny master naszego klastra.
- Załaduj do domowego katalogu wszystkie pliki Twojego projektu, a także plik z definicją przepływu



- Sprawdź czy wszystkie pliki znalazły się na serwerze master.

```
jankiewicz_krzysztof@cluster-0999-m:~$ ls
airflow combiner.py mapper.py projekt1.py reducer.py transform5.pig
```

- Utwórz w katalogu ~/airflow podkatalog dags, z którego będą załadowywane definicje przepływów. Przy okazji, utwórz w nim podkatalog project\_files, który będzie zawierał wszystkie składowe naszego projektu

```
mkdir -p ~/airflow/dags/project_files
```

- Przenieś plik z definicją przepływu do katalogu dags. Pozostałe pliki przenieś do katalogu project\_files

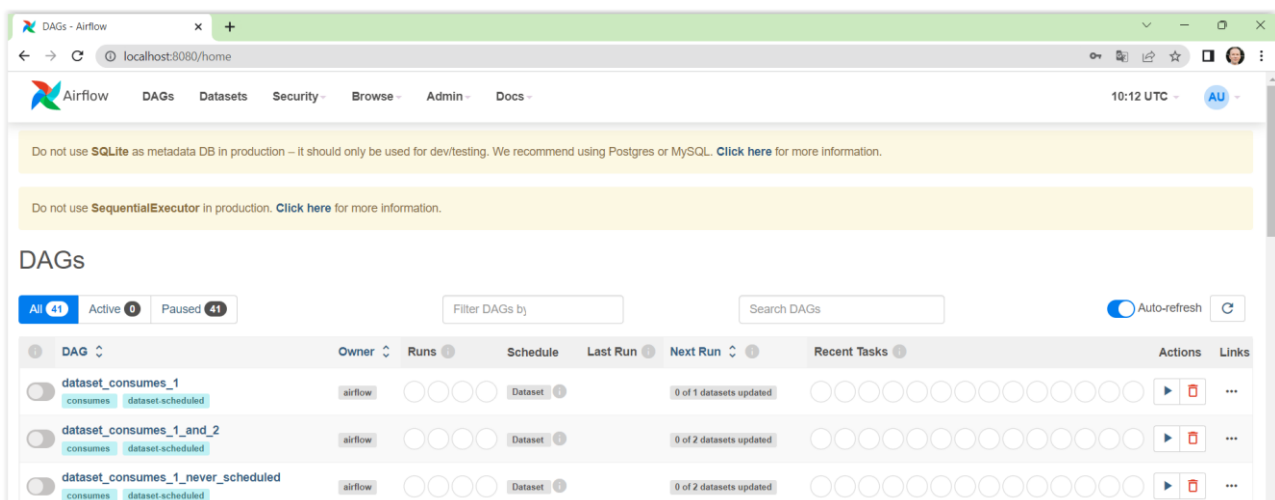
```
mv projekt1.py ~/airflow/dags/
mv *.* ~/airflow/dags/project_files
```

- Utwórz tunel do portu 8080 naszego serwera master. Sposób definicji tunelu został opisany na stronie kursu

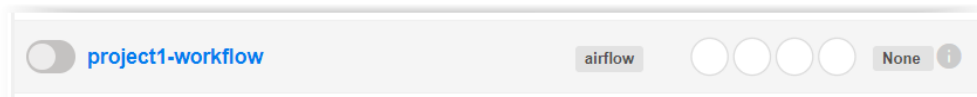


Konfiguracja środowiska zajęć - część 4 - udostępnienie połączeń dla zewnętrznego klienta SSH (PuTTY)

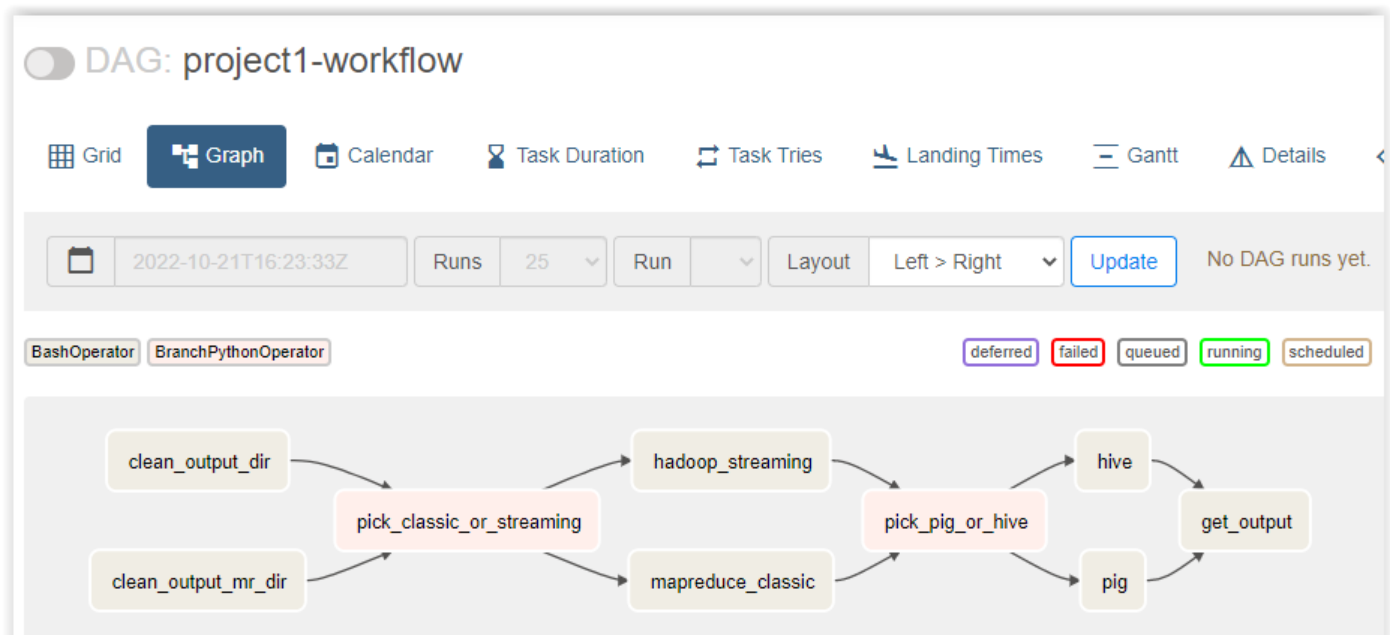
- Otwórz przeglądarkę pod adresem <http://localhost:8080>. Zaloguj się zgodnie z poznanymi danymi uwierzytelniającymi.



15. Znajdź wśród wielu przykładowych przepływów przepływ przeznaczony dla naszego projektu



16. Wybierz go, a następnie przyglądnij się jego graficznej reprezentacji



17. Zglądnij także do kodu tego przepływu i zapoznaj się z jego budową.

```

1 from airflow import DAG
2 from airflow.models.param import Param
3 from airflow.operators.bash import BashOperator
4 from airflow.operators.python import BranchPythonOperator
5 from datetime import datetime
6
7 with DAG(
8     "project1-workflow",
9     start_date=datetime(2015, 12, 1),
10    schedule_interval=None,
11    params={
12        "dags_home": Param("/home/username/airflow/dags", type="string"),
13        "input_dir": Param("/project1/input", type="string"),
14        "output_mr_dir": Param("/project1/output_mr3", type="string"),
15        "output_dir": Param("/project1/output6", type="string"),
16        "classic_or_streaming": Param("streaming", enum=["classic", "streaming"]),
17        "pig_or_hive": Param("pig", enum=["hive", "pig"]),
18    },
19    render_template_as_native_obj=True
20 ) as dag:
21     clean_output_mr_dir = BashOperator(
22         task_id="clean_output_mr_dir",
23         bash_command="""if $(hadoop fs -test -d {{ params.output_mr_dir }}) ; then hadoop fs -rm -f -r {{ params.output_mr_dir }}; fi""",
24     )
25     clean_output_dir = BashOperator(
26         task_id="clean_output_dir",
27         bash_command="""if $(hadoop fs -test -d {{ params.output_dir }}) ; then hadoop fs -rm -f -r {{ params.output_dir }}; fi""",
28     )
29
30     def _pick_classic_or_streaming():
31         if dag.params['classic_or_streaming'] == "classic":
32             return "mapreduce_classic"
33         else:
34             return "hadoop_streaming"
35
36     pick_classic_or_streaming = BranchPythonOperator(
37         task_id="pick_classic_or_streaming", python_callable=_pick_classic_or_streaming

```

18. Zwróć uwagę na:

- a. to, że jest on oparty o dwa typy operatorów
  - i. BashOperator – który uruchamia polecenia systemu operacyjnego, oraz
  - ii. BranchPythonOperator, stosowany do warunkowego uruchamiania kolejnych operatorów
- b. parametry oraz ich wartości domyślne,
- c. polecenia zawarte w operatorach BashOperator oraz wykorzystanie w nich wartości parametrów.

19. Aby nasz przepływ mógł zadziałać konieczne są poprawki. W terminalu SSH wywołaj edytor tekstowy

```
nano ~/airflow/dags/projekt1.py
```

20. A następnie:

- a. zmień domyślne parametry classic\_or\_streaming oraz pig\_or\_hive na zgodne z Twoim wariantem projektu. Dla przykładu:

```
"output_dir": Param("/project1/output6", type="string"),
"classic_or_streaming": Param("streaming", enum=["classic", "streaming"]),
"pig_or_hive": Param("pig", enum=["hive", "pig"]),
},
```

- b. nie zmieniaj domyślnych katalogów wejściowych oraz wyjściowych – mogłyby one zdradzić Twoje personalia
- c. popraw polecenia wywołujące elementy składowe Twojego projektu. Wykorzystaj te polecenia, które były przez Ciebie wykorzystane podczas testów. Uzupełnij je o odwołania do odpowiednich parametrów. W przypadku wariantu wykorzystywanego w tym warsztacie konieczna będzie modyfikacja:
  - i. polecenia w operatorze hadoop\_streaming. Przykładowe niekompletne polecenie wywołujące przetwarzanie *Hadoop Streaming*:

```
hadoop_streaming = BashOperator(
    task_id="hadoop_streaming",
    bash_command="""mapred streaming \
-files {{ params.dags_home }}/project_files/mapper.py,\
{{ params.dags_home }}/project_files/combiner.py,\
{{ params.dags_home }}/project_files/reducer.py \
-input {{ params.input_dir }}/datasource1 \
-mapper mapper.py \
-combiner combiner.py \
-reducer reducer.py \
-output {{ params.output_mr_dir }} \
```

- ii. polecenia w operatorze pig. Przykładowe (również niekompletnie) polecenie uruchamiające skrypt obsługiwany za pomocą *Piga*:

```
pig = BashOperator(
    task_id="pig",
    bash_command="""pig -f {{ params.dags_home }}/project_files/transform5.pig \
-param input_dir4={{ params.input_dir }}/datasource4 \
-param input_dir3={{ params.output_mr_dir }} \
-param output_dir6={{ params.output_dir }}""",
)
```

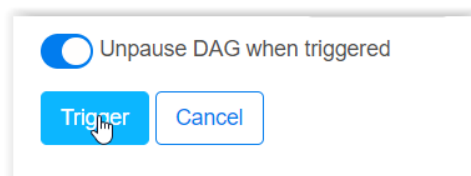
Pamiętaj, aby dokonać modyfikacji zgodnych z Twoim wariantem projektu i Twoimi składowymi. Powyższe fragmenty to tylko przykład dla wariantu wykorzystywanego w tym warsztacie.

21. Po zakończeniu wszystkich poprawek zapisz zmiany w pliku definicji przepływu.

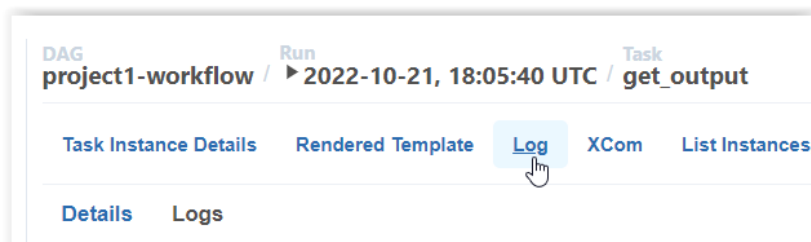
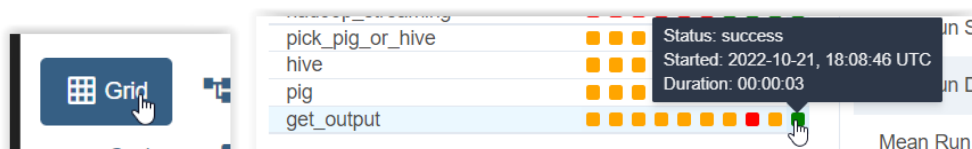
22. Powróć do interfejsu sieciowego Apache Airflow. Odśwież kod przepływu, sprawdź czy wszystkie wprowadzone zmiany są poprawne. Jeśli nie, dokonaj odpowiednich poprawek.
23. Jeśli Twój przepływ jest zgodny z oczekiwaniami. Czas na jego uruchomienie. Nasze domyślne ścieżki nie odpowiadają rzeczywistości, dlatego uruchomimy nasz przepływ w taki sposób, który pozwala nam skorygować domyślną konfigurację. Popraw odpowiednie ścieżki.



24. Uruchom przepływ za pomocą przycisku *Trigger*.



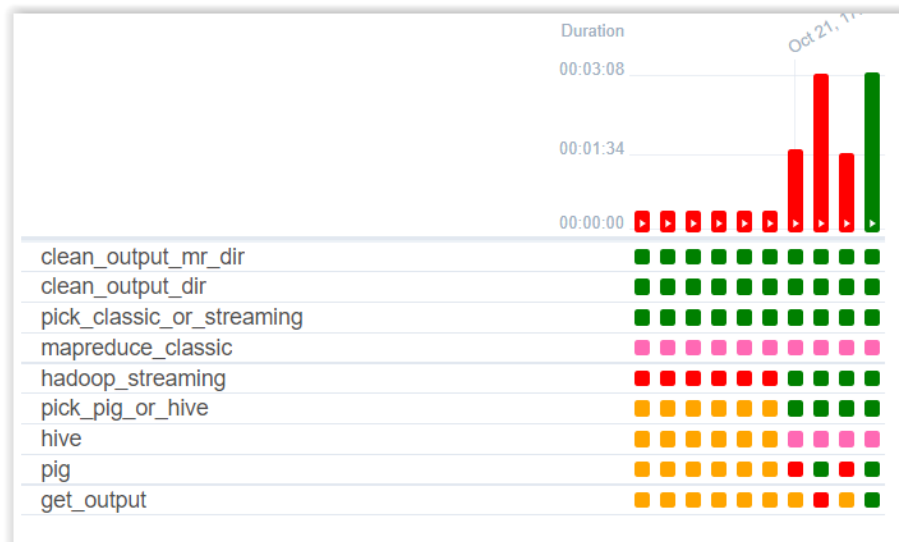
25. Poczekaj aż Twój przepływ zakończy swoje przetwarzanie. Podglądaj przebieg przepływu na zakładce *Grid*. Jeśli Twój przepływ za pierwszym zakończy się sukcesem... no cóż, przyjmij gratulacje. To zazwyczaj się nie zdarza. Zglądnij do logów ostatniego z zadań (get\_output) – znajdziesz w nim zaprezentowany końcowy wynik Twojego przetwarzania. Jeśli wszystko jest analogiczne jak podczas Twoich testów, mamy gotową całość. Wystarczy ją przygotować i zarejestrować jako rozwiązanie projektu.



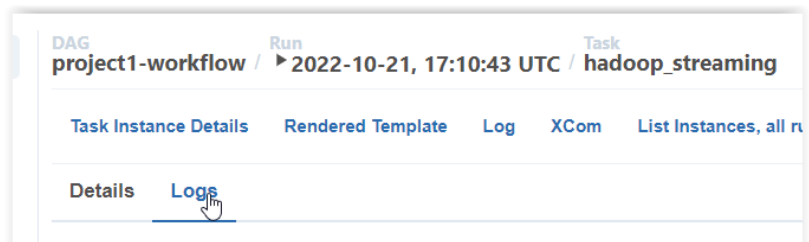
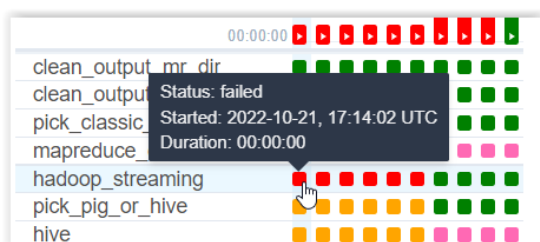
```

[2022-10-21, 18:08:47 UTC] {subprocess.py:75} INFO - Running command: ['/usr/bin/bash', '-c', 'h
[2022-10-21, 18:08:47 UTC] {subprocess.py:86} INFO - Output:
[2022-10-21, 18:08:49 UTC] {subprocess.py:93} INFO - street      person_type      killed injured
[2022-10-21, 18:08:49 UTC] {subprocess.py:93} INFO - {"fields":[{"name":"street","type":55,"desc
```

26. Jeśli jednak coś poszło nie tak (nie każdy jest nieomylny), przykładowo wielokrotnie poniżej...



...wówczas Apache Airflow jest nieoceniony. Każde z niepoprawnie wykonanych zadań można zaznaczyć, a następnie zaglądnąć do wygenerowanych logów zapisanych przez Apache Airflow. Z nich można wyczytać wszystko.



```

jprocess.py:75} INFO - Running command: ['/usr/bin/bash', '-c', 'mapred streaming -files mapper.py, combiner.py, red
jprocess.py:86} INFO - Output:
jprocess.py:93} INFO - WARNING: HADOOP_JOB_HISTORYSERVER_OPTS has been replaced by MAPRED_HISTORYSERVER_OPTS. Using
jprocess.py:93} INFO - Exception in thread "main" java.io.FileNotFoundException: File mapper.py does not exist
jprocess.py:93} INFO - at org.apache.hadoop.fs.RawLocalFileSystem.deprecatedGetFileStatus(RawLocalFileSystem.java
jprocess.py:93} INFO - at org.apache.hadoop.fs.RawLocalFileSystem.getFileLinkStatusInternal(RawLocalFileSystem.java
    
```



27. Napraw kolejno wszystkie błędy, doprowadzając swój przepływ do etapu, w którym będzie działał on bez zarzutu i powtarzalnie.

## Przygotowanie projektu do rejestracji

28. Jeśli Twój projekt jest gotowy, należy przygotować go do rejestracji. Przejdź do katalogu ~/airflow, a następnie spakuj katalog dags wraz z jego zawartością do pliku projekt1.zip

```

cd ~/airflow
rm -r dags/__pycache__
zip -r ~/projekt1.zip dags
    
```

29. Skopiuj tak utworzony plik do swojego zasobnika. Dzięki temu przeżyje on usunięcie klastra.

```

hadoop fs -copyFromLocal ~/projekt1.zip gs://nazwa_zasobnika/projekt1
    
```

30. Pobierz plik z zasobnika i zarejestruj go jako rozwiązanie swojego projektu. Osoba, która będzie go oceniała wystarczy, że rozpakuje Twój plik w katalogu ~/airflow i już będzie mogła go przetestować.