1



File path: Documents › MIT › project-23 ›    Search project-23

| Name | Status | Date modified | Type | Size |
|---|---|---|---|---|
| airflow-docker | ↻ | 2/13/2023 7:47 AM | File folder | |
| code visualization | ↻ | 2/13/2023 7:48 AM | File folder | |
| assignment | ↻ | 2/13/2023 7:51 AM | Python Source File | 0 KB |

2

assignment.py ×

```python
assignment.py
1   from airflow import DAG
2   from datetime import timedelta
3
4   from  airflow.operators.bash import BashOperator
5   from airflow.utils.dates import days_ago
6   from airflow.operators.python import PythonOperator
7
8   import urllib.request
9   import time
10  import glob
11  import os
12  import json
13
14  def catalog():
15
16      def pull(url):
17          pass
18
19      def store(data, file):
20          pass
```

3

```python
14  def catalog():
15
16      def pull(url):
17          response = urllib.request.urlopen(url).read()
18          data = response.decode('utf-8')
19          return data
20
21      def store(data, file):
22          with open("data/" + file, 'w+') as file:
23              file.write(data)
```

4

```python
14 ∨ def catalog():
15
16 ∨     def pull(url):
17             response = urllib.request.urlopen(url).read()
18             data = response.decode('utf-8')
19             return data
20
21 ∨     def store(data, file):
22 ∨         with open("data/" + file, 'w+') as file:
23                 file.write(data)
24
25
26 ∨     with open("00_urls.txt", 'r') as file:
27         lines = file.readlines()
28         urls = [line.strip() for line in lines]
29
30 ∨     for url in urls:
31         data = pull(url)
32
33         index = url.rfind('/') + 1
34         file = url[index:]
35         store(data, file)
36
37         print('pulled: ' + file)
38         print('--- waiting ---')
39         time.sleep(15)
```

5

```python
41 ∨ def combine():
42 ∨     with open('combo.txt', 'w+') as outfile:
43 ∨         for file in glob.glob("data\*.html"):
44 ∨             with open(file) as infile:
45                 outfile.write(infile.read())
```

6

```python
47 ∨ def titles():
48        from bs4 import BeautifulSoup
49
50 ∨      def store_json(data,file):
51 ∨          with open(file, 'w', encoding='utf-8') as f:
52                  json.dump(data, f, ensure_ascii=False, indent=4)
53                  print('wrote file: ' + file)
54
55 ∨      with open('combo.txt', 'r') as html:
56
57              html = html.read().replace('\n', ' ').replace('\r', '')
58
59              #the following creates an html parser
60              soup = BeautifulSoup(html, "html.parser")
61              results = soup.find_all('h3')
62              titles = []
63
64              # tag inner text
65 ∨            for item in results:
66                  titles.append(item.text)
67
68              store_json(titles, 'titles.json')
```

7

```python
70 ∨ def clean():
71 ∨      def store_json(data,file):
72 ∨          with open(file, 'w', encoding='utf-8') as f:
73                  json.dump(data, f, ensure_ascii=False, indent=4)
74                  print('wrote file: ' + file)
75
76 ∨      with open("titles.json") as file:
77              titles = json.load(file)
78
79              # remove punctuation/numbers
80 ∨            for index, title in enumerate(titles):
81                  punctuation= '''!()-[]{};:'"\,<>./?@#$%^&*_~1234567890'''
82                  translationTable= str.maketrans("","",punctuation)
83                  clean = title.translate(translationTable)
84                  titles[index] = clean
85
86              # remove one character words
87 ∨            for index, title in enumerate(titles):
88                  clean = ' '.join( [word for word in title.split() if len(word)>1] )
89                  titles[index] = clean
90
91              store_json(titles, 'titles_clean.json')
```

8

```python
93    def count_words():
94        from collections import Counter
95
96        def store_json(data,file):
97            with open(file, 'w', encoding='utf-8') as f:
98                json.dump(data, f, ensure_ascii=False, indent=4)
99                print('wrote file: ' + file)
100
101        with open("titles_clean.json") as file:
102            titles = json.load(file)
103            words = []
104
105            # extract words and flatten
106            for title in titles:
107                words.extend(title.split())
108
109            # count word frequency
110            counts = Counter(words)
111            store_json(counts, 'words.json')
```

9

```python
114 ∨ with DAG(
115       "assignment",
116      start_date=days_ago(1),
117      schedule_interval="@daily",catchup=False,
118  ) as dag:
119
120      # ts are tasks
121 ∨    t0 = BashOperator(
122          task_id='task_zero',
123          bash_command='pip install beautifulsoup4',
124          retries=2
125      )
126 ∨    t1 = PythonOperator(
127          task_id='task_one',
128          depends_on_past=False,
129          python_callable=catalog
130      )
131
132 ∨    t2 = PythonOperator(
133          task_id='task_two',
134          depends_on_past=False,
135          python_callable=combine
136      )
137
138 ∨    t3 = PythonOperator(
139          task_id='task_three',
140          depends_on_past=False,
141          python_callable=titles
142      )
143
144 ∨    t4 = PythonOperator(
145          task_id='task_four',
146          depends_on_past=False,
147          python_callable=clean
148      )
149
150 ∨    t5 = PythonOperator(
151          task_id='task_five',
152          depends_on_past=False,
153          python_callable=count_words
154      )
155
156      t0>>t1>>t2>>t3>>t4>>t5
```
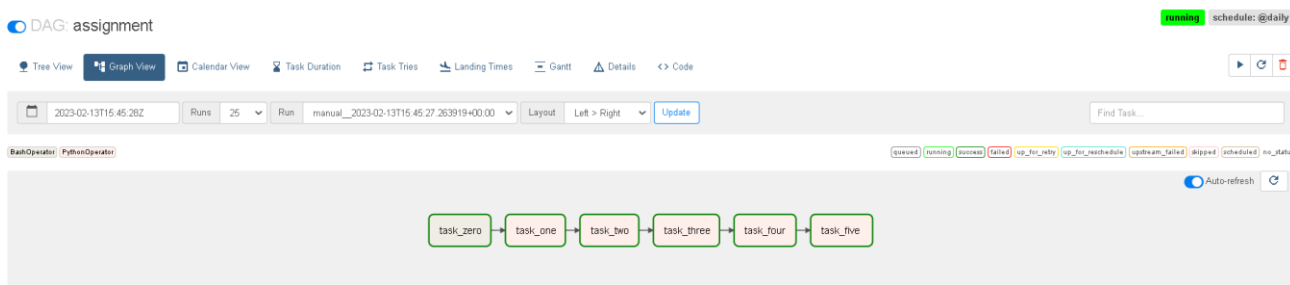
10

## Containers  Give Feedback

A container packages up code and its dependencies so the application runs quickly and reliably from one computing environment to another. Learn more

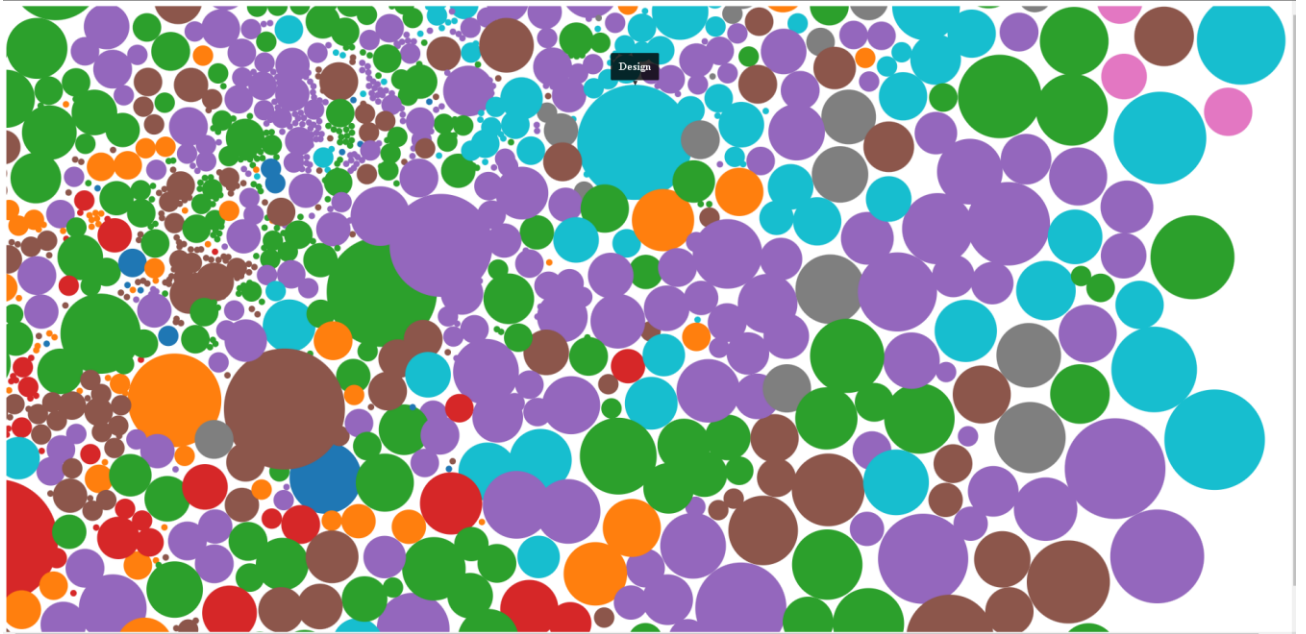Showing 8 items                                                                                                    🔍 Search            ⋮

| | NAME | IMAGE | STATUS | PORT(S) | STARTED ↓ | ACTIONS |
|---|---|---|---|---|---|---|
| ☐ ⊗ ⬙ | **airflow-docker**<br>7 containers | - | Running (7/7) | - | | ■ ⋮ 🗑 |
| ☐ 🟢 | **airflow-worker-1**<br>fa027b7d5394 📋 | apache/airflow:2.1.1 | Running | - | 3 minutes ago | ■ ⋮ 🗑 |
| ☐ 🟢 | **airflow-scheduler-1**<br>f625663d186c 📋 | apache/airflow:2.1.1 | Running | - | 3 minutes ago | ■ ⋮ 🗑 |
| ☐ 🟢 | **airflow-webserver-1**<br>3b4b36b556d1 📋 | apache/airflow:2.1.1 | Running | 8080 | 3 minutes ago | ■ ⋮ 🗑 |
| ☐ 🟢 | **flower-1**<br>d4e003d9af65 📋 | apache/airflow:2.1.1 | Running | 5555 | 3 minutes ago | ■ ⋮ 🗑 |
| ☐ 🟢 | **airflow-init-1**<br>3fbe9fff1b27 📋 | apache/airflow:2.1.1 | Running | - | 3 minutes ago | ■ ⋮ 🗑 |
| ☐ 🟢 | **postgres-1**<br>7a1d86c47b97 📋 | postgres:13 | Running | - | 3 minutes ago | ■ ⋮ 🗑 |
| ☐ 🟢 | **redis-1**<br>ab2788a02896 📋 | redis:latest | Running | 6379 | 3 minutes ago | ■ ⋮ 🗑 |

11

⬤ DAG: assignment                                                                    **running**  schedule: @daily

| ● Tree View | ▪▪ Graph View | 📅 Calendar View | ⏱ Task Duration | ⇄ Task Tries | ⤴ Landing Times | ☰ Gantt | ⚠ Details | <> Code |

| 📅 2023-02-13T15:45:28Z | Runs | 25 ⌄ | Run | manual__2023-02-13T15:45:27.263919+00:00 ⌄ | Layout | Left > Right ⌄ | Update |    Find Task... |

BashOperator  PythonOperator          queued running success failed up_for_retry up_for_reschedule upstream_failed skipped scheduled no_status

                                                                                                        ⬤ Auto-refresh  ↻

task_zero → task_one → task_two → task_three → task_four → task_five

12

```python
scores = {
    "Molecule": 1,
    "Builders": 1,
    "Engineering": 346,
    "Molecular": 54,
    "Marvels": 1,
    "Careers": 3,
    "and": 1187,
    "ChemE": 1,
    "at": 14,
    "MIT": 8,
    "Ethics": 13,
    "for": 190,
    "Engineers": 19,
    "Foundations": 13,
    "of": 418,
    "Entrepreneurship": 26,
    "Advances": 4,
    "in": 781,
    "Biomanufacturing": 8,
    "Philosophical": 2,
    "History": 36,
    "Energy": 83,
    "Foundational": 1,
    "Analyses": 3,
    "Problems": 22,
    "the": 139,
    "Environment": 21,
    "Advanced": 145,
    "Topics": 107,
    "Debating": 1,
    "About": 6,
    "Society": 13,
    "Cultural": 1,
    "Studies": 49,
    "Chemical": 48,
    "Graduate": 56,
    "Students": 2,
    "Models": 21,
```

13



14