

# Airflow Zero to Hero

郭益華

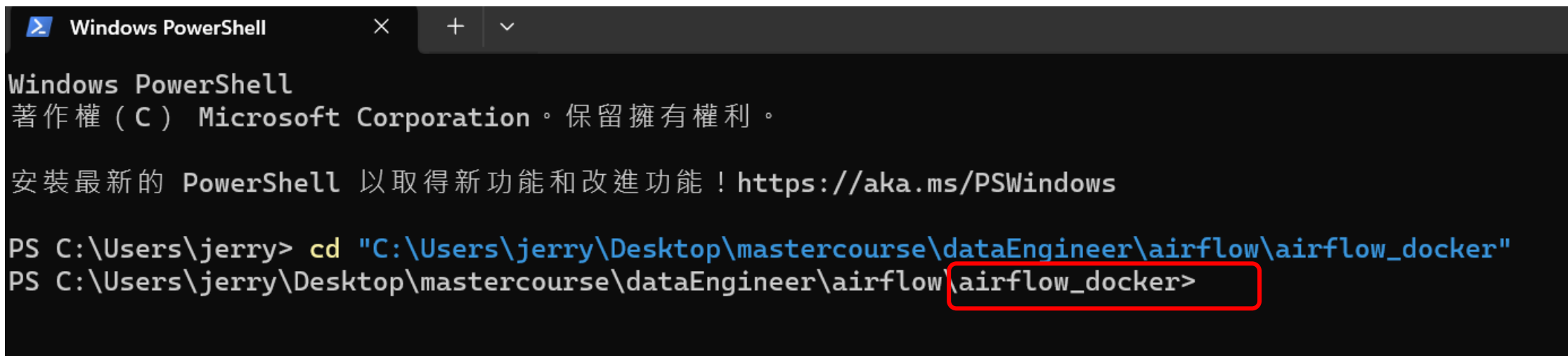
**GitHub**

# 目錄

- [Run Airflow in Docker](#)
- [Airflow DAG with Bash Operator](#)
- [Airflow DAG with Python Operator](#)
- [Data Sharing via Airflow Xcoms](#)
- [Airflow Task Flow API](#)
- [Airflow Catch-Up and Backfill](#)
- [Airflow Scheduler with Cron Expression](#)

# Run Airflow in Docker

# 建立一個 airflow\_docker 資料夾



```
Windows PowerShell
著作權 (C) Microsoft Corporation。保留擁有權利。

安裝最新的 PowerShell 以取得新功能和改進功能！https://aka.ms/PSWindows

PS C:\Users\jerry> cd "C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker"
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker>
```

# 搜尋 airflow docker compose



The image shows a Google search interface. The search bar contains the text "airflow docker compose" and is highlighted with a red rectangle. Below the search bar, there are tabs for "全部" (All), "圖片" (Images), "影片" (Videos), "購物" (Shopping), "書籍" (Books), and "更多" (More). The "全部" tab is selected. Below the tabs, it says "約有 1,340,000 項結果 (搜尋時間 : 0.28 秒)". A提示 (Tip) says: "提示：限制搜尋繁體中文的結果。瞭解詳情如何依語言過濾結果". Below the tip, there is a search result for "Apache Airflow" with the URL "https://airflow.apache.org › stable › howto". The result title is "Running Airflow in Docker — Airflow Documentation" and is highlighted with a red rectangle. Below the title, it says "The Docker Compose file uses the latest Airflow image (apache/airflow). If you need to install a new Python library or system library, you can customize and ...".

Google

airflow docker compose

全部 圖片 影片 購物 書籍 更多 工具

約有 1,340,000 項結果 (搜尋時間 : 0.28 秒)

提示：限制搜尋繁體中文的結果。瞭解詳情如何依語言過濾結果

Apache Airflow  
https://airflow.apache.org › stable › howto

Running Airflow in Docker — Airflow Documentation

The Docker Compose file uses the latest Airflow image (apache/airflow). If you need to install a new Python library or system library, you can customize and ...

# 複製此指令

Fetching `docker-compose.yaml`

To deploy Airflow on Docker Compose, you should fetch [docker-compose.yaml](https://airflow.apache.org/docs/apache-airflow/2.7.2/docker-compose.yaml).

```
curl -Lf0 'https://airflow.apache.org/docs/apache-airflow/2.7.2/docker-compose.yaml'
```

This file contains several service definitions:

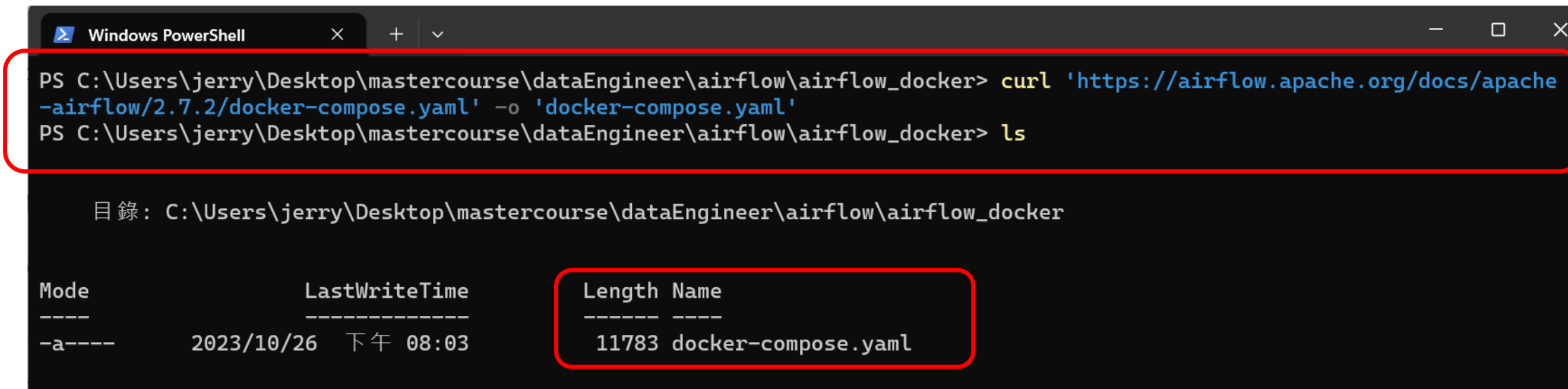
**Windows Powershell :**

```
curl 'https://airflow.apache.org/docs/apache-airflow/2.7.2/docker-compose.yaml' -o 'docker-compose.yaml'
```

**Linux:**

```
curl -Lf0 'https://airflow.apache.org/docs/apache-airflow/2.7.2/docker-compose.yaml'
```

# 下載



The screenshot shows a Windows PowerShell terminal window with the title bar 'Windows PowerShell'. The terminal content is as follows:

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> curl 'https://airflow.apache.org/docs/apache-airflow/2.7.2/docker-compose.yaml' -o 'docker-compose.yaml'
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> ls
```

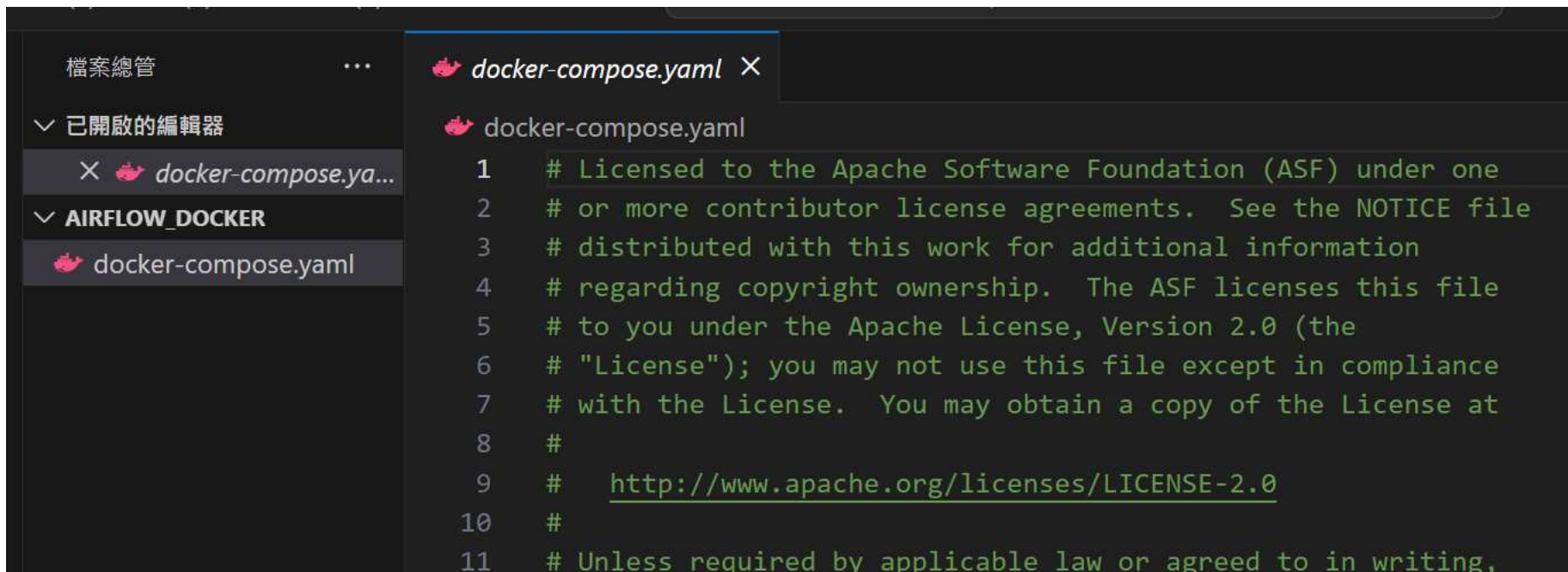
Below the commands, the directory path is shown:

目錄: C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow\_docker

A table listing the files in the directory is displayed below. The 'Length' and 'Name' columns are highlighted with a red box.

Mode	LastWriteTime	Length	Name
-a----	2023/10/26 下午 08:03	11783	docker-compose.yaml

# 查看docker-compose.yaml



The screenshot shows a code editor interface with a dark theme. On the left, a sidebar contains a file explorer. The top of the sidebar is labeled '檔案總管' (File Manager) with a three-dot menu icon. Below it, a section '已開啟的編輯器' (Opened Editors) lists two files: 'docker-compose.ya...' and 'AIRFLOW\_DOCKER'. Under 'AIRFLOW\_DOCKER', the file 'docker-compose.yaml' is selected. The main editor area displays the content of 'docker-compose.yaml', which is a license notice. The text is as follows:

```
1  # Licensed to the Apache Software Foundation (ASF) under one
2  # or more contributor license agreements.  See the NOTICE file
3  # distributed with this work for additional information
4  # regarding copyright ownership.  The ASF licenses this file
5  # to you under the Apache License, Version 2.0 (the
6  # "License"); you may not use this file except in compliance
7  # with the License.  You may obtain a copy of the License at
8  #
9  # http://www.apache.org/licenses/LICENSE-2.0
10 #
11 # Unless required by applicable law or agreed to in writing,
```



# 可看到有dag, log, plugins

```
75 volumes:
76   - ${AIRFLOW_PROJ_DIR:-.}/dags:/opt/airflow/dags
77   - ${AIRFLOW_PROJ_DIR:-.}/logs:/opt/airflow/logs
78   - ${AIRFLOW_PROJ_DIR:-.}/config:/opt/airflow/config
79   - ${AIRFLOW_PROJ_DIR:-.}/plugins:/opt/airflow/plugins
80 user: "${AIRFLOW_UID:-50000}:0"
```

# 在airflow\_docker中新增三個資料夾

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> mkdir dags
```

目錄: C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow\_docker

Mode	LastWriteTime	Length	Name
d-----	2023/10/26 下午 08:06		dags

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> mkdir logs
```

目錄: C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow\_docker

Mode	LastWriteTime	Length	Name
d-----	2023/10/26 下午 08:06		logs

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> mkdir plugins
```

目錄: C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow\_docker

# 新增成功

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> ls
```

目錄: C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow\_docker

Mode	LastWriteTime	Length	Name
d-----	2023/10/26 下午 08:06		dags
d-----	2023/10/26 下午 08:06		logs
d-----	2023/10/26 下午 08:06		plugins
-a-----	2023/10/26 下午 08:03	11783	docker-compose.yaml

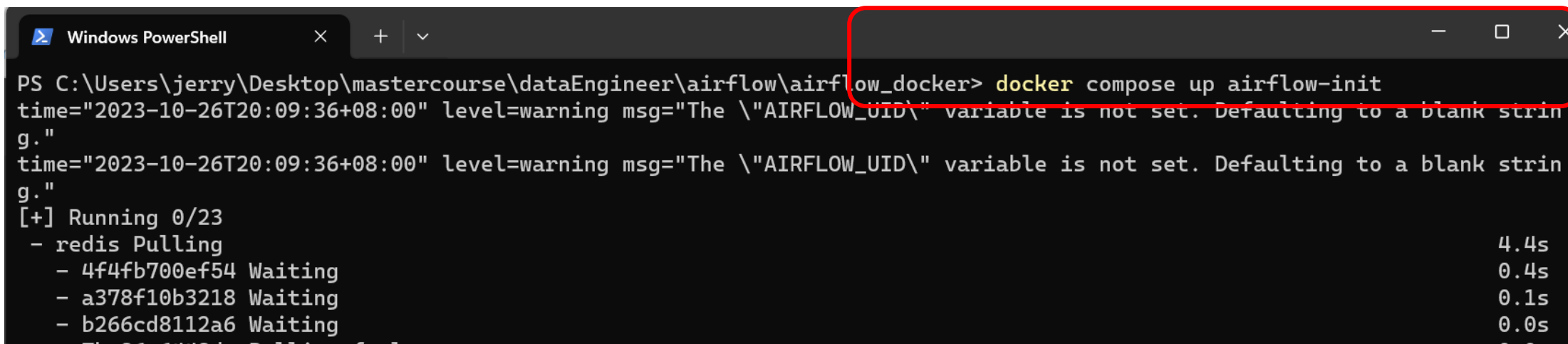
# 查看預設的airflow帳號密碼

```
252 # yamllint enable rule:line-length
253 environment:
254     <<: *airflow-common-env
255     _AIRFLOW_DB_MIGRATE: 'true'
256     _AIRFLOW_WWW_USER_CREATE: 'true'
257     _AIRFLOW_WWW_USER_USERNAME: ${_AIRFLOW_WWW_USER_USERNAME:-airflow}
258     _AIRFLOW_WWW_USER_PASSWORD: ${_AIRFLOW_WWW_USER_PASSWORD:-airflow}
259     _PIP_ADDITIONAL_REQUIREMENTS: ''
260 user: "0:0"
```

Username: airflow  
Password: airflow

# 建立docker airflow 初始化

Command:  
docker compose up airflow-init



```
Windows PowerShell
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> docker compose up airflow-init
time="2023-10-26T20:09:36+08:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
time="2023-10-26T20:09:36+08:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
[+] Running 0/23
- redis Pulling 4.4s
  - 4f4fb700ef54 Waiting 0.4s
  - a378f10b3218 Waiting 0.1s
  - b266cd8112a6 Waiting 0.0s
```

# 建立完成畫面

```
airflow_docker-airflow-init-1 | [2023-10-26T12:14:33.937+0000] {manager.py:555} INFO - Added Permission %s to role %s
airflow_docker-airflow-init-1 | [2023-10-26T12:14:33.968+0000] {manager.py:499} INFO - Created Permission View: %s
airflow_docker-airflow-init-1 | [2023-10-26T12:14:33.977+0000] {manager.py:555} INFO - Added Permission %s to role %s
airflow_docker-airflow-init-1 | [2023-10-26T12:14:33.993+0000] {manager.py:499} INFO - Created Permission View: %s
airflow_docker-airflow-init-1 | [2023-10-26T12:14:34.001+0000] {manager.py:555} INFO - Added Permission %s to role %s
airflow_docker-airflow-init-1 | [2023-10-26T12:14:34.808+0000] {manager.py:211} INFO - Added user %s
airflow_docker-airflow-init-1 | User "airflow" created with role "Admin"
airflow_docker-airflow-init-1 | 2.7.2
airflow_docker-airflow-init-1 exited with code 0
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> |
```

# 啟動docker compose

Command:  
docker compose up

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> docker compose up
time="2023-10-26T20:15:11+08:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
time="2023-10-26T20:15:11+08:00" level=warning msg="The \"AIRFLOW_UID\" variable is not set. Defaulting to a blank string."
[+] Running 7/7
 - Container airflow_docker-redis-1      Recreated      0.5s
 - Container airflow_docker-postgres-1   Recreated      0.5s
 - Container airflow_docker-airflow-init-1 Recreated      0.1s
 - Container airflow_docker-airflow-webserver-1 Created        0.1s
 - Container airflow_docker-airflow-worker-1 Created        0.1s
 - Container airflow_docker-airflow-scheduler-1 Created        0.1s
 - Container airflow_docker-airflow-triggerer-1 Created        0.1s
Attaching to airflow_docker-airflow-init-1, airflow_docker-airflow-scheduler-1, airflow_docker-airflow-triggerer-1, airflow_docker-airflow-webserver-1, airflow_docker-airflow-worker-1, airflow_docker-postgres-1, airflow_docker-redis-1
airflow_docker-redis-1 | 1:C 26 Oct 2023 12:15:12.601 # WARNING Memory overcommit must be enabled! Without it, a background save or replication may fail under low memory condition. Being disabled, it can also cause failures without low memory condition, see https://github.com/jemalloc/jemalloc/issues/1328. To fix this issue add 'vm.overcommit_memory = 1' to /etc/sysctl.conf and then reboot or run the command 'sysctl vm.overcommit_memory=1' for this to take effect
```



# 查看當前container狀態，確定啟動成功

Command:  
docker ps

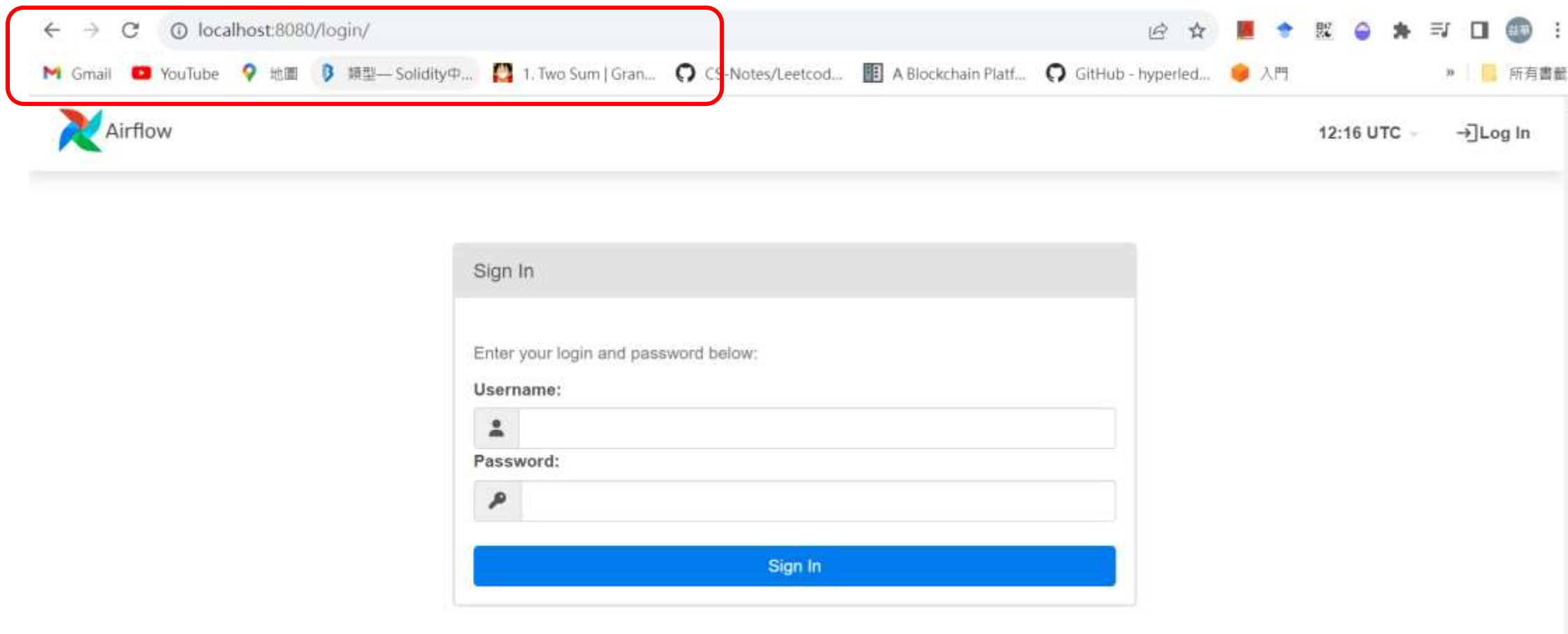
```
PS C:\Users\jerry> docker ps
```

CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS
362a77edfa0a	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	35 seconds ago	Up 5 seconds (health: starting)	8080/t
cp	airflow_docker-airflow-scheduler-1				
e21ca6c0df98	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	35 seconds ago	Up 5 seconds (health: starting)	8080/t
cp	airflow_docker-airflow-worker-1				
d1bd24aeef4	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	35 seconds ago	Up 5 seconds (health: starting)	0.0.0.
0:8080->8080/tcp	airflow_docker-airflow-webserver-1				
884754c6a479	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	35 seconds ago	Up 5 seconds (health: starting)	8080/t
cp	airflow_docker-airflow-triggerer-1				
14b2ed3ba694	postgres:13	"docker-entrypoint.s..."	35 seconds ago	Up 33 seconds (healthy)	5432/t
cp	airflow_docker-postgres-1				
a91b8584285d	redis:latest	"docker-entrypoint.s..."	35 seconds ago	Up 33 seconds (healthy)	6379/t
cp	airflow_docker-redis-1				

```
PS C:\Users\jerry>
```



# 瀏覽器輸入 localhost:8080



The screenshot shows a web browser window with the address bar highlighted in red, displaying the URL `localhost:8080/login/`. The browser's address bar also shows navigation icons (back, forward, refresh) and a search icon. Below the address bar, there are several tabs open, including Gmail, YouTube, 地圖, 類型—Solidity中..., 1. Two Sum | Gran..., CS-Notes/Leetcod..., A Blockchain Platf..., GitHub - hyperled..., and 入門. The main content area of the browser shows the Airflow login page. The page has a header with the Airflow logo on the left, the time `12:16 UTC` in the center, and a `Log In` button on the right. The main body of the page contains a `Sign In` form. The form has a title `Sign In` and a prompt `Enter your login and password below:`. It includes two input fields: `Username:` and `Password:`, each with a corresponding icon (a person icon for the username and a key icon for the password). Below the input fields is a blue `Sign In` button.


# 輸入帳號密碼

**Username: airflow**  
**Password: airflow**


Sign In

Enter your login and password below:

Username:

 airflow

Password:

 .....

Sign In

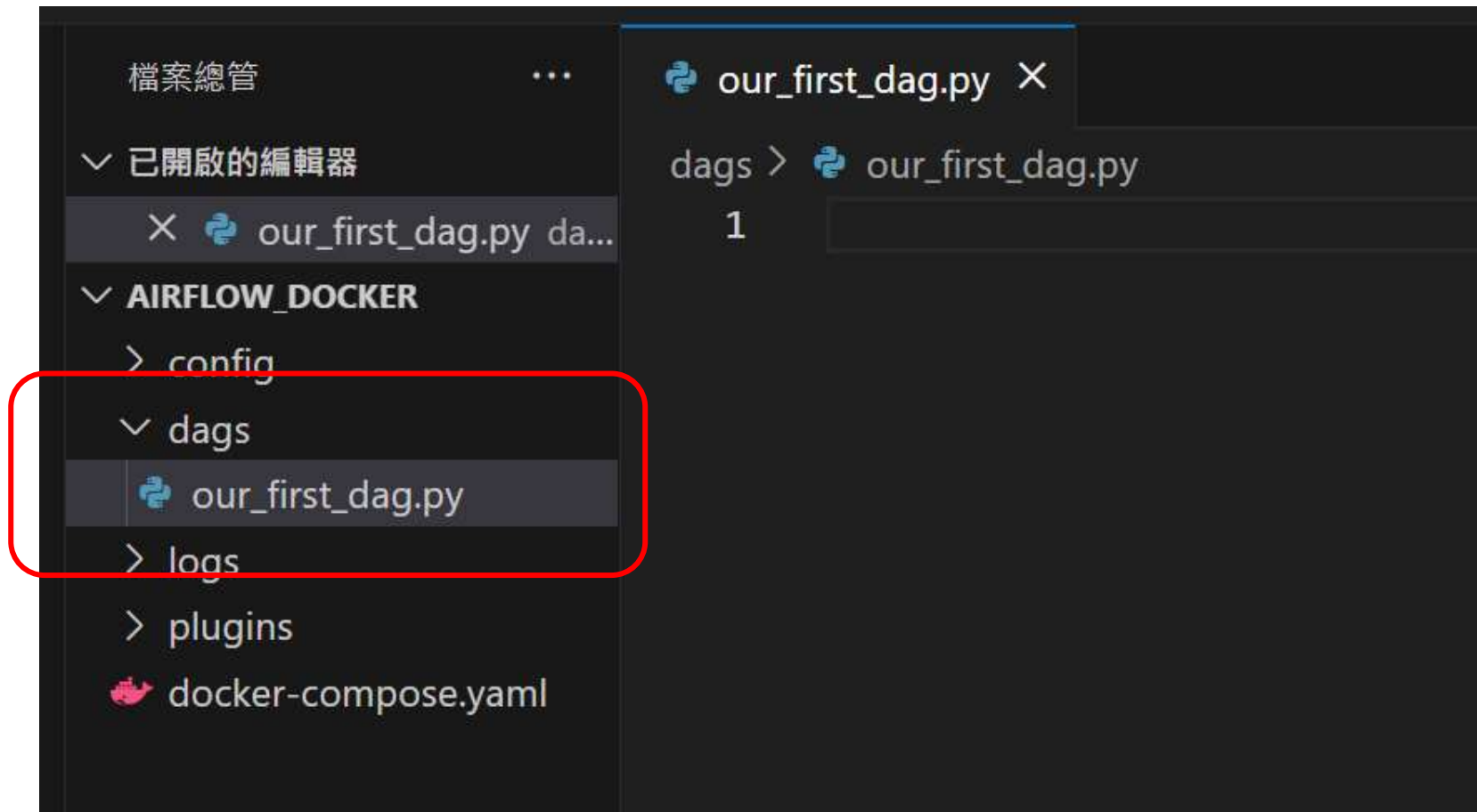
# 成功登入畫面

The screenshot shows the Apache Airflow web interface in a browser. The address bar indicates the URL is `localhost:8080/home`. The top navigation bar includes links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. The current time is 12:17 UTC. The main section is titled "DAGs" and features a filter bar with buttons for "All 53", "Active 0", "Paused 53", "Running 0", and "Failed 0". There is also a "Filter DAGs by tag" input field, a "Search DAGs" search bar, and an "Auto-refresh" toggle switch. Below the filter bar is a table listing DAGs with columns for DAG name, Owner, Runs, Schedule, Last Run, Next Run, and Recent Tasks. The table contains five rows of DAGs, all owned by "airflow".

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
<input type="checkbox"/> dataset_consumes_1 consumes dataset-scheduled	airflow	0	Dataset		On s3://dag1/output_1.txt	0
<input type="checkbox"/> dataset_consumes_1_and_2 consumes dataset-scheduled	airflow	0	Dataset		0 of 2 datasets updated	0
<input type="checkbox"/> dataset_consumes_1_never_scheduled consumes dataset-scheduled	airflow	0	Dataset		0 of 2 datasets updated	0
<input type="checkbox"/> dataset_consumes_unknown_never_scheduled dataset-scheduled	airflow	0	Dataset		0 of 2 datasets updated	0
<input type="checkbox"/> dataset_produces_1						

# Airflow DAG with Bash Operator

# 於 dags 資料夾中建立一個.py檔



# 載入相關套件並定義參數

```
our_first_dag.py X
dags > our_first_dag.py
1  from datetime import datetime, timedelta
2  from airflow import DAG
3  from airflow.operator.bash import BashOperator
4
5  # 設定重置時間5分鐘&延遲2分鐘(每次等待重置的時間)
6  default_args = {
7      'owner': 'coder2j',
8      'retries': 5,
9      'retry_delay': timedelta(minute=2)
10
11  }
12
```

# 撰寫程式碼

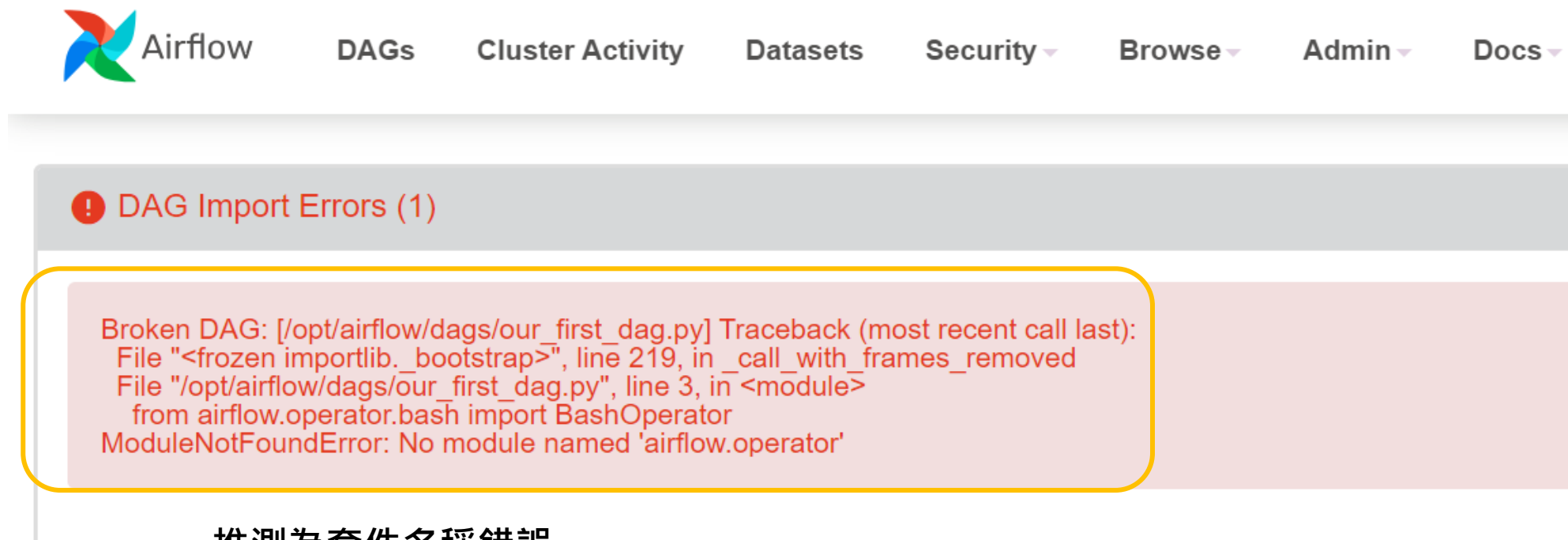
設定DAG中的相關資訊

```
13
14     with DAG(
15         dag_id='our_first_dag',
16         default_args=default_args,
17         description='This is our first dag that we write',
18         # 設定從甚麼時候啟動及按時啟動
19         # 2021年7月29日開始每天凌晨兩點啟動
20         start_date=datetime(2021, 7, 29, 2),
21         scheduler_interval='@daily'
22     ) as dag:
23     pass
24
```

使用bash執行task

```
22
23     ) as dag:
24         task1 = BashOperator(
25             task_id='first_task',
26             bash_command="echo Hello World. This is the first task!"
27         )
28
29         task1
```

# 執行發現有ERROR



The screenshot shows the Apache Airflow web interface. At the top is a navigation bar with the Airflow logo and links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs. Below the navigation bar, a grey banner indicates 'DAG Import Errors (1)'. A red-bordered box highlights the error details:

```
Broken DAG: [/opt/airflow/dags/our_first_dag.py] Traceback (most recent call last):  
  File "<frozen importlib._bootstrap>", line 219, in _call_with_frames_removed  
  File "/opt/airflow/dags/our_first_dag.py", line 3, in <module>  
    from airflow.operator.bash import BashOperator  
ModuleNotFoundError: No module named 'airflow.operator'
```

Below the error box, the text '推測為套件名稱錯誤' (Suspected to be a package name error) is displayed.



# 修正Module名稱


原本沒有+s


```
our_first_dag.py X
dags > our_first_dag.py
1 from datetime import datetime, timedelta
2 from airflow import DAG
3 from airflow.operator.bash import BashOperator
```

修改後

```
our_first_dag.py X
dags > our_first_dag.py
1 from datetime import datetime, timedelta
2 from airflow import DAG
3 from airflow.operators.bash import BashOperator
```

# timedelta 語法錯誤

AirflowDAGsCluster ActivityDatasetsSecurity ▾Browse ▾Admin ▾Docs ▾

 DAG Import Errors (1)

Broken DAG: [/opt/airflow/dags/our\_first\_dag.py] Traceback (most recent call last):  
File "<frozen importlib.\_bootstrap>", line 219, in \_call\_with\_frames\_removed  
File "/opt/airflow/dags/our\_first\_dag.py", line 9, in <module>  
    'retry\_delay': timedelta(minute=2)  
TypeError: 'minute' is an invalid keyword argument for new ()

# 修正timedelta語法


原本沒有+s

```
# 設定重置時間5分鐘&延遲2分鐘(每次等待重置的時間)
default_args = {
    'owner': 'coder2j',
    'retries': 5,
    'retry_delay': timedelta(minute=2)
}
```

修改後

```
5 # 設定重置時間5分鐘&延遲2分鐘(每次等待重置的時間)
6 default_args = {
7     'owner': 'coder2j',
8     'retries': 5,
9     'retry_delay': timedelta(minutes=2)
10 }
11
12
```

# scheduler interval 有錯誤



Airflow

DAGs Cluster Activity Datasets Security ▾ Browse ▾ Admin ▾ Docs ▾

! DAG Import Errors (1)

```
Broken DAG: [/opt/airflow/dags/our_first_dag.py] Traceback (most recent call last):
  File "<frozen importlib._bootstrap>", line 219, in _call_with_frames_removed
  File "/opt/airflow/dags/our_first_dag.py", line 14, in <module>
    with DAG(
TypeError: init () got an unexpected keyword argument 'scheduler interval'
```

# 修改錯誤

```
14 with DAG(  
15     dag_id='our_first_dag',  
16     default_args=default_args,  
17     description='This is our first dag that we w  
18     # 設定從甚麼時候啟動及按時啟動  
19     # 2021年7月29日開始每天凌晨兩點啟動  
20     start_date=datetime(2021, 7, 29, 2),  
21     scheduler_interval='@daily'
```

```
13  
14 with DAG(  
15     dag_id='our_first_dag',  
16     default_args=default_args,  
17     description='This is our first dag that we w  
18     # 設定從甚麼時候啟動及按時啟動  
19     # 2021年7月29日開始每天凌晨兩點啟動  
20     start_date=datetime(2021, 7, 29, 2),  
21     schedule_interval='@daily'
```

# 成功! 沒有錯誤訊息

The screenshot shows the Apache Airflow web interface in a browser window. The browser's address bar displays 'localhost:8080/home'. The Airflow navigation bar includes links for DAGs, Cluster Activity, Datasets, Security, Browse, Admin, and Docs, along with the current time '12:33 UTC' and a user profile icon labeled 'AA'.

The main section is titled 'DAGs' and features a filter bar with the following status counts: All (54), Active (0), Paused (54), Running (0), and Failed (0). There are also input fields for 'Filter DAGs by tag' and 'Search DAGs', and a toggle for 'Auto-refresh'.

DAG	Owner	Runs	Schedule	Last Run	Next Run	Recent Tasks
<input type="checkbox"/> dataset_consumes_1 consumes dataset-scheduled	airflow	<div><div></div><div></div><div></div><div></div></div>	Dataset		On s3://dag1/output_1.txt	<div><div></div><div></div><div></div><div></div><div></div></div>

# 可看到我們所新增的dag

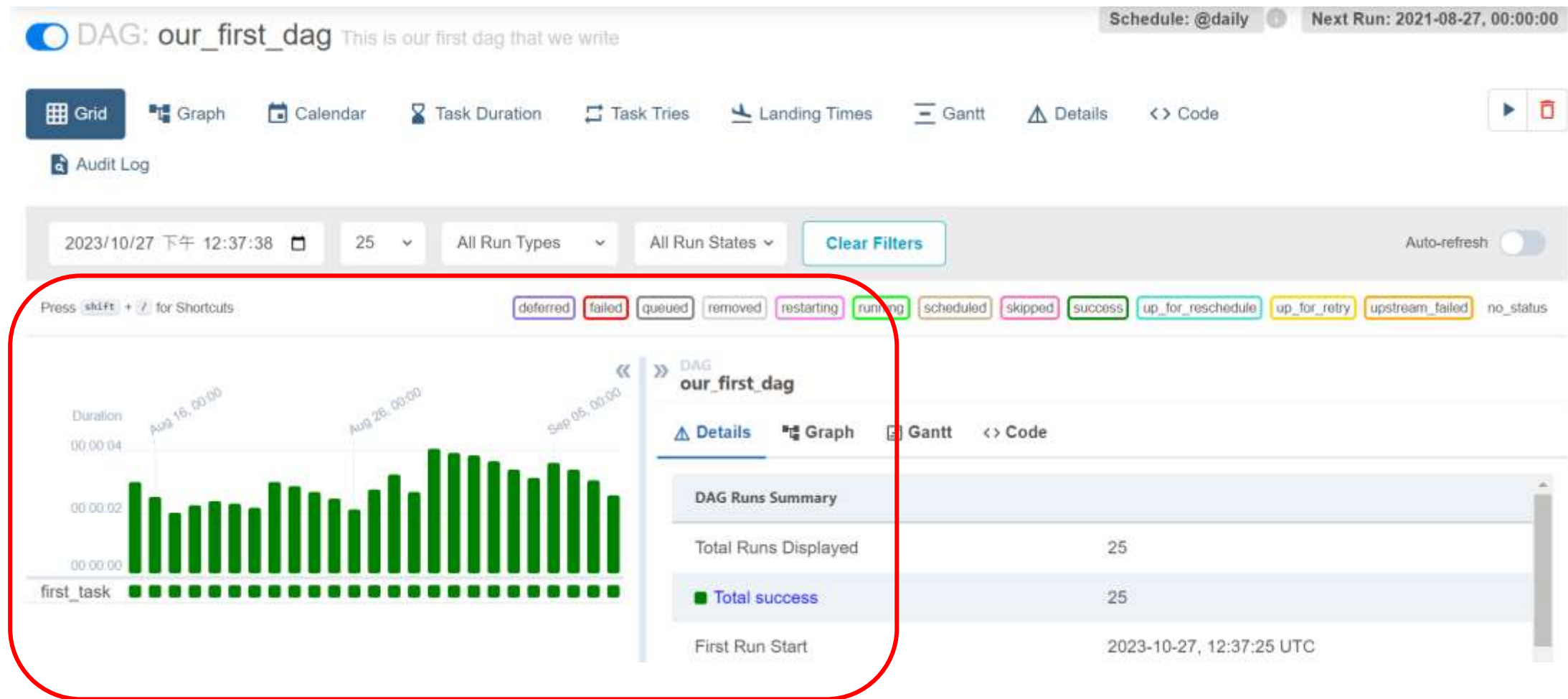
<input type="checkbox"/>	latest_only_with_trigger example3	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	4:00:00	2023-10-27, 08:33:23 <i>i</i>	<input type="checkbox"/>
<input type="checkbox"/>	our_first_dag	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	@daily <i>i</i>	2021-07-30, 00:00:00 <i>i</i>	<input type="checkbox"/>
<input type="checkbox"/>	tutorial example	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	1 day, 0:00:00	2023-10-26, 12:33:19 <i>i</i>	<input type="checkbox"/>
<input type="checkbox"/>	tutorial_dag example	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	None <i>i</i>		<input type="checkbox"/>
<input type="checkbox"/>	tutorial_taskflow_api example	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	None <i>i</i>		<input type="checkbox"/>

# 啟動

<input checked="" type="checkbox"/>	our_first_dag	coder2j	<div><div></div><div></div><div></div><div></div></div>	@daily ⓘ
<input type="checkbox"/>	tutorial example	airflow	<div><div></div><div></div><div></div><div></div></div>	1 day, 0:00:00
<input type="checkbox"/>	tutorial_dag example	airflow	<div><div></div><div></div><div></div><div></div></div>	None ⓘ
<input type="checkbox"/>	tutorial_taskflow_api example	airflow	<div><div></div><div></div><div></div><div></div></div>	None ⓘ



# 綠色為成功啟動



# 新增task2

```
) as dag:
    task1 = BashOperator(
        task_id='first_task',
        bash_command="echo Hello World. This is the first task!"
    )

    task2 = BashOperator(
        task_id='second_task',
        bash_command='echo Hey, I am task2 and will be running after task1!'
    )

# 設定 task2 在 task1 後執行
task1.set_downstream(task2)
```

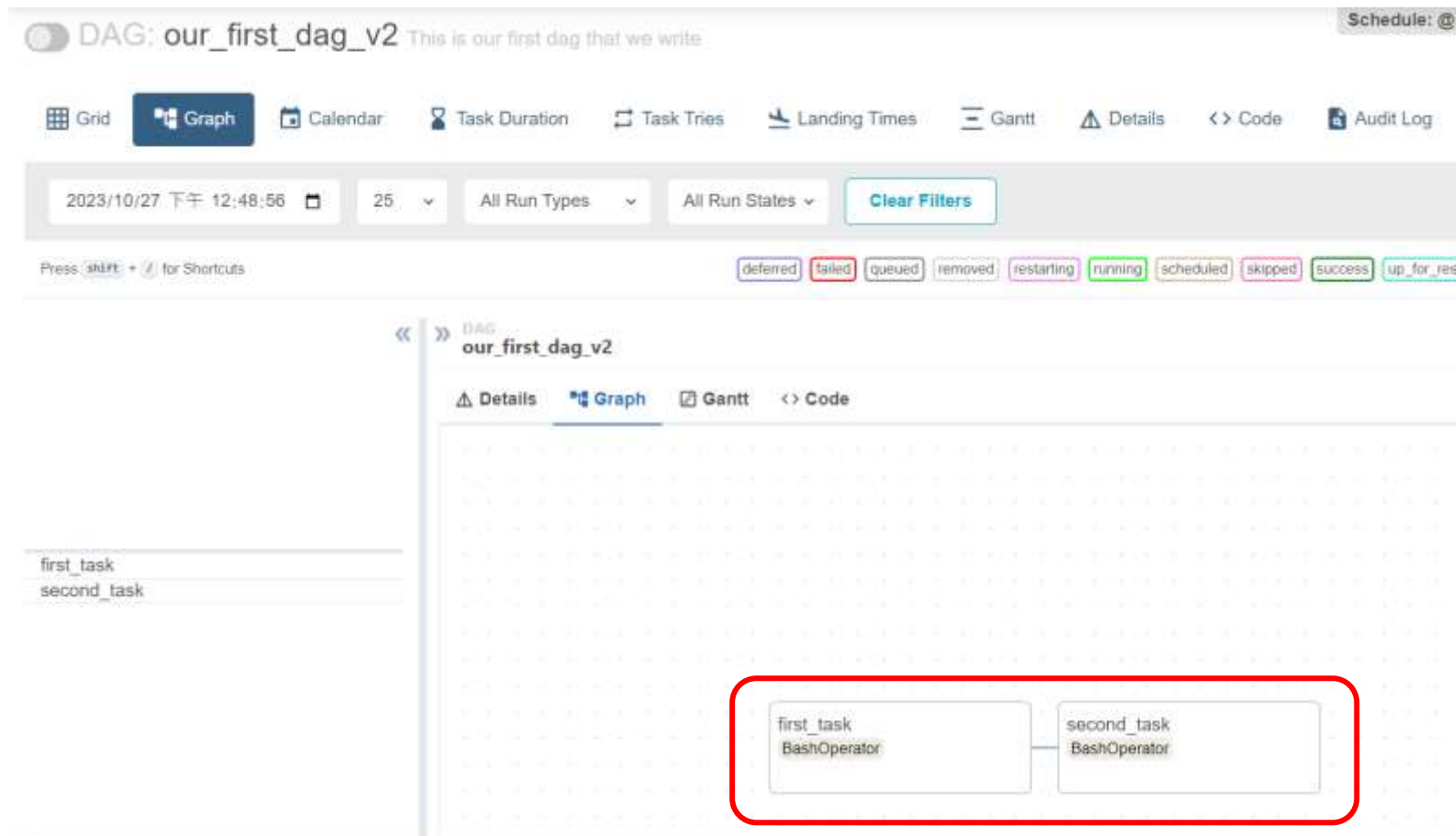
# 修改id避免混淆

```
14  with DAG(  
15      dag_id='our_first_dag_v2',  
16      default_args=default_args,  
17      description='This is our first dag that we write',  
18      # 設定從甚麼時候啟動及按時啟動  
19      # 2021年7月29日開始每天凌晨兩點啟動  
20      start_date=datetime(2021, 7, 29, 2),  
21      schedule_interval='@daily'  
22  )
```

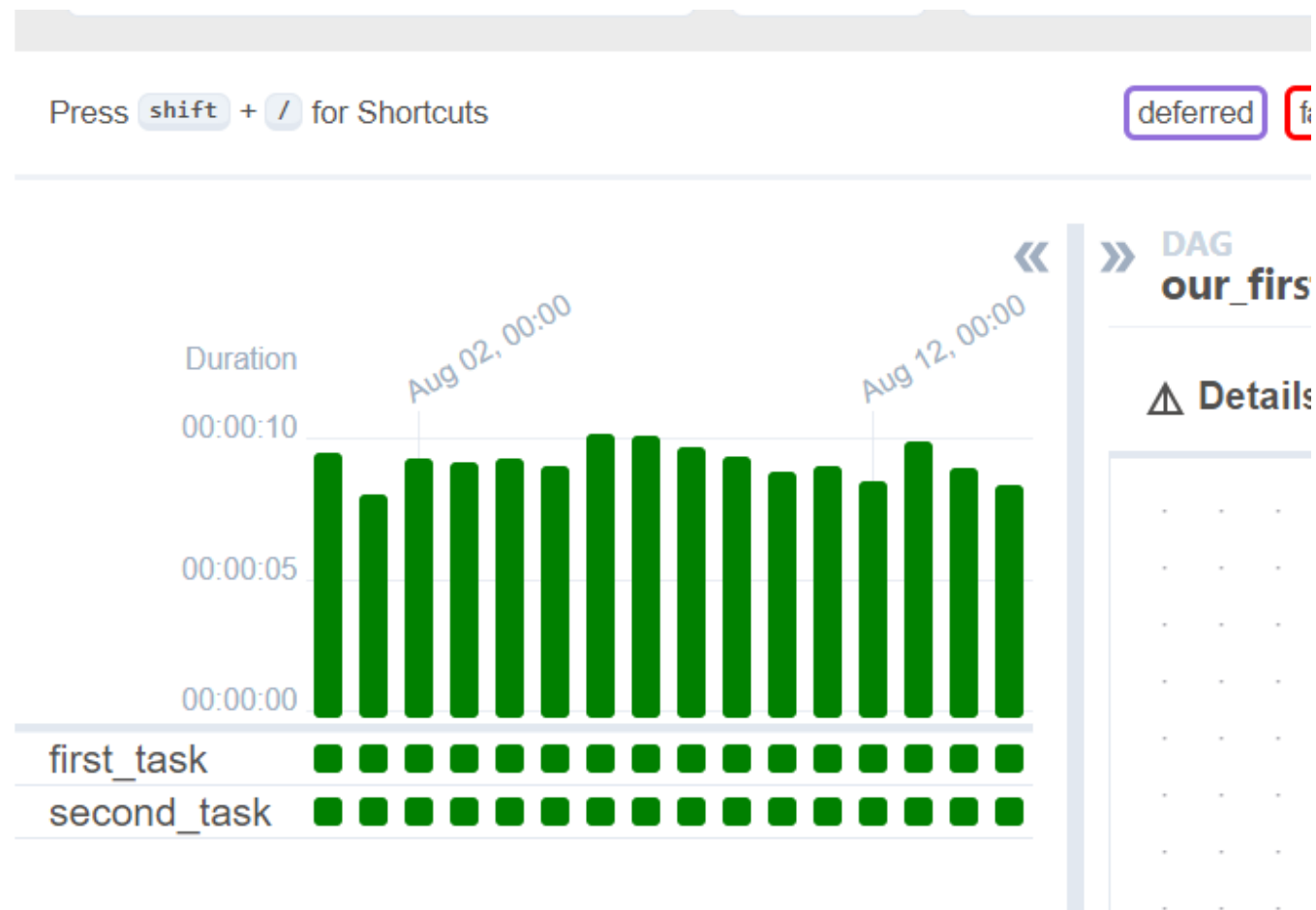
# 成功建立畫面

<input type="checkbox"/> our_first_dag	coder2j	<div><div>445</div><div></div><div></div><div></div></div>	@daily ⓘ	2022-10-17, 00:00:00 ⓘ	2022-10-18, 00:00:00 ⓘ
<input type="checkbox"/> our_first_dag_v2	coder2j	<div><div></div><div></div><div></div><div></div></div>	@daily ⓘ		2021-07-30, 00:00:00 ⓘ
<input type="checkbox"/> tutorial example	airflow	<div><div></div><div></div><div></div><div></div></div>	1 day, 0:00:00		2023-10-26, 12:47:36 ⓘ
<input type="checkbox"/> tutorial_dag example	airflow	<div><div></div><div></div><div></div><div></div></div>	None ⓘ		

# task流程畫面



# 成功執行



# 再新增一個task3

```
task2 = BashOperator(  
    task_id='second_task',  
    bash_command='echo Hey, I am task2 and will be running after task1!'  
)  
  
task3 = BashOperator(  
    task_id='third_task',  
    bash_command='echo Hey, I am task3 and will be running after task1 at the same time as task2!'  
)
```

# 設定task執行順序，第一種方式

```
41  
42     # task3 在task1後與task2同時執行  
43     # 第一種方法  
44     task1.set_downstream(task2)  
45     task1.set_downstream(task3)
```



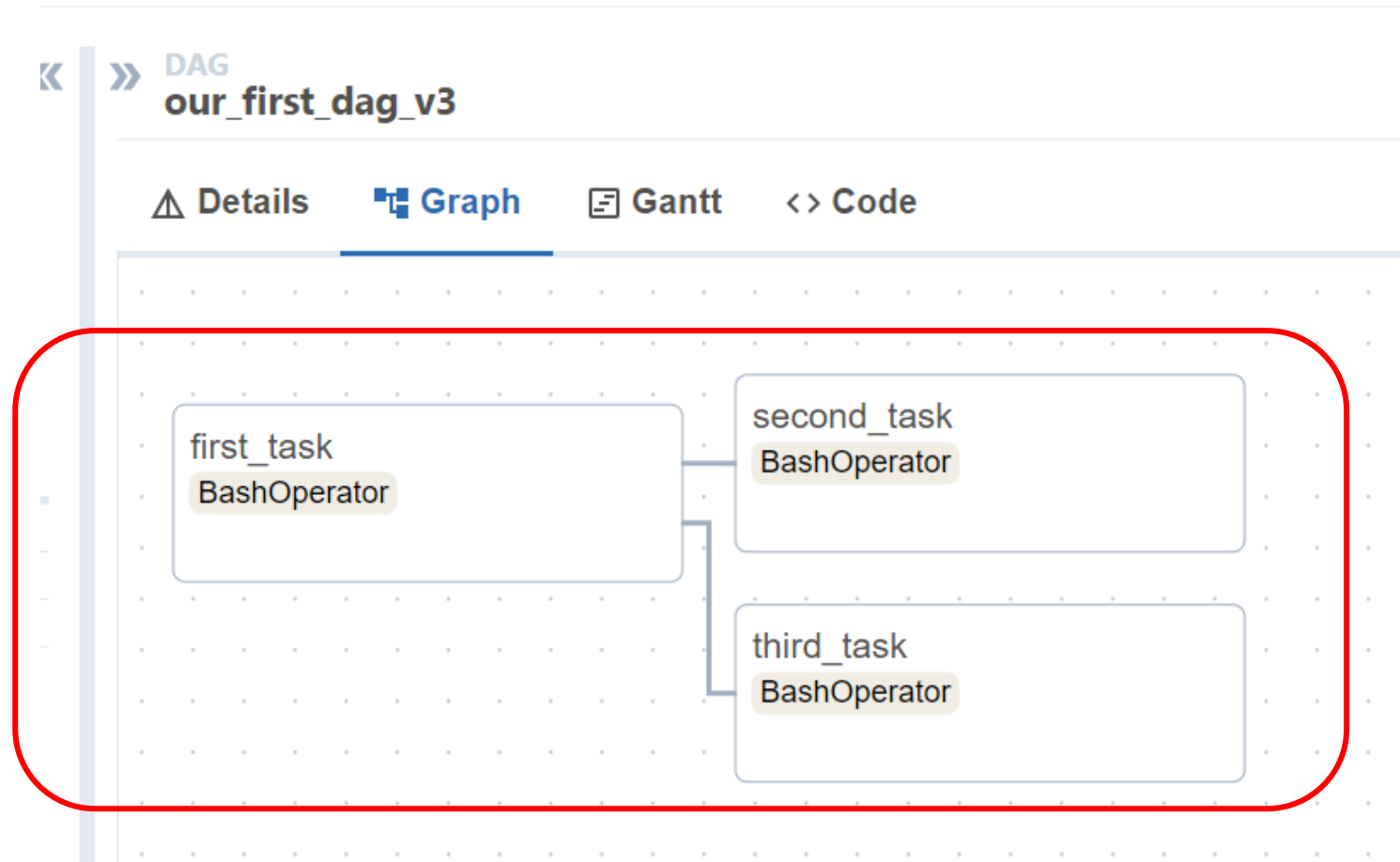
# 修改id避免混淆

```
4 with DAG(  
5     dag_id='our_first_dag_v3',  
6     default_args=default_args,  
7     description='This is our first dag that we write',  
8     # 設定從甚麼時候啟動及按時啟動  
9     # 2021年7月29日開始每天凌晨兩點啟動  
10    start_date=datetime(2021, 7, 29, 2),  
11    schedule_interval='@daily'
```

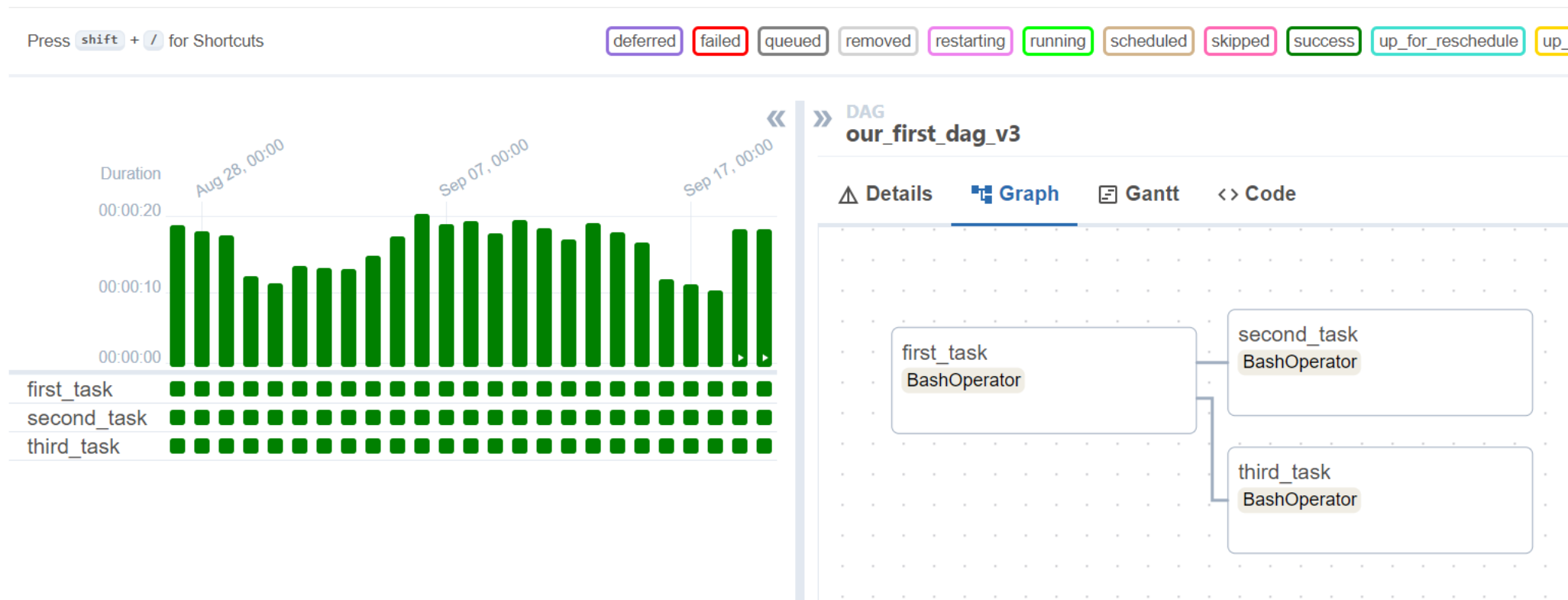
# 成功建立畫面

<input type="checkbox"/>	latest_only_with_trigger example3	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	4:00:00
<input type="checkbox"/>	our_first_dag_v2	coder2j	<input type="checkbox"/> <input checked="" type="checkbox"/> 192 <input checked="" type="checkbox"/> 8 <input type="checkbox"/>	@daily ⓘ 202
<input type="checkbox"/>	our_first_dag_v3	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	@daily ⓘ
<input type="checkbox"/>	tutorial example	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	1 day, 0:00:00

# task流程畫面



# 成功執行



# 設定task執行順序，第二種方式

```
42  # task3 在task1後與task2同時執行
43  # 第一種方法
44  # task1.set_downstream(task2)
45  # task1.set_downstream(task3)
46
47  # 第二種方法
48  task1 >> task2
49  task1 >> task3
50
```

# 設定task執行順序，第三種方式

```
42  # task3 在task1後與task2同時執行
43  # 第一種方法
44  # task1.set_downstream(task2)
45  # task1.set_downstream(task3)
46
47  # 第二種方法
48  # task1 >> task2
49  # task1 >> task3
50
51  # 第三種方法
52  task1 >> [task2, task3]
```

# Airflow DAG with Python Operator

# 改用python操作airflow

載入airflow.python

```
our_first_dag.py  python_operator.py X
dags > python_operator.py
1  from datetime import datetime, timedelta
2  from airflow import DAG
3  from airflow.operators.python import PythonOperator
4
5  default_args = {
6      'owner': 'coder2j',
7      'retries': 5,
8      'retry_delay': timedelta(minutes=2)
9
10 }
11
```



# 使用python定義一個greet function

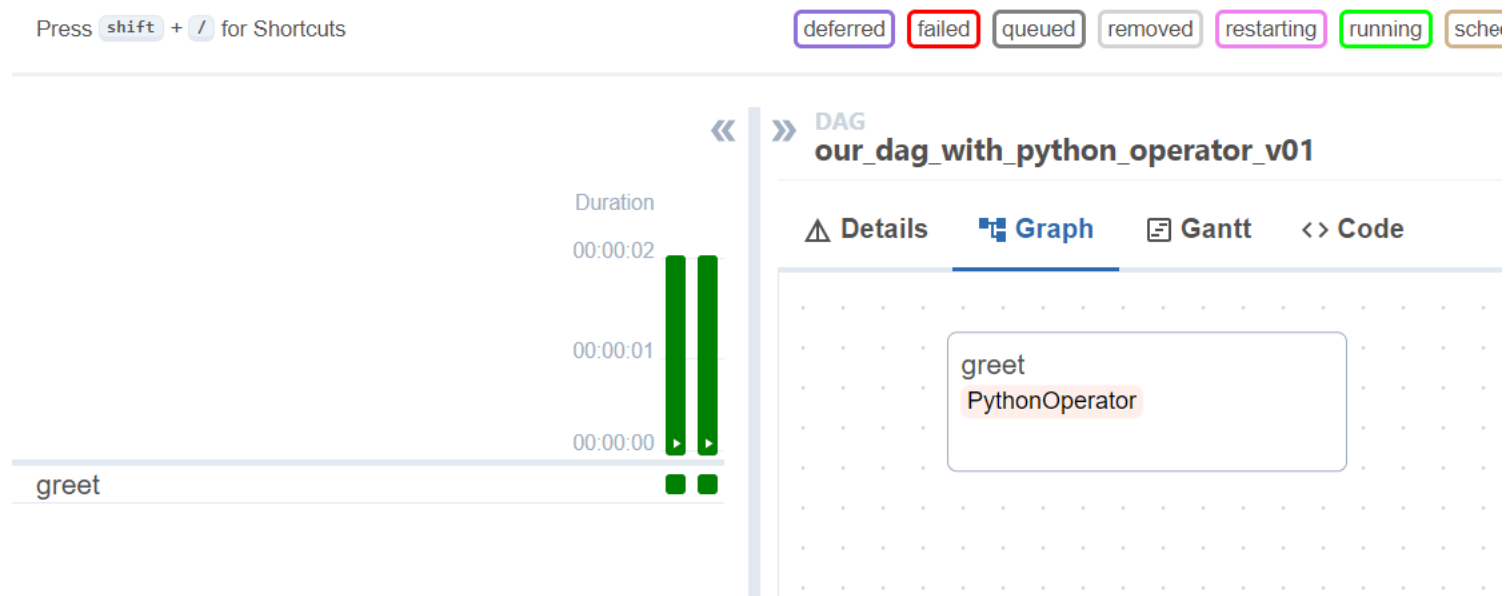
```
12  def greet():
13      print("Hello World!")
14
15  with DAG(
16      default_args=default_args,
17      dag_id='our_dag_with_python_operator_v01',
18      description='Oure first dag using python operator',
19      start_date=datetime(2023, 10, 28),
20      schedule_interval='@daily'
21  ) as dag:
22      task1 = PythonOperator(
23          task_id='greet',
24          python_callable=greet
25      )
26
27  task1
```

並將fun寫入task中

# 成功建立畫面

<input type="checkbox"/>	latest_only_with_trigger example3	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	4:00:00
<input type="checkbox"/>	our_dag_with_python_operator_v01	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	@daily ⓘ
<input type="checkbox"/>	our_first_dag_v5	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	@daily ⓘ
<input type="checkbox"/>	tutorial example	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>	1 day, 0:00:00

# 成功執行



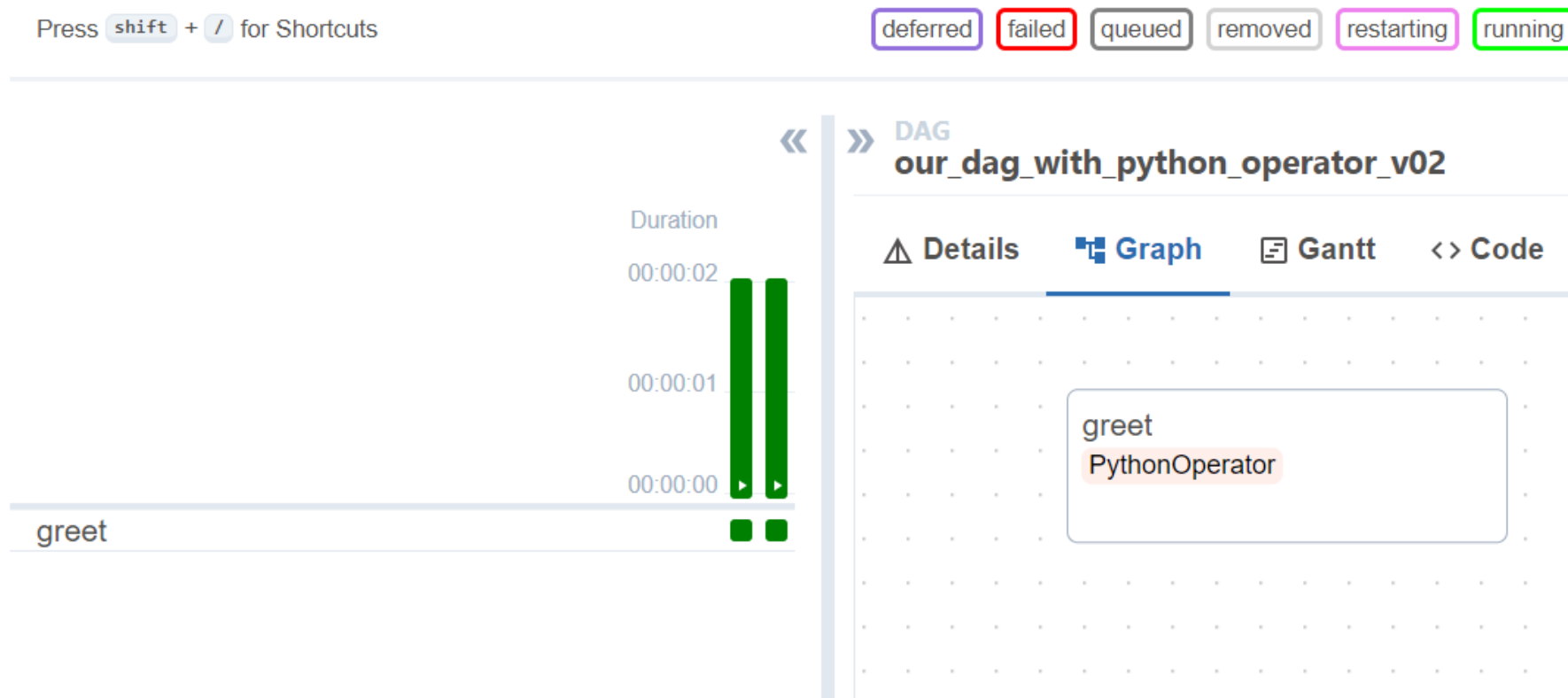
# 將greet()修改為回傳名字、年齡

```
12 def greet(name, age):
13     print(f"Hello World! My name is {name}, "
14           f"and I am {age} years old!")
15
16 with DAG(
17     default_args=default_args,
18     dag_id='our_dag_with_python_operator_v02',
19     description='Our first dag using python operator',
20     start_date=datetime(2023, 10, 28),
21     schedule_interval='@daily'
22 ) as dag:
23     task1 = PythonOperator(
24         task_id='greet',
25         python_callable=greet,
26         op_kwargs={'name': 'Tom', 'age': 20}
27     )
28
29 task1
```

# 成功建立畫面

<input type="checkbox"/>	latest_only_with_trigger example3	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/>	our_dag_with_python_operator_v02	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/>	our_first_dag_v5	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/>	tutorial example	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>

# 成功執行



# Data Sharing via Airflow Xcoms

# 說明

- XComs 是 Airflow 中的一種機制，用於讓不同的 Task 之間進行通信。
- 預設情況下，Task 是完全隔離的，可能在完全不同的機器上運行。XComs 由一個 key（基本上是它的名稱），以及它來自的 task\_id 和 dag\_id 來識別。
- 它們可以有任何（可序列化）值，但它們只設計用於小量數據；
- 不要使用它們來傳遞大量數據



# 新增一個 get\_name() function

```
12  def greet(name, age):  
13      print(f"Hello World! My name is {name}, "  
14          f"and I am {age} years old!")  
15  
16  def get_name():  
17      return 'Jerry'  
18
```

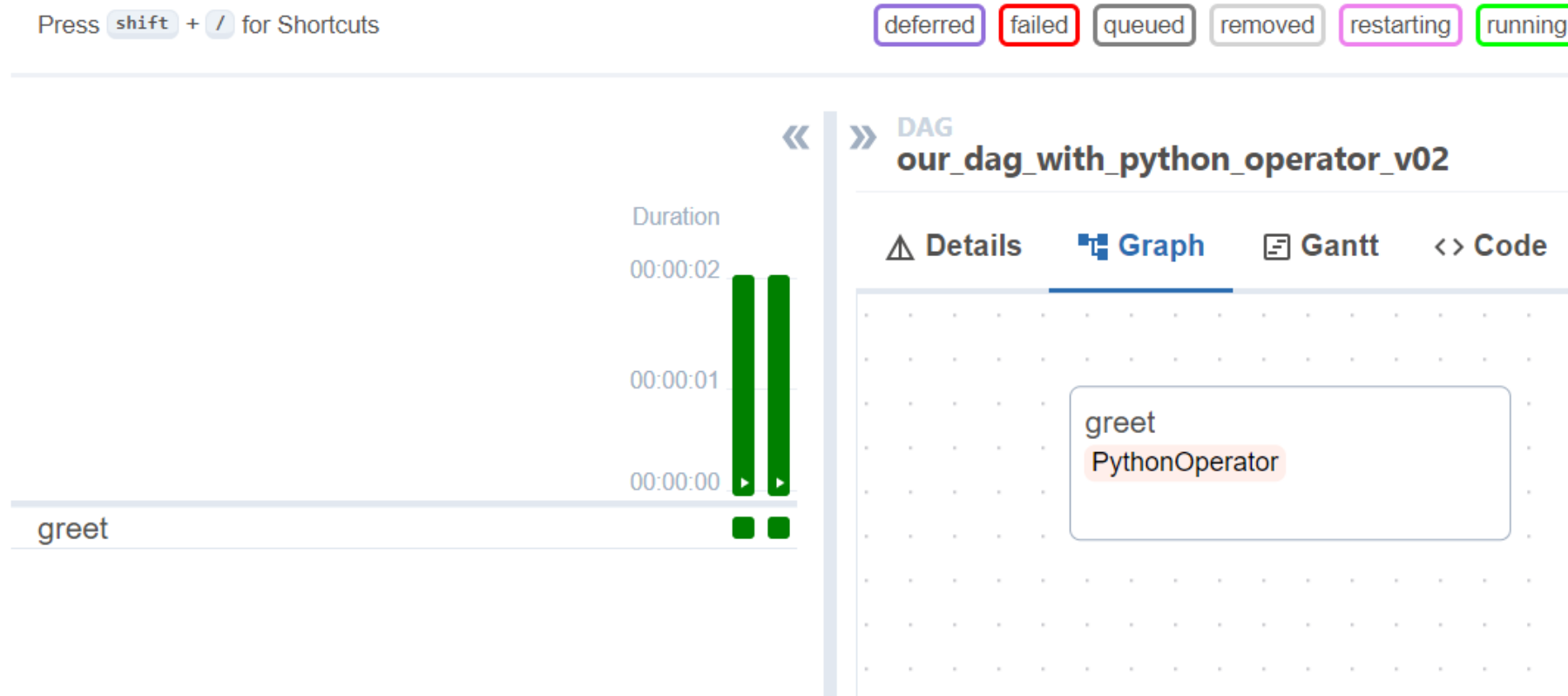
# 新增task2並寫入 get\_name function

```
19 with DAG(  
20     default_args=default_args,  
21     dag_id='our_dag_with_python_operator_v03',  
22     description='Our first dag using python operator',  
23     start_date=datetime(2023, 10, 28),  
24     schedule_interval='@daily'  
25 ) as dag:  
26     # task1 = PythonOperator(  
27     #     task_id='greet',  
28     #     python_callable=greet,  
29     #     op_kwargs={'name': 'Tom', 'age': 20}  
30     # )  
31  
32     task2 = PythonOperator(  
33         task_id='get_name',  
34         python_callable=get_name  
35     )  
36  
37 # task1  
38 task2
```

# 成功建立畫面

<input type="checkbox"/>	our_dag_with_python_operator_v02	coder2j	<div><div></div><div>3</div><div></div><div></div></div>	@daily	
<input type="checkbox"/>	our_dag_with_python_operator_v03	coder2j	<div><div></div><div></div><div></div><div></div></div>	@daily	
<input type="checkbox"/>	our_first_dag_v5	coder2j	<div><div></div><div></div><div></div><div></div></div>	@daily	
<input type="checkbox"/>	tutorial example	airflow	<div><div></div><div></div><div></div><div></div></div>	1 day, 0:00:00	

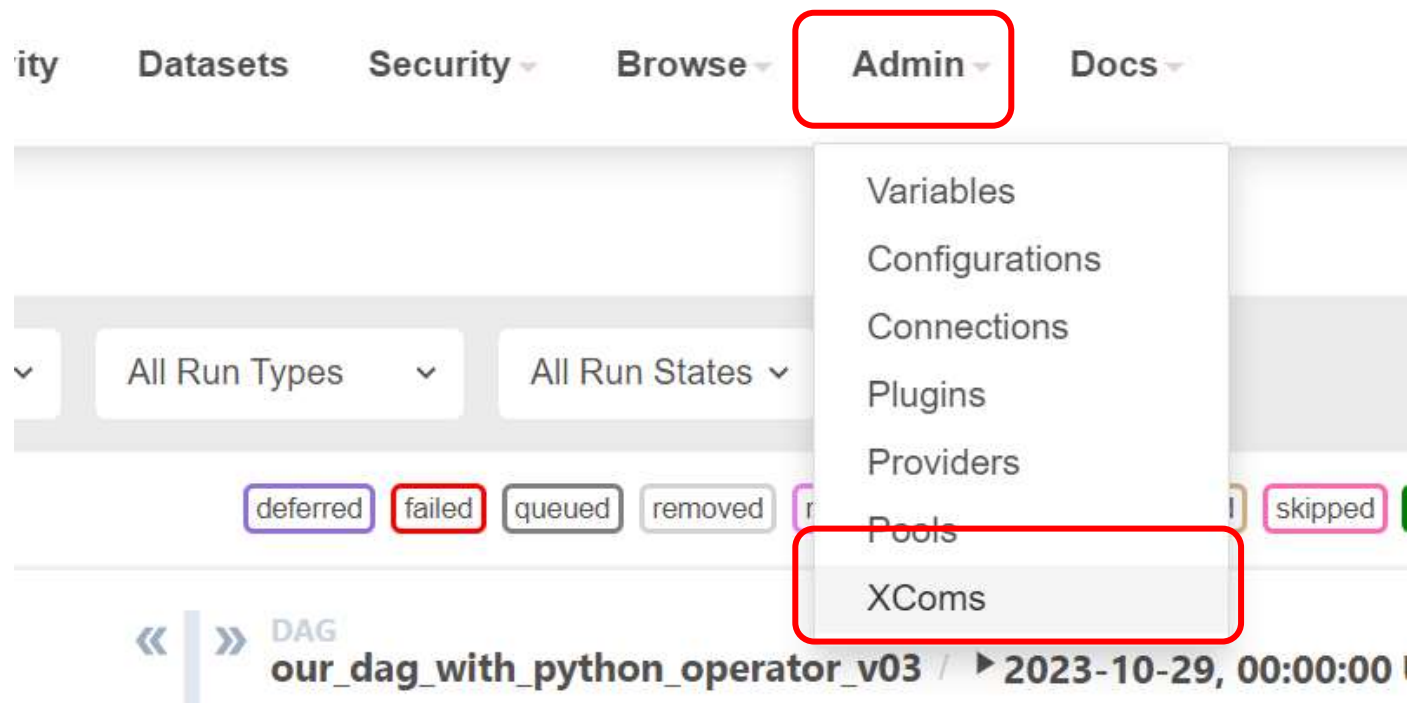
# 成功執行




# 查看log確實回傳Jerry

```
[2023-10-29, 10:43:14 UTC] {standard_task_runner.py:84} INFO - Running: ['***', 'tasks', 'run', 'ou  
[2023-10-29, 10:43:14 UTC] {standard_task_runner.py:85} INFO - Job 1222: Subtask get_name  
[2023-10-29, 10:43:14 UTC] {task_command.py:416} INFO - Running <TaskInstance: our_dag_with_python_  
[2023-10-29, 10:43:14 UTC] {taskinstance.py:1662} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER=  
[2023-10-29, 10:43:14 UTC] {python.py:194} INFO - Done. Returned value was: Jerry  
[2023-10-29, 10:43:14 UTC] {taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=our_dag_wi  
[2023-10-29, 10:43:14 UTC] {local_task_job_runner.py:228} INFO - Task exited with return code 0  
[2023-10-29, 10:43:14 UTC] {taskinstance.py:2778} INFO - 0 downstream tasks scheduled from follow-c
```

# 查看Xcoms



<input type="checkbox"/>		Key ↕	Value ↕	Timestamp ↕	Dag Id ↕	Task Id ↕
<input type="checkbox"/>		return_value	Jerry	2023-10-29, 10:43:14	our_dag_with_python_operator_v03	get_name ▼
<input type="checkbox"/>		return_value	Jerry	2023-10-29, 10:43:14	our_dag_with_python_operator_v03	get_name ▼
<input type="checkbox"/>		return_value	Jerry	2023-10-29, 10:43:13	our_dag_with_python_operator_v03	get_name ▼
Hello						

# 將greet修改為 ti.xcoms\_pull

也就說ti.xcom\_pull會回傳get\_name()中的值

```
11
12 def greet(age, ti):
13     name = ti.xcoms_pull(task_ids='get_name')
14     print(f"Hello World! My name is {name},")
15         f"and I am {age} years old!")
16
17 def get_name():
18     return 'Jerry'
19
```

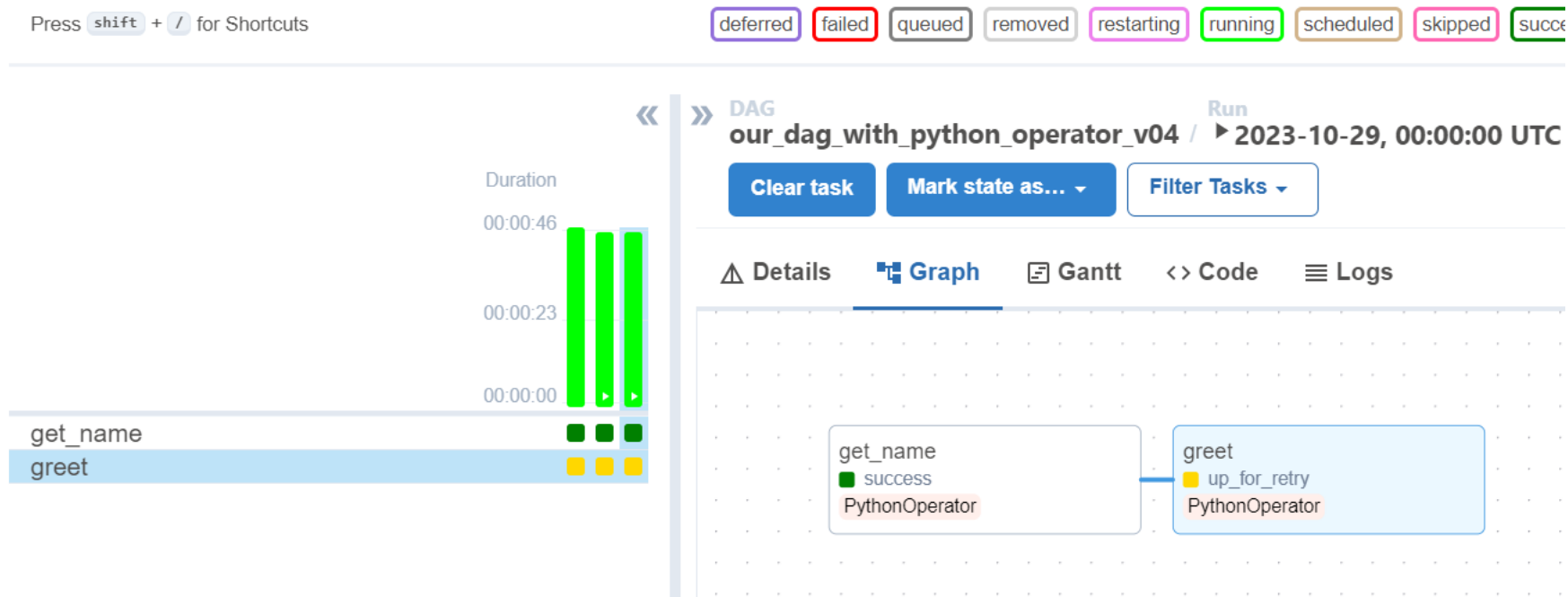


```
20 with DAG(  
21     default_args=default_args,  
22     dag_id='our_dag_with_python_operator_v04',  
23     description='Our first dag using python operator',  
24     start_date=datetime(2023, 10, 28),  
25     schedule_interval='@daily'  
26 ) as dag:  
27     task1 = PythonOperator(  
28         task_id='greet',  
29         python_callable=greet,  
30         op_kwargs={'age': 20}  
31     )  
32  
33     task2 = PythonOperator(  
34         task_id='get_name',  
35         python_callable=get_name  
36     )  
37  
38 # task1  
39 task2 >> task1
```

# 成功建立畫面

<input type="checkbox"/>	latest_only_with_trigger example3	airflow	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/>	our_dag_with_python_operator_v03	coder2j	<input type="checkbox"/> <input checked="" type="checkbox"/> 3 <input type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/>	our_dag_with_python_operator_v04	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>
<input type="checkbox"/>	our_first_dag_v5	coder2j	<input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/>

# 執行發現error



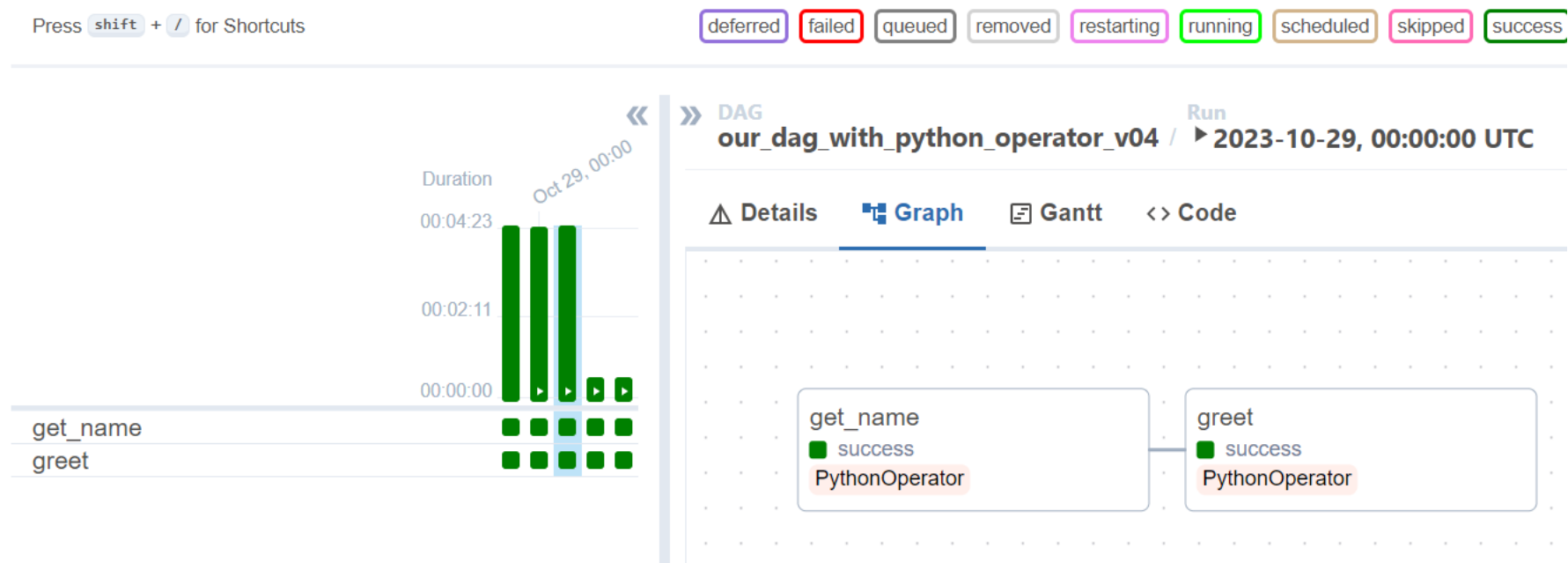
# 推測 xcoms\_pull 語法錯誤

```
[2023-10-29, 10:53:22 UTC] {taskinstance.py:1662} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='coder2j' AIRFLOW_C
[2023-10-29, 10:53:22 UTC] {taskinstance.py:1937} ERROR - Task failed with exception
Traceback (most recent call last):
  File "/home/airflow/.local/lib/python3.8/site-packages/airflow/operators/python.py", line 192, in execute
    return_value = self.execute_callable()
  File "/home/airflow/.local/lib/python3.8/site-packages/airflow/operators/python.py", line 209, in execute_callable
    return self.python_callable(*self.op_args, **self.op_kwargs)
  File "/opt/airflow/dags/python_operator.py", line 13, in greet
    name = ti.xcoms_pull(task_ids='get_name')
AttributeError: 'TaskInstance' object has no attribute 'xcoms_pull'
[2023-10-29, 10:53:22 UTC] {taskinstance.py:1400} INFO - Marking task as UP FOR RETRY. dag id=our dag with python_oper
[2023-10-29, 10:53:22 UTC] {standard_task_runner.py:104} ERROR - Failed to execute job 1227 for task greet ('TaskInsta
```

# Xcom後面沒有s

```
10     }
11
12     def greet(age, ti):
13         name = ti.xcom_pull(task_ids='get_name')
14         print(f"Hello World! My name is {name}, "
15               f"and I am {age} years old!")
16
17     def get_name():
18         return 'Jerry'
```

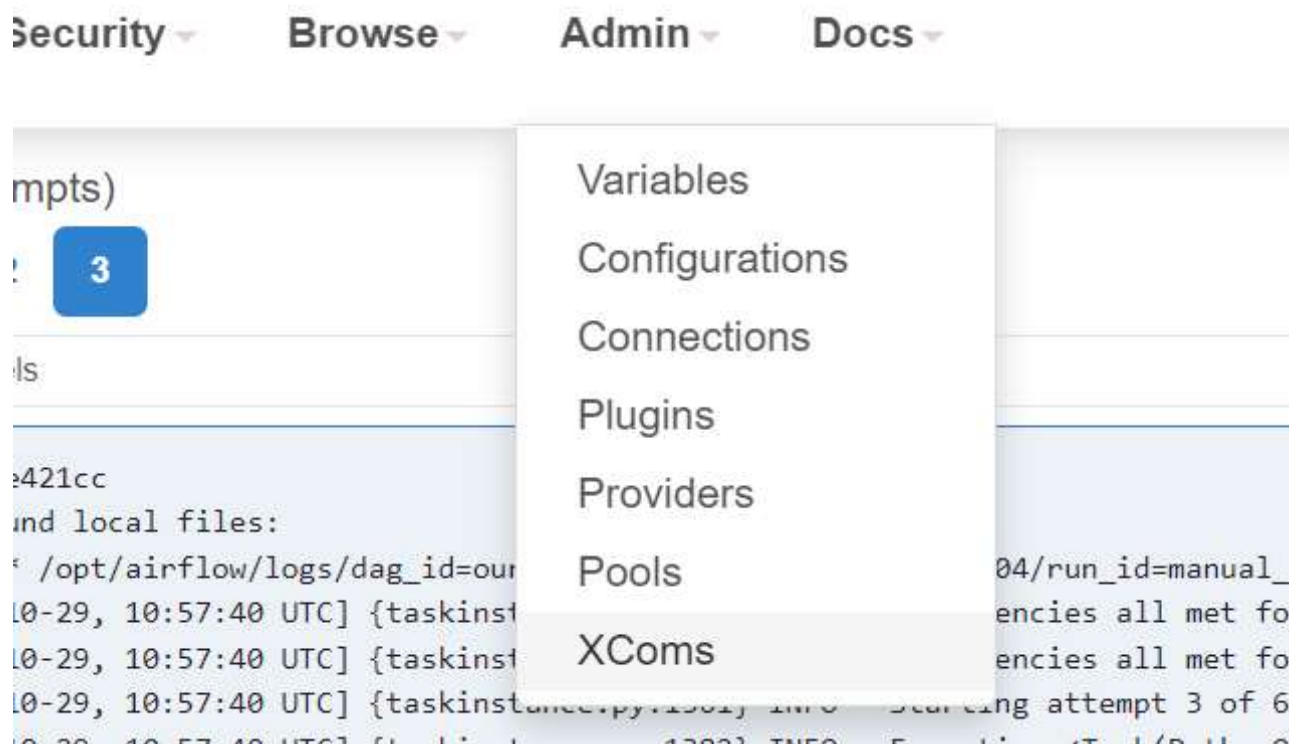
# 成功執行






# 查看logs

```
[2023-10-29, 10:57:40 UTC] {standard_task_runner.py:85} INFO - Job 1234: Subtask greet
[2023-10-29, 10:57:40 UTC] {task_command.py:416} INFO - Running <TaskInstance: our_dag_with_python_operator_v04.greet
[2023-10-29, 10:57:41 UTC] {taskinstance.py:1662} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='coder2j' AIRFLOW_C
[2023-10-29, 10:57:41 UTC] {logging_mixin.py:151} INFO - Hello World! My name is Jerry, and I am 20 years old!
[2023-10-29, 10:57:41 UTC] {python.py:194} INFO - Done. Returned value was: None
[2023-10-29, 10:57:41 UTC] {taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=our_dag_with_python_operator_
[2023-10-29, 10:57:41 UTC] {local_task_job_runner.py:228} INFO - Task exited with return code 0
[2023-10-29, 10:57:41 UTC] {taskinstance.py:2778} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

# 查看Xcoms





<div> <div>« &lt; 1 2 3 4 5 6 7 &gt; »</div> <div>Page size▼</div> <div>Actions▼</div> <div>←</div> </div>						
<input type="checkbox"/>		Key ↑	Value ↑	Timestamp ↑	Dag Id ↑	Task Id ↑
<input type="checkbox"/>		return_value	Jerry	2023-10-29, 10:57:42	our_dag_with_python_operator_v04	get_name ▼
<input type="checkbox"/>		return_value	Jerry	2023-10-29, 10:57:42	our_dag_with_python_operator_v04	get_name ▼
<input type="checkbox"/>		return_value	Jerry	2023-10-29, 10:53:20	our_dag_with_python_operator_v04	get_name ▼

# Xcom push 多個值

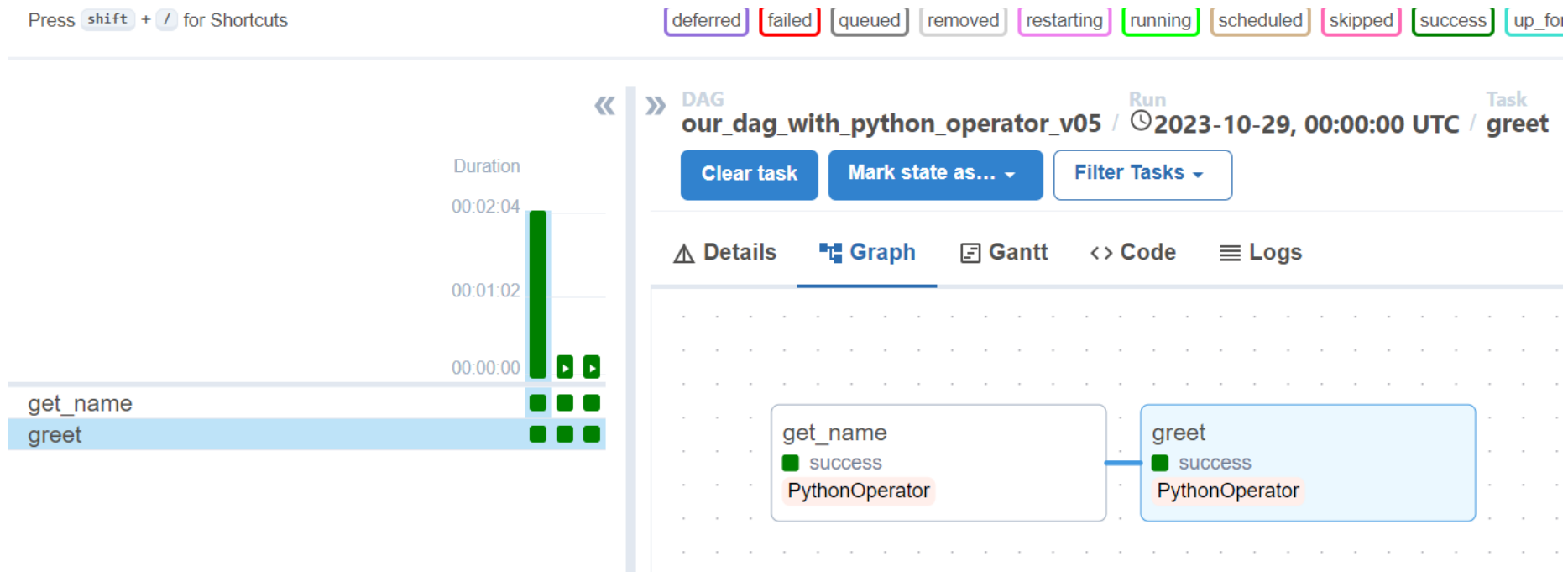
```
2 def greet(age, ti):
3     first_name = ti.xcom_pull(task_ids='get_name', key='first_name')
4     last_name = ti.xcom_pull(task_ids='get_name', key='last_name')
5     print(f"Hello World! My name is {first_name} {last_name}, "
6           f"and I am {age} years old!")
7
8 def get_name(ti):
9     ti.xcom_push(key='first_name', value='Jerry')
10    ti.xcom_push(key='last_name', value='Fridman')
```

```
22 with DAG(  
23     default_args=default_args,  
24     dag_id='our_dag_with_python_operator_v05',  
25     description='Our first dag using python operator',  
26     start_date=datetime(2023, 10, 28),  
27     schedule_interval='@daily'  
28 ) as dag:  
29     task1 = PythonOperator(  
30         task_id='greet',  
31         python_callable=greet,  
32         op_kwargs={'age': 20}  
33     )  
34  
35     task2 = PythonOperator(  
36         task_id='get_name',  
37         python_callable=get_name  
38     )  
39  
40 # task1  
41 task2 >> task1
```

# 成功建立畫面

<input type="checkbox"/>	our_dag_with_python_operator_v04	coder2j	<input type="radio"/>	<input checked="" type="radio"/>	<input type="radio"/>	<input type="radio"/>
<input type="checkbox"/>	our_dag_with_python_operator_v05	coder2j	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
<input type="checkbox"/>	our_first_dag_v5	coder2j	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
<input type="checkbox"/>	tutorial example	airflow	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>





# 成功執行



# 查看logs

```
[2023-10-29, 11:12:27 UTC] {standard_task_runner.py:85} INFO - Job 1245: Subtask greet
[2023-10-29, 11:12:27 UTC] {task_command.py:416} INFO - Running <TaskInstance: our_dag_with_python_operator_v05.greet scheduled__2
[2023-10-29, 11:12:27 UTC] {taskinstance.py:1662} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='coder2j' AIRFLOW_CTX_DAG_ID='c
[2023-10-29, 11:12:27 UTC] {logging_mixin.py:151} INFO - Hello World! My name is Jerry Fridman,and I am 20 years old!
[2023-10-29, 11:12:27 UTC] {python.py:194} INFO - Done. Returned value was: None
[2023-10-29, 11:12:27 UTC] {taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=our_dag_with_python_operator_v05, task_id
[2023-10-29, 11:12:27 UTC] {local_task_job_runner.py:228} INFO - Task exited with return code 0
[2023-10-29, 11:12:27 UTC] {taskinstance.py:2778} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

# 查看xcom

<div>« &lt; 1 2 3 4 5 6 7 &gt; » Page size Actions ←</div>					
<input type="checkbox"/>	Key ↑	Value ↑	Timestamp ↑	Dag Id ↑	Task Id ↑
<input type="checkbox"/> 	last_name	Fridman	2023-10-29, 11:11:55	our_dag_with_python_operator_v05	get_name ▼
<input type="checkbox"/> 	first_name	Jerry	2023-10-29, 11:11:55	our_dag_with_python_operator_v05	get_name ▼
<input type="checkbox"/> 	last_name	Fridman	2023-10-29, 11:11:55	our_dag_with_python_operator_v05	get_name ▼
<input type="checkbox"/> 	first_name	Jerry	2023-10-29, 11:11:55	our_dag_with_python_operator_v05	get_name ▼

# 新增 get\_age() function

```
12 def greet(ti):
13     first_name = ti.xcom_pull(task_ids='get_name', key='first_name')
14     last_name = ti.xcom_pull(task_ids='get_name', key='last_name')
15     age = ti.xcom_pull(task_ids='get_age', key='age')
16     print(f"Hello World! My name is {first_name} {last_name}, "
17           f"and I am {age} years old!")
18
19 def get_name(ti):
20     ti.xcom_push(key='first_name', value='Jerry')
21     ti.xcom_push(key='last_name', value='Fridman')
22
23 def get_age(ti):
24     ti.xcom_push(key='age', value=19)
```



# 修改id避免混淆

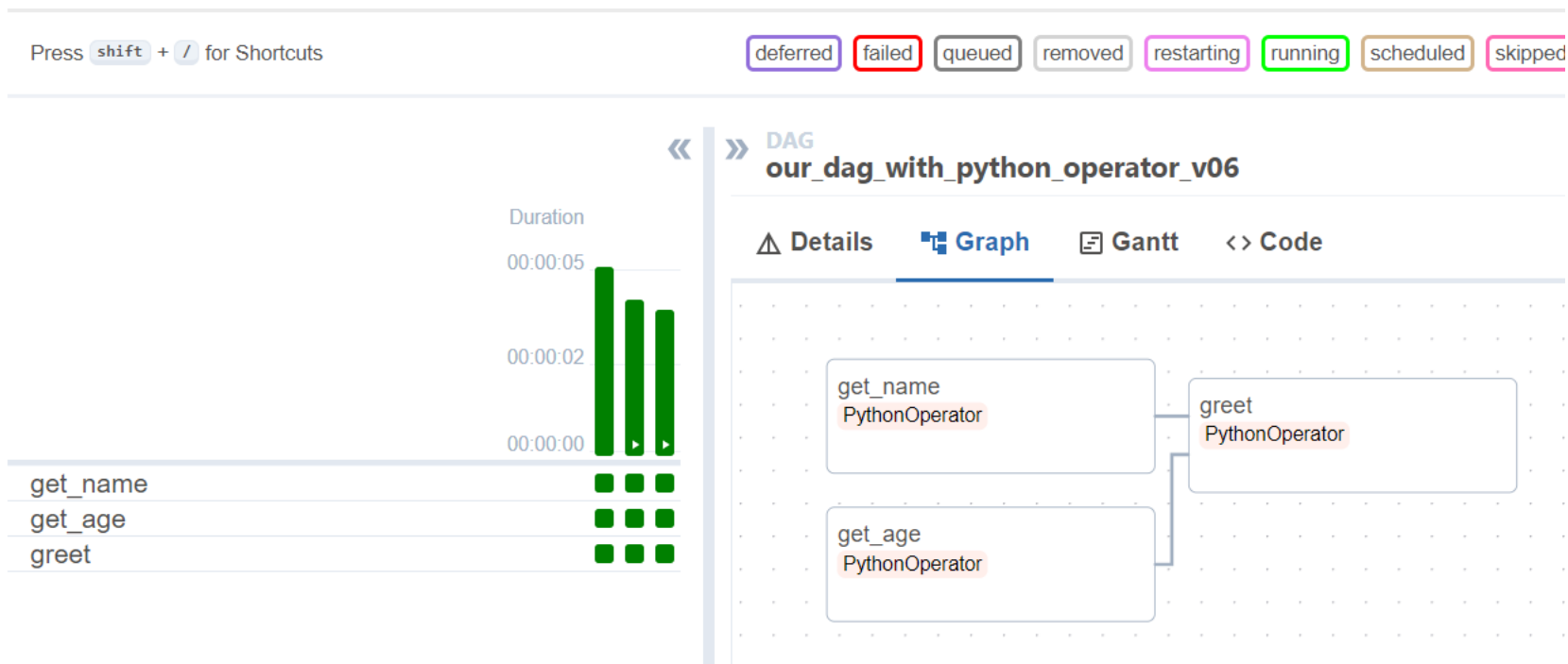
```
26 with DAG(  
27     default_args=default_args,  
28     dag_id='our_dag_with_python_operator_v06',  
29     description='Our first dag using python operator',  
30     start_date=datetime(2023, 10, 28),  
31     schedule_interval='@daily'  
32 ) as dag:  
33     task1 = PythonOperator(  
34         task_id='greet',  
35         python_callable=greet  
36         # op_kwargs={'age': 20}  
37     )  
38
```

```
39     task2 = PythonOperator(  
40         task_id='get_name',  
41         python_callable=get_name  
42     )  
43  
44     task3 = PythonOperator(  
45         task_id='get_age',  
46         python_callable=get_age  
47     )  
48  
49     # task1  
50     # task2 >> task1  
51     [task2, task3] >> task1
```

# 成功建立畫面

<input type="checkbox"/>	latest_only_with_trigger example3	airflow	<input type="radio"/>	<input type="radio"/>
<input type="checkbox"/>	our_dag_with_python_operator_v06	coder2j	<input type="radio"/>	<input type="radio"/>
<input type="checkbox"/>	our_first_dag_v5	coder2j	<input type="radio"/>	<input type="radio"/>
<input type="checkbox"/>	tutorial example	airflow	<input type="radio"/>	<input type="radio"/>


# 成功執行



# 查看logs

```
[2023-10-29, 11:22:34 UTC] {task_command.py:416} INFO - Running <TaskInstance: our_dag_with_python_operator_v06.greet manu
[2023-10-29, 11:22:35 UTC] {taskinstance.py:1662} INFO - Exporting env vars: AIRFLOW CTX DAG OWNER='coder2j' AIRFLOW CTX_I
[2023-10-29, 11:22:35 UTC] {logging_mixin.py:151} INFO - Hello World! My name is Jerry Fridman,and I am 19 years old!
[2023-10-29, 11:22:35 UTC] {python.py:194} INFO - Done. Returned value was: None
[2023-10-29, 11:22:35 UTC] {taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=our_dag_with_python_operator_v06,
[2023-10-29, 11:22:35 UTC] {local_task_job_runner.py:228} INFO - Task exited with return code 0
[2023-10-29, 11:22:35 UTC] {taskinstance.py:2778} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

# 查看xcom

<input type="checkbox"/>	Key ↕	Value ↕	Timestamp ↕	Dag Id ↕	Task Id ↕
<input type="checkbox"/> 	last_name	Fridman	2023-10-29, 11:22:33	our_dag_with_python_operator_v06	get_name ▼
<input type="checkbox"/> 	first_name	Jerry	2023-10-29, 11:22:33	our_dag_with_python_operator_v06	get_name ▼
<input type="checkbox"/> 	age	19	2023-10-29, 11:22:33	our_dag_with_python_operator_v06	get_age ▼
<input type="checkbox"/> 	last_name	Fridman	2023-10-29, 11:22:33	our_dag_with_python_operator_v06	get_name ▼
<input type="checkbox"/> 	first_name	Jerry	2023-10-29, 11:22:33	our_dag_with_python_operator_v06	get_name ▼
<input type="checkbox"/> 	age	19	2023-10-29, 11:22:33	our_dag_with_python_operator_v06	get_age ▼

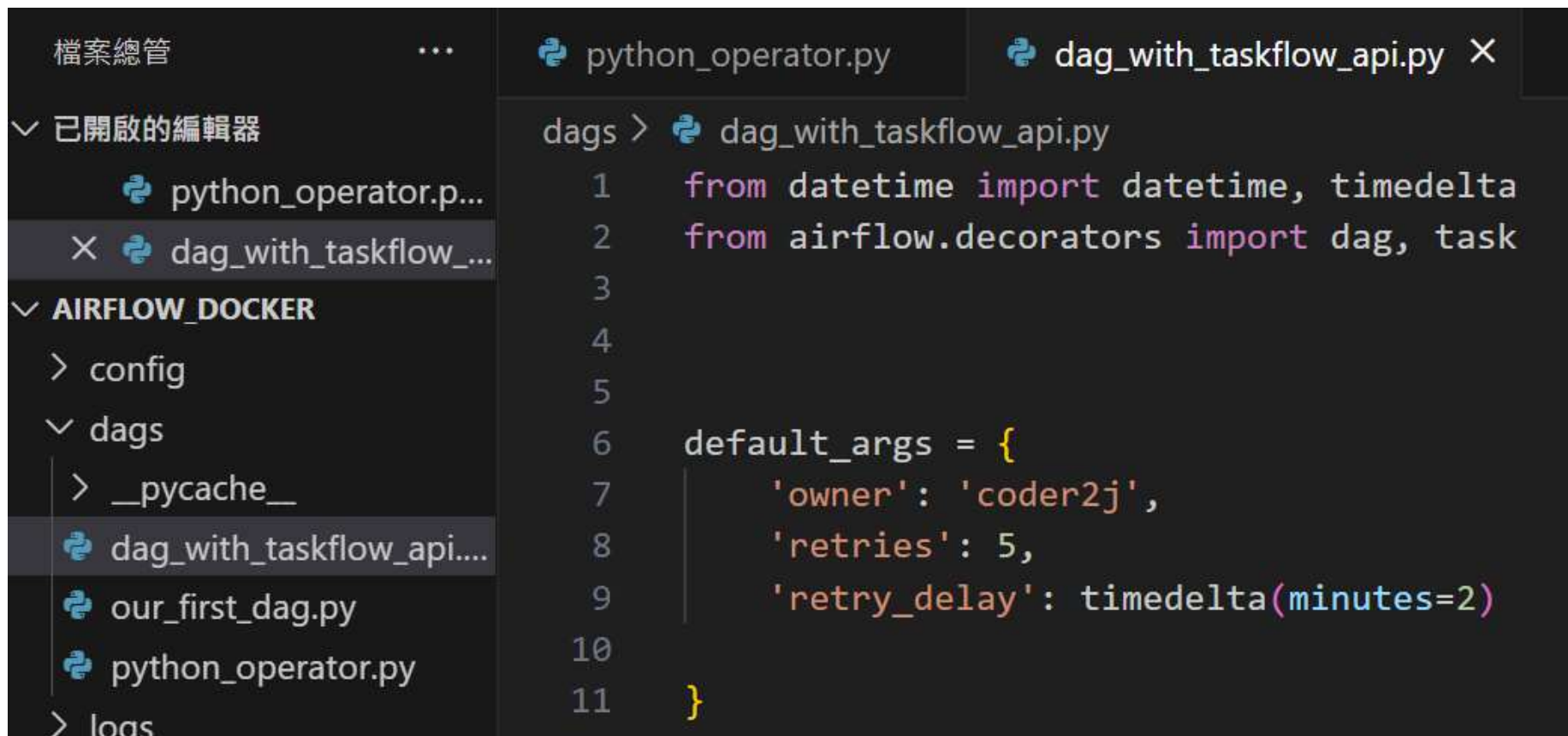
# Airflow Task Flow API

# 說明

- Airflow TaskFlow API 是 Apache Airflow 2.0 中引入的一種撰寫數據流水線的新範例。它是一種函數式 API，允許您顯式聲明消息傳遞，同時隱式聲明任務依賴關係。
- TaskFlow API 與傳統的 DAGs 撰寫範例相比，更加靈活，易於編寫和維護。它採用了一種基於任務的編程模型，其中任務是數據流水線中的基本構建塊。TaskFlow API 通過提供一組簡單的裝飾器和函數，使得編寫數據流水線變得更加容易。
- 簡單來說使用 TaskFlow API 撰寫會較於傳統 DAGs 更容易更清楚



# 新增一個 dag\_with\_taskflow\_api.py



```
檔案總管  ... python_operator.py dag_with_taskflow_api.py X
```

已開啟的編輯器

- python\_operator.p...
- X dag\_with\_taskflow\_...

AIRFLOW\_DOCKER

- > config
- > dags
  - > \_\_pycache\_\_
  - dag\_with\_taskflow\_api....
  - our\_first\_dag.py
  - python\_operator.py
- > logs

```
dags > dag_with_taskflow_api.py
1  from datetime import datetime, timedelta
2  from airflow.decorators import dag, task
3
4
5
6  default_args = {
7      'owner': 'coder2j',
8      'retries': 5,
9      'retry_delay': timedelta(minutes=2)
10
11 }
```

# 使用 taskflow api撰寫

```
12
13     @dag(dag_id='dag_with_taskflow_api_v01',
14           default_args=default_args,
15           start_date=datetime(2023, 10, 28),
16           schedule_interval='@daily')
17     def hello_world_etl():
18         @task()
19         def get_name():
20             return "Jerry"
21         @task()
22         def get_age():
23             return 19
24         @task()
25         def greet(name, age):
26             print(f"Hello World! My name is {name}"
27                   f"and I am {age} years old!")
28         name = get_name()
29         age = get_age()
30         greet(name=name, age=age)
31
32     greet_dag = hello_world_etl()
33
```

# 成功建立畫面

UI elements at the top:

- Buttons: All 65, Active 1, Paused 64, Running 0, Failed 0
- Text: Filter DAGs by tag

Table header:

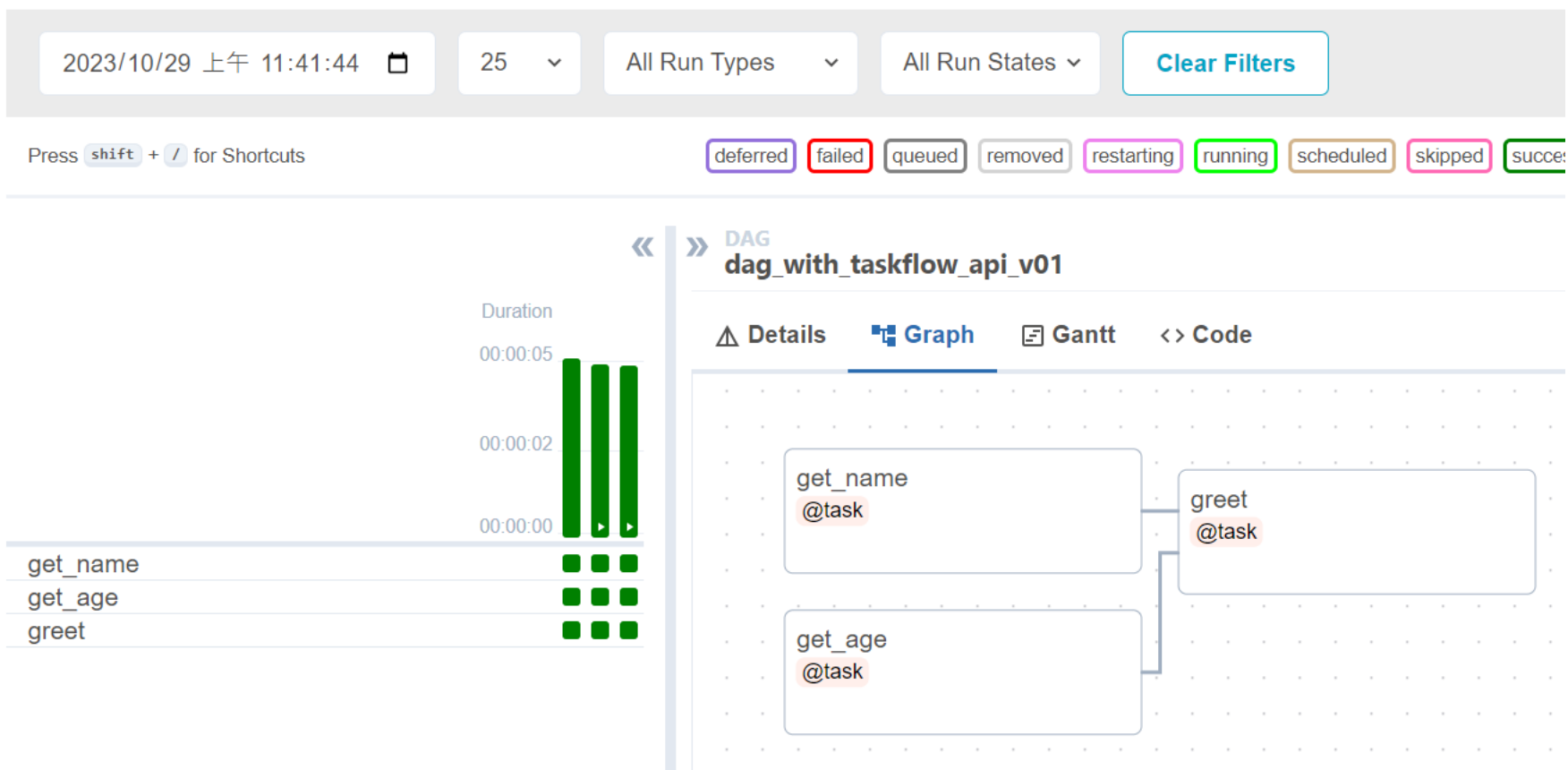
DAG	Owner	Runs
-----	-------	------

Table body:

<input type="checkbox"/> dag_with_taskflow_api_v01	coder2j	<div><div></div><div></div><div></div><div></div></div>
<input type="checkbox"/> dataset_consumes_1 consumes dataset-scheduled	airflow	<div><div></div><div></div><div></div><div></div></div>
<input type="checkbox"/> dataset_consumes_1_and_2 consumes dataset-scheduled	airflow	<div><div></div><div></div><div></div><div></div></div>

Tooltip: Show DAGs with failed latest DAG run





# 成功執行



# 查看logs

```
{standard_task_runner.py:85} INFO - Job 1263: Subtask greet
{task_command.py:416} INFO - Running <TaskInstance: dag_with_taskflow_api_v01.greet manual__2023-10-2
{taskinstance.py:1662} INFO - Exporting env vars: AIRFLOW CTX DAG OWNER='coder2j' AIRFLOW CTX DAG ID=
{logging_mixin.py:151} INFO - Hello World! My name is Jerryand I am 19 years old!
{python.py:194} INFO - Done. Returned value was: None
{taskinstance.py:1400} INFO - Marking task as SUCCESS. dag_id=dag_with_taskflow_api_v01, task_id=gree
{local_task_job_runner.py:228} INFO - Task exited with return code 0
{taskinstance.py:2778} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

# 查看xcom

<input type="checkbox"/>	Key ↕	Value ↕	Timestamp ↕	Dag Id ↕	Task Id ↕
<input type="checkbox"/> 	return_value	Jerry	2023-10-29, 11:41:46	dag_with_taskflow_api_v01	get_name ▼
<input type="checkbox"/> 	return_value	19	2023-10-29, 11:41:46	dag_with_taskflow_api_v01	get_age ▼
<input type="checkbox"/> 	return_value	Jerry	2023-10-29, 11:41:46	dag_with_taskflow_api_v01	get_name ▼
<input type="checkbox"/> 	return_value	19	2023-10-29, 11:41:46	dag_with_taskflow_api_v01	get_age ▼

# 撰寫更複雜的taskflow api

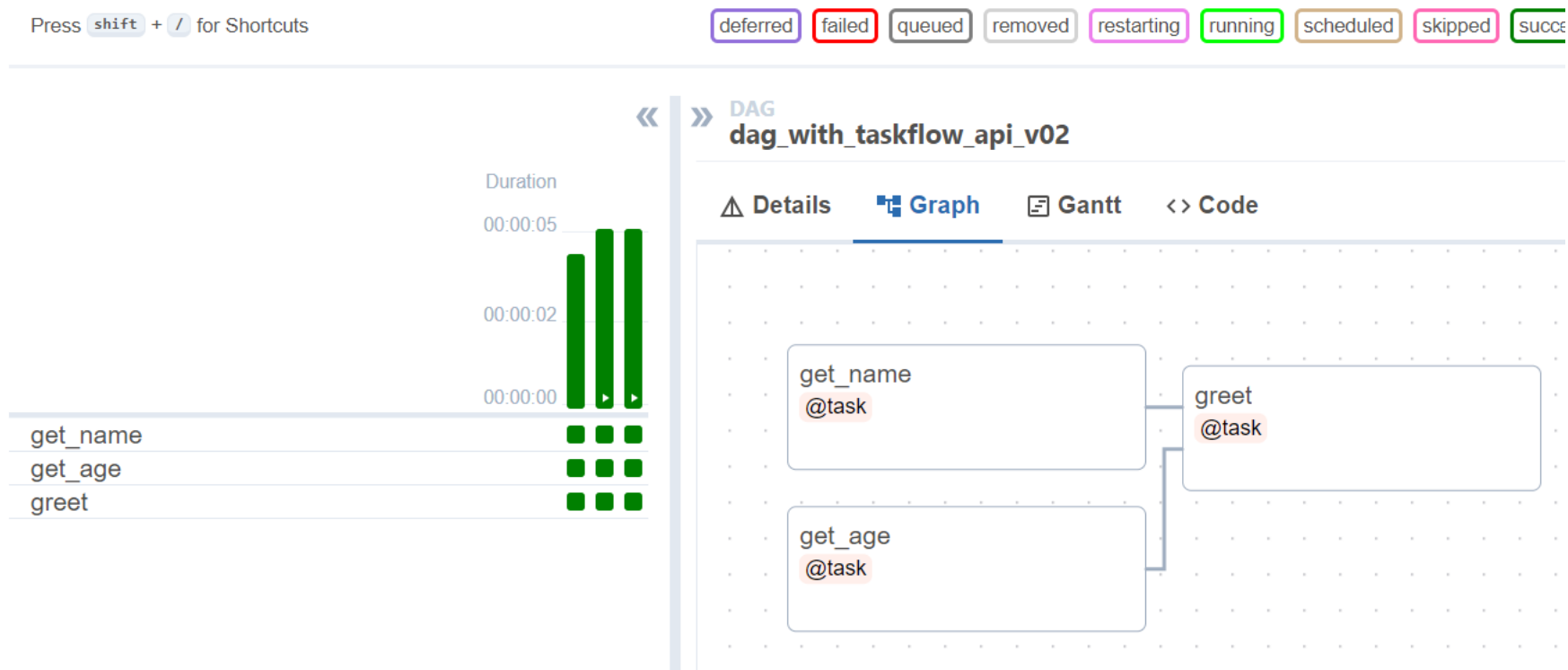
```
13  @dag(dag_id='dag_with_taskflow_api_v02',
14      |  default_args=default_args,
15      |  start_date=datetime(2023, 10, 28),
16      |  schedule_interval='@daily')
17  def hello_world_etl():
18      |  @task(multiple_outputs=True)
19      |  def get_name():
20      |  |  return {
21      |  |  |  'first_name': 'Jerry',
22      |  |  |  'last_name': 'Fridman'
23      |  |  }
```



```
24     @task()
25     def get_age():
26         return 19
27     @task()
28     def greet(first_name, last_name, age):
29         print(f"Hello World! My name is {first_name} {last_name}"
30               f"and I am {age} years old!")
31     name_dict = get_name()
32     age = get_age()
33     greet(first_name=name_dict['first_name'],
34           last_name=name_dict['last_name'],
35           age=age)
36
37 greet_dag = hello_world_etl()
38
```






# 成功啟動



# 查看logs

```
sk_runner.py:85} INFO - Job 1271: Subtask greet  
d.py:416} INFO - Running <TaskInstance: dag_with_taskflow_api_v02.greet manual__2023-10-  
e.py:1662} INFO - Exporting env vars: AIRFLOW_CTX_DAG_OWNER='coder2j' AIRFLOW_CTX_DAG_I  
n.py:151} INFO - Hello World! My name is Jerry Fridmanand I am 19 years old!  
94} INFO - Done. Returned value was: None  
e.py:1400} INFO - Marking task as SUCCESS. dag_id=dag_with_taskflow_api_v02, task_id=gre  
job_runner.py:228} INFO - Task exited with return code 0  
e.py:2778} INFO - 0 downstream tasks scheduled from follow-on schedule check
```

# 查看xcom

<input type="checkbox"/>	Key ↕	Value ↕	Timestamp ↕	Dag Id ↕	Task Id ↕
<input type="checkbox"/>	 return_value	{'first_name': 'Jerry', 'last_name': 'Fridman'}	2023-10-29, 11:51:15	dag_with_taskflow_api_v02	get_name ▼
<input type="checkbox"/>	 last_name	Fridman	2023-10-29, 11:51:15	dag_with_taskflow_api_v02	get_name ▼
<input type="checkbox"/>	 first_name	Jerry	2023-10-29, 11:51:15	dag_with_taskflow_api_v02	get_name ▼

# 可觀察到兩者程式碼行數有明顯差異

使用原本DAGs撰寫

```
39     task2 = PythonOperator(  
40         task_id='get_name',  
41         python_callable=get_name  
42     )  
43  
44     task3 = PythonOperator(  
45         task_id='get_age',  
46         python_callable=get_age  
47     )  
48  
49     # task1  
50     # task2 >> task1  
51     [task2, task3] >> task1
```

使用taskflow API撰寫

```
dags > dag_with_taskflow_api.py  
27     @task()  
28     def greet(first_name, last_name, age):  
29         print(f"Hello World! My name is {first_name} {last_name}"  
30             f"and I am {age} years old!")  
31     name_dict = get_name()  
32     age = get_age()  
33     greet(first_name=name_dict['first_name'],  
34           last_name=name_dict['last_name'],  
35           age=age)  
36  
37     greet_dag = hello_world_etl()  
38
```

# Airflow Catch-Up and Backfill

# 說明

- Airflow 中的 Catch-Up 和 Backfill 是兩個不同的概念。
- Catch-Up:  
是指在啟用 DAG 時，將所有過去的 DAG Runs 都執行一次，以便追上當前時間點。
- Backfill:  
則是指手動重新執行過去某個時間段內的 DAG Runs。

# 再新增一個.py

```
python_operator.py  dag_with_catchup_and_backfill.py X
dags > dag_with_catchup_and_backfill.py
1  from datetime import datetime, timedelta
2  from airflow import DAG
3  from airflow.operators.bash import BashOperator
4
5  default_args = {
6      'owner': 'coder2j',
7      'retries': 5,
8      'retry_delay': timedelta(minutes=2)
9
10 }
```

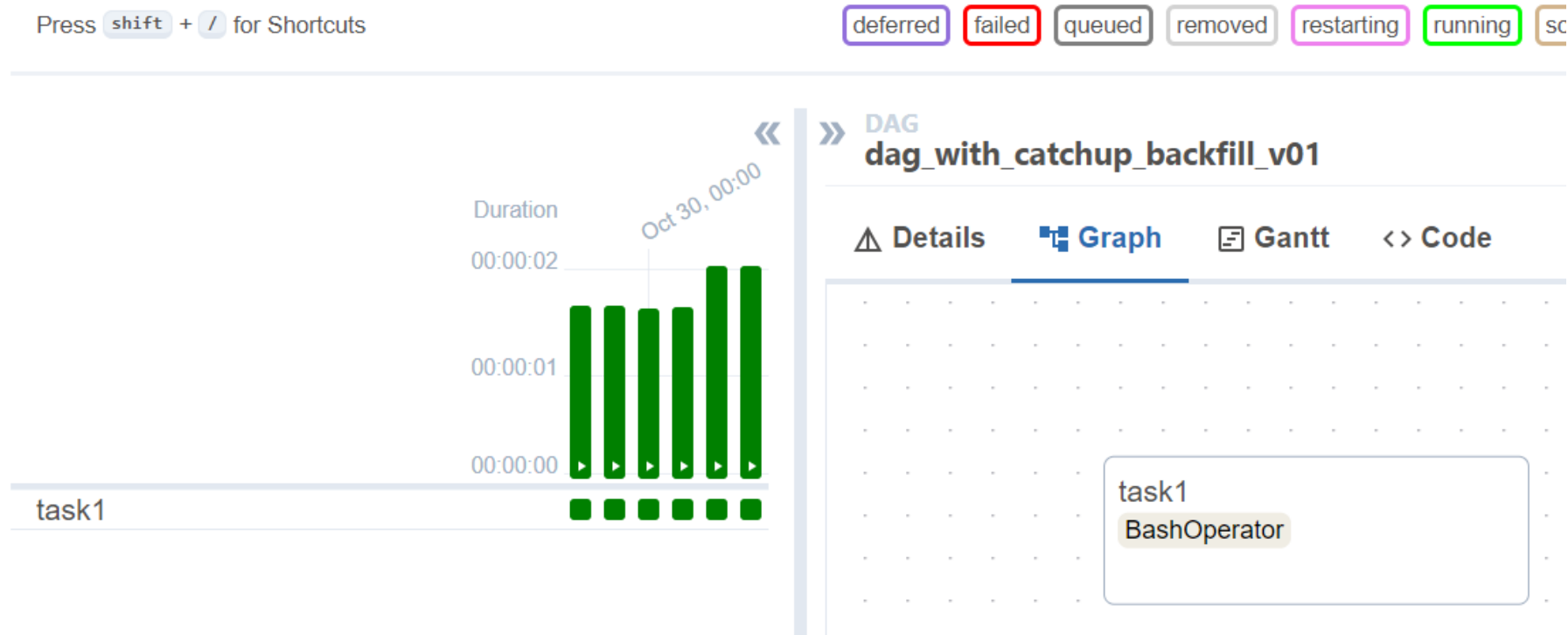
```
11
12     with DAG(
13         dag_id='dag_with_catchup_backfill_v01',
14         default_args=default_args,
15         start_date=datetime(2023, 10, 30),
16         schedule_interval='@daily',
17         catchup=True
18     ) as dag:
19         task1 = BashOperator(
20             task_id='task1',
21             bash_command='echo This is a simple bash command!'
22         )
```



# 成功建立畫面

DAG 	Owner 	Runs 
 dag_with_catchup_backfill_v01	coder2j	   
 dag_with_taskflow_api_v02	coder2j	   
 dataset_consumes_1 consumes dataset-scheduled	airflow	   

# 成功執行



```
12  with DAG(  
13      dag_id='dag_with_catchup_backfill_v02',  
14      default_args=default_args,  
15      start_date=datetime(2023, 10, 30),  
16      schedule_interval='@daily',  
17      catchup=False  
18  ) as dag:  
19      task1 = BashOperator(  
20          task_id='task1',  
21          bash_command='echo This is a simple bash command!'  
22      )
```

 DAG 	Owner 	Runs 
<input type="checkbox"/> dag_with_catchup_backfill_v01	coder2j	 6  
<input type="checkbox"/> dag_with_catchup_backfill_v02	coder2j	   
<input type="checkbox"/> dag_with_taskflow_api_v02	coder2j	 4  
<input type="checkbox"/> dataset_consumes_1 <span>consumes</span> <span>dataset-scheduled</span>	airflow	   

# 成功執行

for Shortcuts

deferred

failed

queued

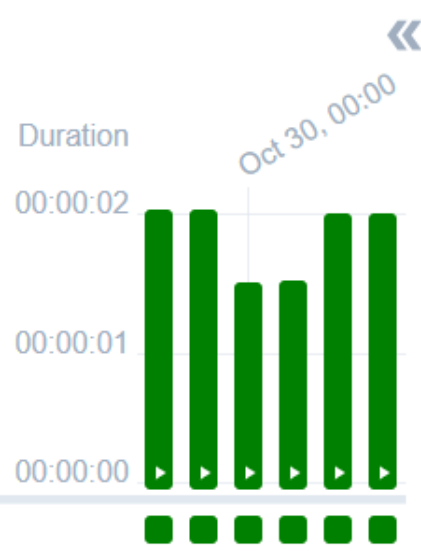
removed

restarting

running

scheduled

skipped



» DAG  
dag\_with\_catchup\_backfill\_v02

Details

Graph

Gantt

Code

task1  
BashOperator

# 查看執行中的container

```
PS C:\Users\jerry> cd "C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker"
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> docker ps
```

CONTAINER ID	IMAGE	COMMAND	CREATED	STATUS	PORTS
NAMES					
16896a554b5f	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	19 minutes ago	Up 18 minutes (healthy)	8080/tcp
airflow_docker-airflow-scheduler-1					
8a76374c9b9f	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	19 minutes ago	Up 18 minutes (healthy)	8080/tcp
airflow_docker-airflow-triggerer-1					
75878a44a270	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	19 minutes ago	Up 18 minutes (healthy)	0.0.0.0:8080->
8080/tcp airflow_docker-airflow-webserver-1					
94b67b9093cd	apache/airflow:2.7.2	"/usr/bin/dumb-init ..."	19 minutes ago	Up 18 minutes (healthy)	8080/tcp
airflow_docker-airflow-worker-1					
d65609074636	redis:latest	"docker-entrypoint.s..."	19 minutes ago	Up 19 minutes (healthy)	6379/tcp
airflow_docker-redis-1					
f82968dfffe8e	postgres:13	"docker-entrypoint.s..."	19 minutes ago	Up 19 minutes (healthy)	5432/tcp
airflow_docker-postgres-1					

Command:

Docker exec -it <container name> bash

```
PS C:\Users\jerry\Desktop\mastercourse\dataEngineer\airflow\airflow_docker> docker exec -it 16896a554b5f bash
airflow@16896a554b5f:/opt/airflow$ airflow dags backfill -s 2023-10-30 -e 2023-11-08 dag_with_catchup_backfill_v02
/home/airflow/.local/lib/python3.8/site-packages/airflow/cli/commands/dag_command.py:129 RemovedInAirflow3Warning: --ignore-first-depends-on-past is deprecated as the value is always set to True
[2023-10-30T00:42:25.366+0000] {dagbag.py:536} INFO - Filling up the DagBag from /opt/airflow/dags
/home/airflow/.local/lib/python3.8/site-packages/airflow/models/dagbag.py:342 RemovedInAirflow3Warning: Param `schedule_interval` is deprecated and will be removed in a future release. Please use `schedule` instead.
[2023-10-30T00:42:26.232+0000] {executor_loader.py:117} INFO - Loaded executor: CeleryExecutor
[2023-10-30T00:42:26.610+0000] {base_executor.py:144} INFO - Adding to queue: ['airflow', 'tasks', 'run', 'dag_with_catchup_backfill_v02', 'task1', 'backfill__2023-10-30T00:00:00+00:00', '--depends-on-past', 'ignore', '--local', '--pool', 'default_pool', '--subdir', 'DAGS_FOLDER/dag_with_catchup_and_backfill.py']
[2023-10-30T00:42:31.423+0000] {backfill_job_runner.py:412} INFO - [backfill progress] | finished run 0 of 10 | tasks waiting: 9 | succeeded: 0 | running: 1 | failed: 0 | skipped: 0 | deadlocked: 0 | not ready: 9
[2023-10-30T00:42:36.361+0000] {dagrun.py:653} INFO - Marking run <DagRun dag_with_catchup_backfill_v02 @ 2023-10-30T00:00:00+00:00: backfill__2023-10-30T00:00:00+00:00, state:running, queued_at: None. externally triggered: False> successful
1
```

# 成功執行Catch-Up畫面(1/2)

These tasks have succeeded:

DAG ID	Task ID	Run ID	Try number
dag_with_catchup_backfill_v02	task1	backfill__2023-10-30T00:00:00+00:00	1

These tasks are running:

DAG ID	Task ID	Run ID	Try number
--------	---------	--------	------------

These tasks have failed:

DAG ID	Task ID	Run ID	Try number
--------	---------	--------	------------

These tasks are skipped:

DAG ID	Task ID	Run ID	Try number
--------	---------	--------	------------



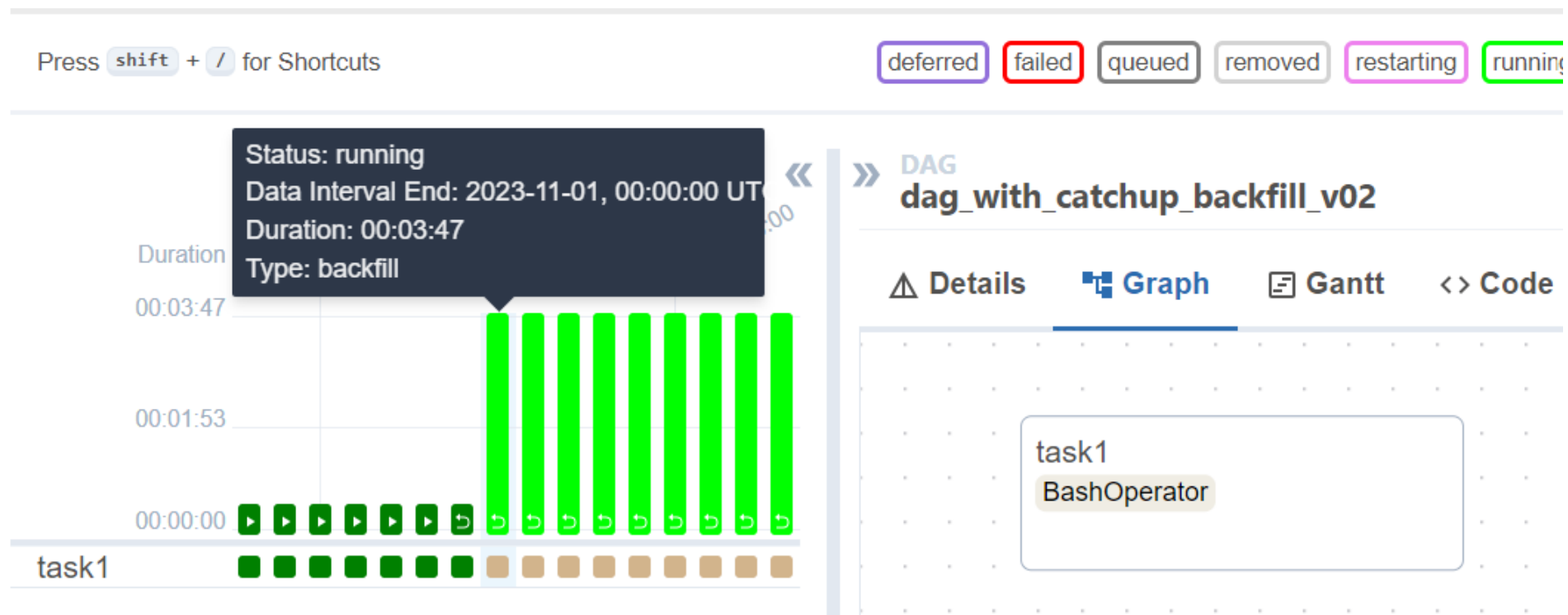
# 成功執行Catch-Up畫面(2/2)

These tasks are deadlocked:

DAG ID	Task ID	Run ID	Try number
dag_with_catchup_backfill_v02	task1	backfill__2023-10-31T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-01T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-02T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-03T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-04T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-05T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-06T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-07T00:00:00+00:00	1
dag_with_catchup_backfill_v02	task1	backfill__2023-11-08T00:00:00+00:00	1

airflow@16896a554b5f:/opt/airflow\$

# 成功執行



# Airflow Scheduler with Cron Expression

# 說明

- Airflow 中的 Cron Expression 是一種用於調度 DAG 執行的表達式。它是一種基於時間的表達式，可以指定 DAG 的運行時間。Cron Expression 由五個或六個字段組成，分別表示分鐘、小時、日期、月份、星期和年份（可選）。
- Airflow 中的 Cron Expression 與傳統的 Cron 表達式略有不同。Airflow 的 Cron Expression 支援秒級精度，而傳統的 Cron 表達式只支持分鐘級別的精度。
- 此外，Airflow 的 Cron Expression 還支援使用 @yearly、@monthly、@weekly、@daily、@hourly 等預定義的時間間隔來指定 DAG 的運行時間。

# 新增一個.py

```
python_operator.py  dag_with_cron_expression.py ×
dags > dag_with_cron_expression.py
1  from datetime import datetime, timedelta
2  from airflow import DAG
3  from airflow.operators.bash import BashOperator
4
5  default_args = {
6      'owner': 'coder2j',
7      'retries': 5,
8      'retry_delay': timedelta(minutes=2)
9
10 }
```

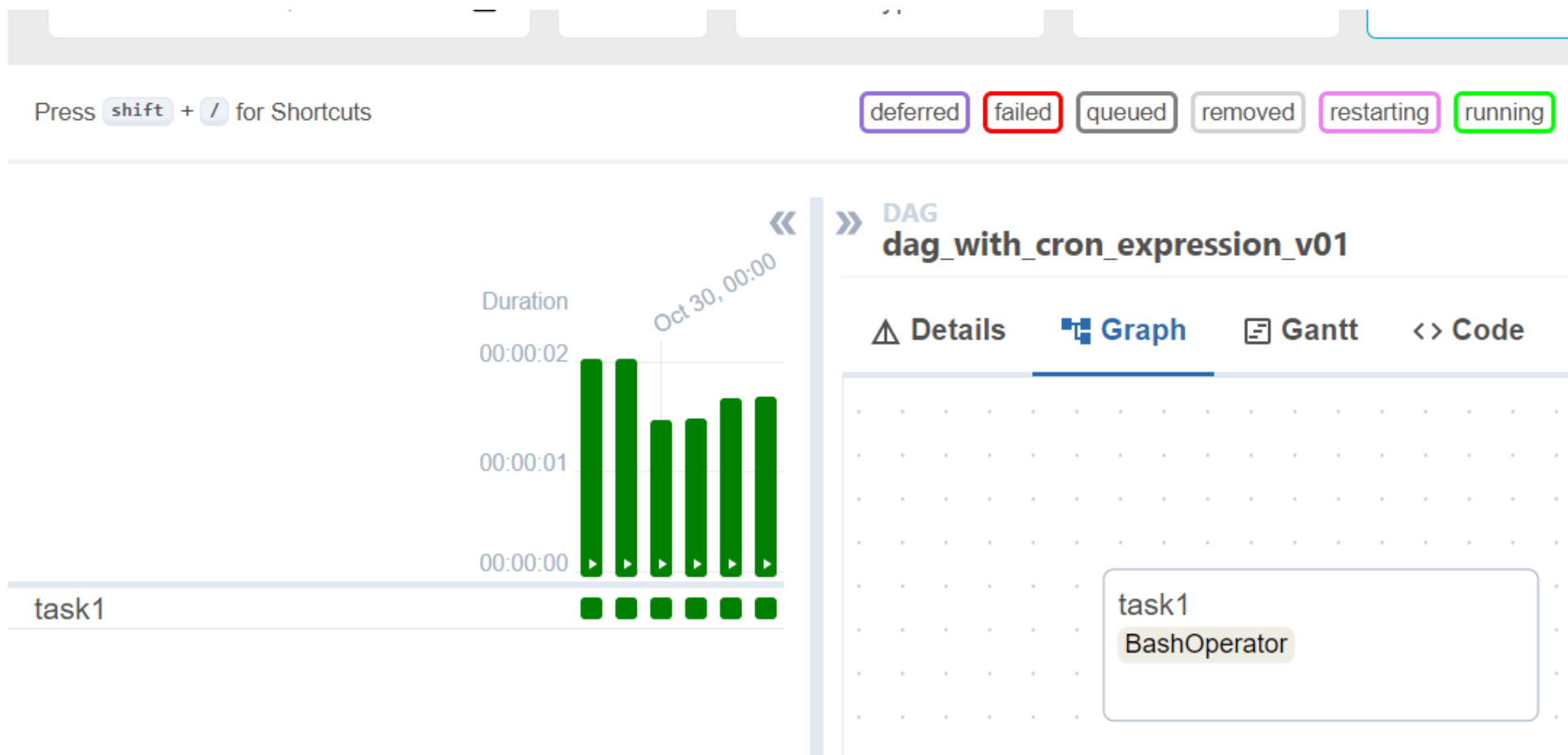
# 使用 Cron Expression

```
11
12     with DAG(
13         dag_id='dag_with_cron_expression_v01',
14         default_args=default_args,
15         start_date=datetime(2023, 10, 30),
16         schedule_interval='@daily'
17     ) as dag:
18         task1 = BashOperator(
19             task_id='task1',
20             bash_command='echo dag with cron expression!'
21         )
22     task1
```

# 成功建立畫面

All 69		Active 2	Paused 67	Running 9	Failed 0	Filter DAGs by tag	
DAG		Owner	Runs				
<input type="checkbox"/> dag_with_catchup_backfill_v02		coder2j	<input type="checkbox"/>	<input checked="" type="checkbox"/> 7	<input checked="" type="checkbox"/> 9	<input type="checkbox"/>	
<input type="checkbox"/> dag_with_cron_expression_v01		coder2j	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
<input type="checkbox"/> dag_with_taskflow_api_v02		coder2j	<input type="checkbox"/>	<input checked="" type="checkbox"/> 4	<input type="checkbox"/>	<input type="checkbox"/>	

# 成功執行





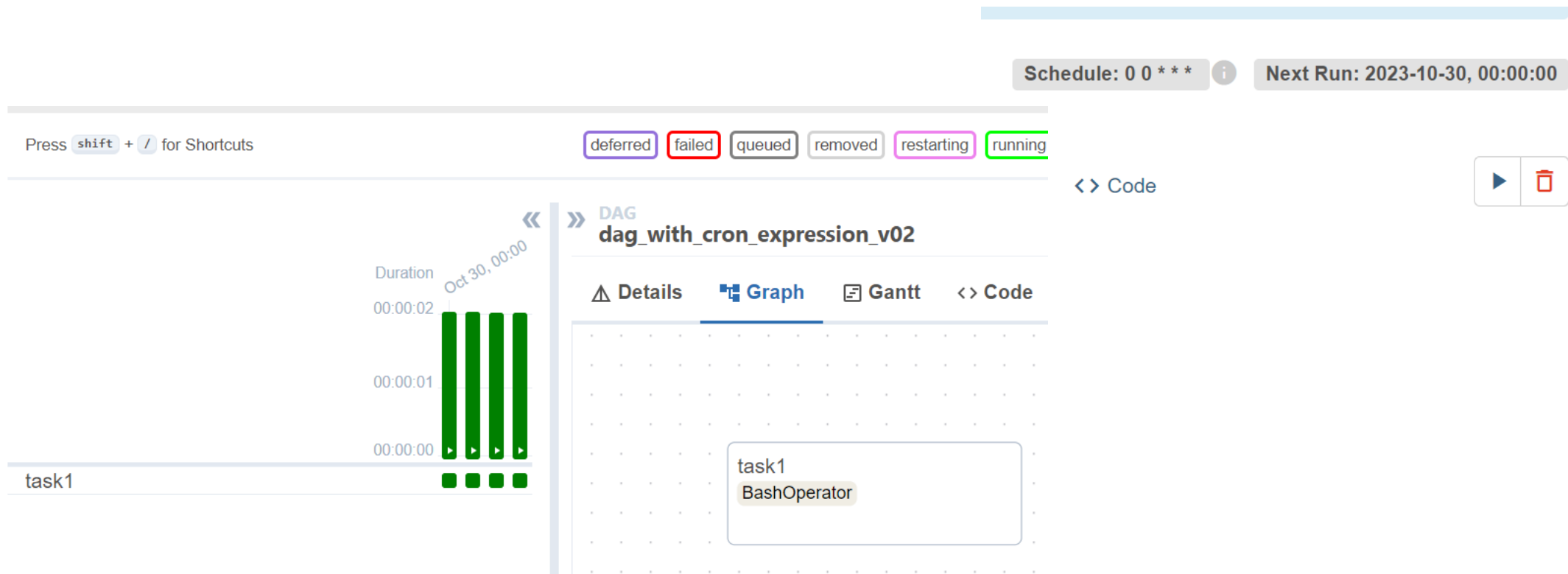
# 修改 schedule\_interval

```
12     with DAG(  
13         dag_id='dag_with_cron_expression_v02',  
14         default_args=default_args,  
15         start_date=datetime(2023, 10, 30),  
16         schedule_interval='0 0 * * *'  
17     ) as dag:  
18         task1 = BashOperator(  
19             task_id='task1',  
20             bash_command='echo dag with cron expression!'  
21         )  
22     task1
```

# 成功建立畫面

All 70		Active 3	Paused 67	Running 9	Failed 0	Filter DAGs by tag	
DAG		Owner	Runs				
<input type="checkbox"/> dag_with_catchup_backfill_v02		coder2j	<input type="checkbox"/>	<input checked="" type="checkbox"/> 7	<input checked="" type="checkbox"/> 9	<input type="checkbox"/>	
<input type="checkbox"/> dag_with_cron_expression_v01		coder2j	<input type="checkbox"/>	<input checked="" type="checkbox"/> 6	<input type="checkbox"/>	<input type="checkbox"/>	
<input type="checkbox"/> dag_with_cron_expression_v02		coder2j	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	
<input type="checkbox"/> dag_with_taskflow_api_v02		coder2j	<input type="checkbox"/>	<input checked="" type="checkbox"/> 4	<input type="checkbox"/>	<input type="checkbox"/>	

# 成功執行



**End**