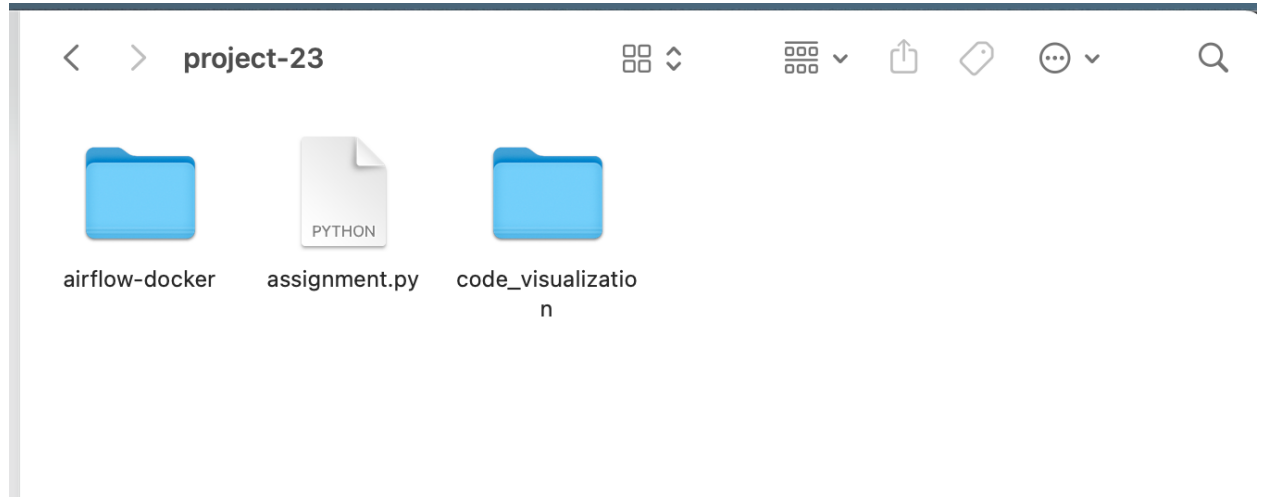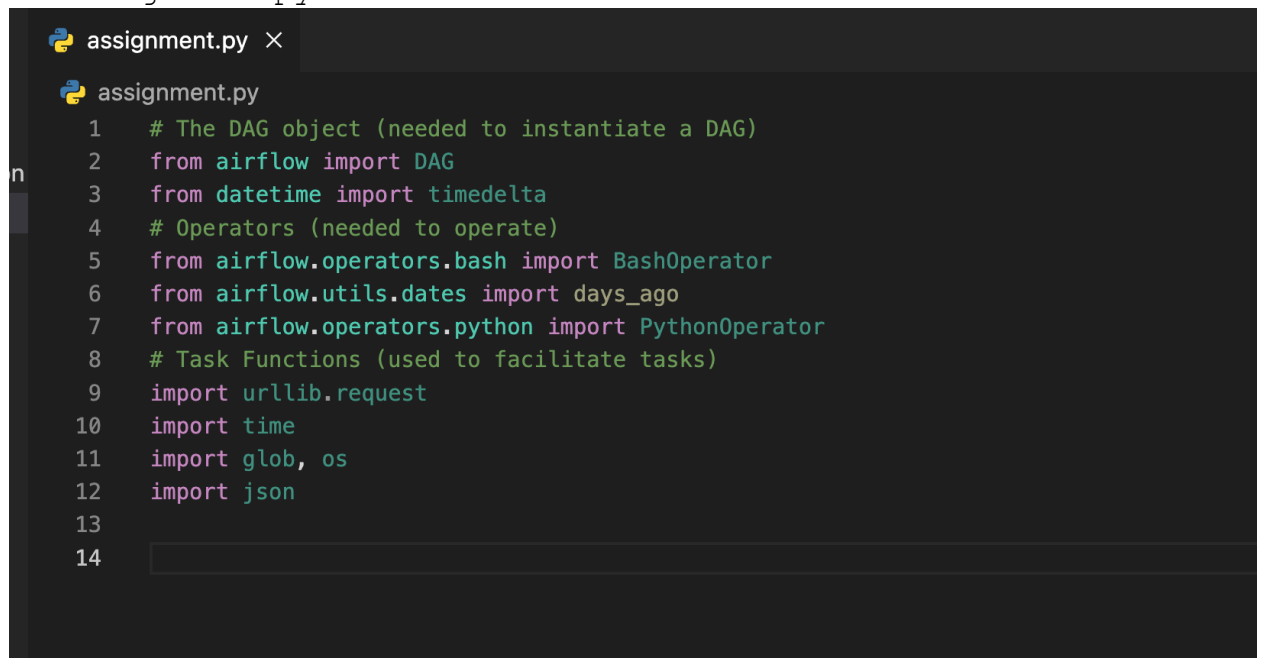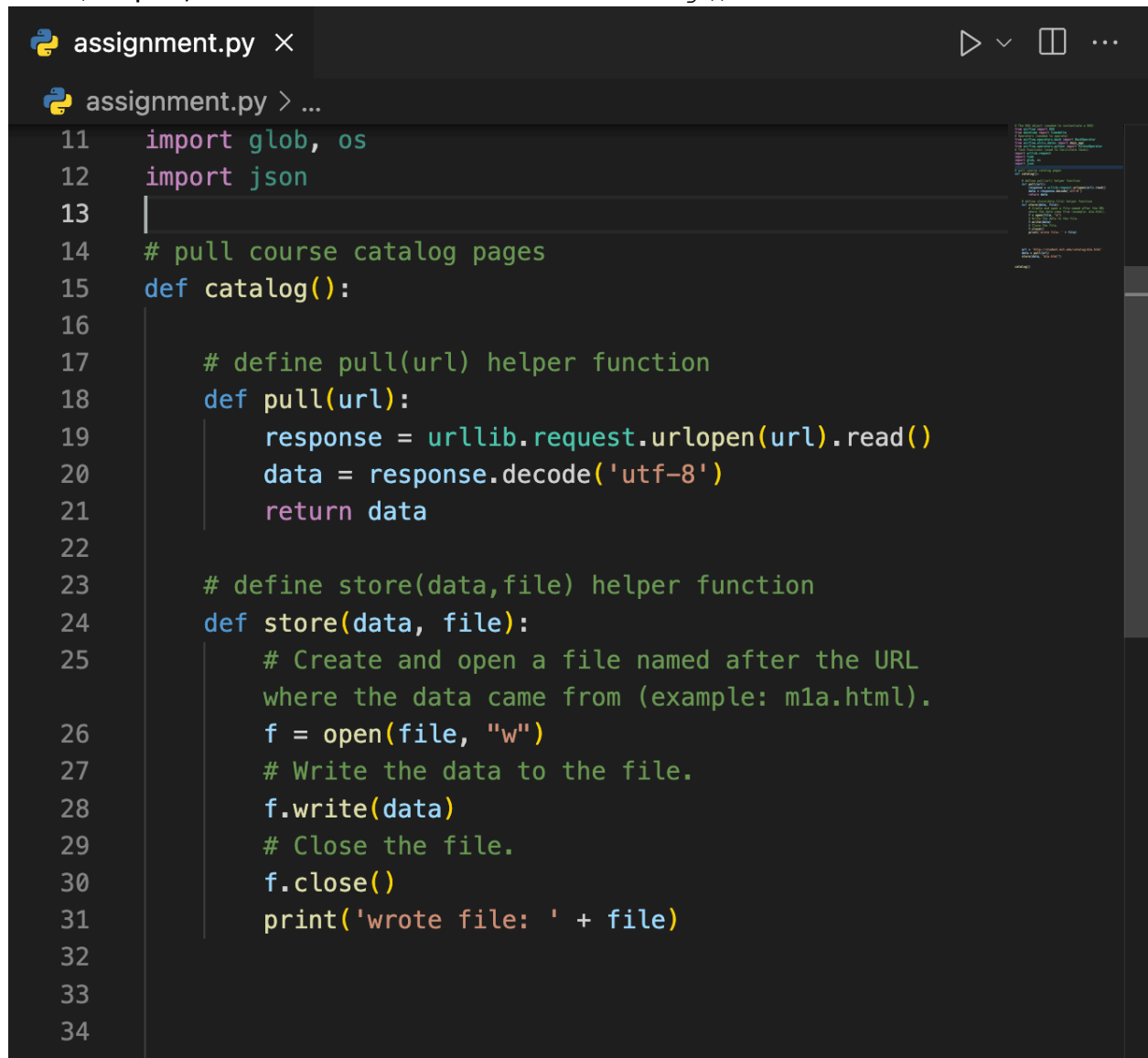**Part 1: Code Development**

1.  Provide a screenshot of the `project-23` folder with the `code visualization` folder, `airflow-docker` folder, and `assignment.py` file within it.



2.  Provide a screenshot to show that you have imported the DAG object, the *operators*, and all of the necessary *task functions* into the `assignment.py` file.



```python
# The DAG object (needed to instantiate a DAG)
from airflow import DAG
from datetime import timedelta
# Operators (needed to operate)
from airflow.operators.bash import BashOperator
from airflow.utils.dates import days_ago
from airflow.operators.python import PythonOperator
# Task Functions (used to facilitate tasks)
import urllib.request
import time
import glob, os
import json
```
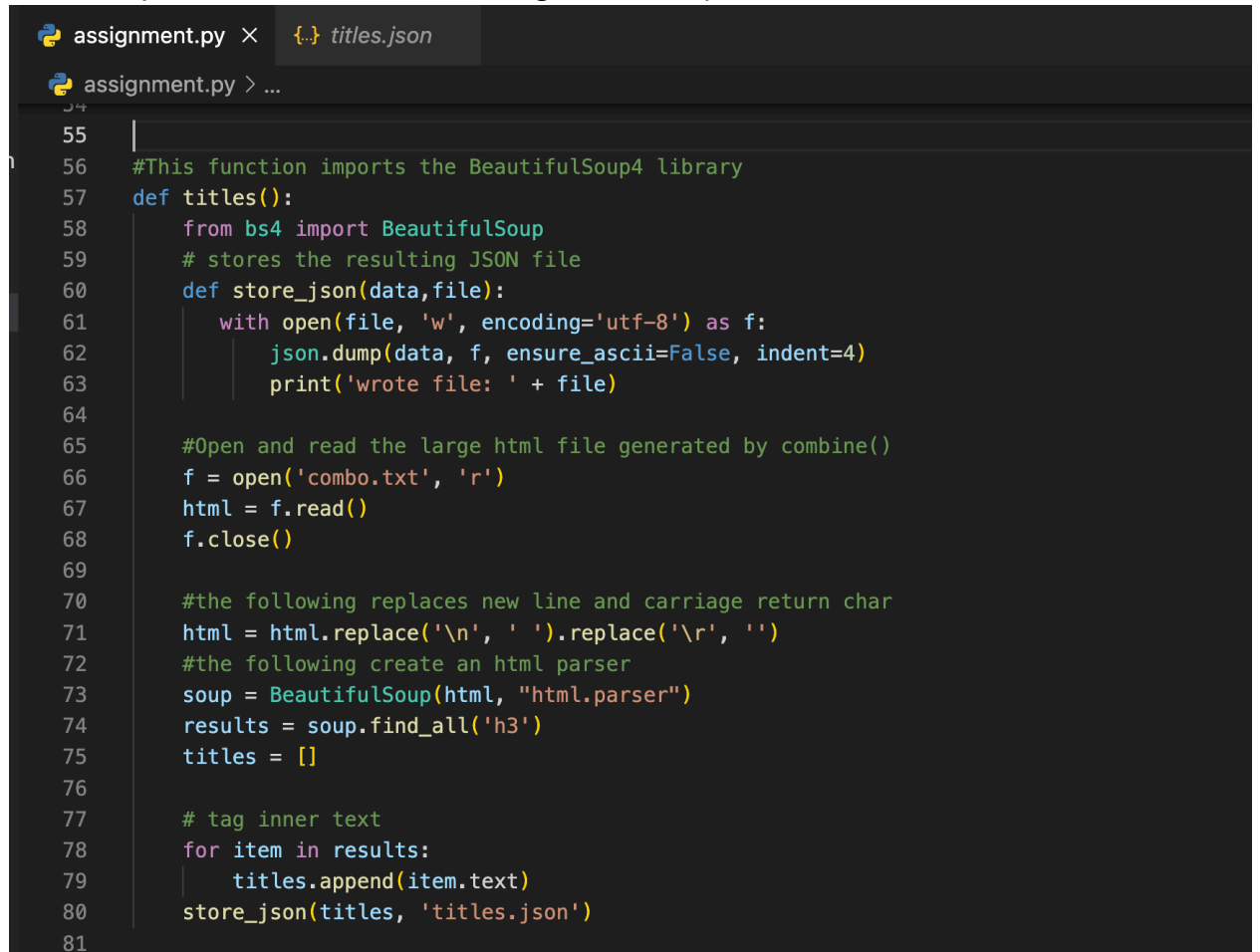
3. Provide a screenshot of the `pull(url)` and `store(data, file)` helper *functions* defined inside of the `catalog()` *task*.

```python
import glob, os
import json

# pull course catalog pages
def catalog():

    # define pull(url) helper function
    def pull(url):
        response = urllib.request.urlopen(url).read()
        data = response.decode('utf-8')
        return data

    # define store(data,file) helper function
    def store(data, file):
        # Create and open a file named after the URL
        # where the data came from (example: m1a.html).
        f = open(file, "w")
        # Write the data to the file.
        f.write(data)
        # Close the file.
        f.close()
        print('wrote file: ' + file)
```

4. Provide a screenshot of the entire `catalog()` *function*, including the `urls` *list* and the `for` *loop* that you just implemented.

```python
29          # Close the file.
30          f.close()
31          print('wrote file: ' + file)
32
33      urls = ['http://student.mit.edu/catalog/m1a.html','http://student.mit.edu/catalog/m1b.html','http://student.mit.edu/catalog/m1c.
        html','http://student.mit.edu/catalog/m2a.html','http://student.mit.edu/catalog/m2b.html','http://student.mit.edu/catalog/m2c.
        html','http://student.mit.edu/catalog/m3a.html','http://student.mit.edu/catalog/m3b.html','http://student.mit.edu/catalog/m4a.
        html','http://student.mit.edu/catalog/m4b.html','http://student.mit.edu/catalog/m4c.html','http://student.mit.edu/catalog/m4d.
        html','http://student.mit.edu/catalog/m4e.html','http://student.mit.edu/catalog/m4f.html','http://student.mit.edu/catalog/m4g.
        html','http://student.mit.edu/catalog/m5a.html','http://student.mit.edu/catalog/m5b.html','http://student.mit.edu/catalog/m6a.
        html','http://student.mit.edu/catalog/m6b.html','http://student.mit.edu/catalog/m6c.html','http://student.mit.edu/catalog/m7a.
        html','http://student.mit.edu/catalog/m8a.html','http://student.mit.edu/catalog/m8b.html','http://student.mit.edu/catalog/m9a.
        html','http://student.mit.edu/catalog/m9b.html','http://student.mit.edu/catalog/m10a.html','http://student.mit.edu/catalog/m10b.
        html','http://student.mit.edu/catalog/m11a.html','http://student.mit.edu/catalog/m11b.html','http://student.mit.edu/catalog/m11c.
        html','http://student.mit.edu/catalog/m12a.html','http://student.mit.edu/catalog/m12b.html','http://student.mit.edu/catalog/m12c.
        html','http://student.mit.edu/catalog/m14a.html','http://student.mit.edu/catalog/m14b.html','http://student.mit.edu/catalog/m15a.
        html','http://student.mit.edu/catalog/m15b.html','http://student.mit.edu/catalog/m15c.html','http://student.mit.edu/catalog/m16a.
        html','http://student.mit.edu/catalog/m16b.html','http://student.mit.edu/catalog/m18a.html','http://student.mit.edu/catalog/m18b.
        html','http://student.mit.edu/catalog/m20a.html','http://student.mit.edu/catalog/m22a.html','http://student.mit.edu/catalog/m22b.
        html','http://student.mit.edu/catalog/m22c.html']
34
35      for url in urls:
36          index = url.rfind('/') + 1
37          #call pull function
38          data = pull(url)
39          file = url[index:]
40          #call store function
41          store(data, file)
42          print('pulled: ' + file)
43          print('--- waiting ---')
44          time.sleep(15)
45
```

5. Provide a screenshot of the `combine()` *method* with the correct code to combine the files.

```python
44              print('--- waiting ---')
45              time.sleep(15)
46
47      # combine all of the unstructured data files into one large file
48      def combine():
49          # Concatenate all files
50          with open('combo.txt', "w") as outfile:
51              for file in glob.glob("*.html"):
52                  with open(file) as infile:
53                      outfile.write(infile.read())
54
55
```

6. Provide a screenshot of the completed `titles()` *method* with the correct code to open and read the HTML file generated by the `combine()` *function*.

```python
assignment.py ×    {..} titles.json

assignment.py > ...
54
55   |
56   #This function imports the BeautifulSoup4 library
57   def titles():
58       from bs4 import BeautifulSoup
59       # stores the resulting JSON file
60       def store_json(data,file):
61           with open(file, 'w', encoding='utf-8') as f:
62               json.dump(data, f, ensure_ascii=False, indent=4)
63               print('wrote file: ' + file)
64
65       #Open and read the large html file generated by combine()
66       f = open('combo.txt', 'r')
67       html = f.read()
68       f.close()
69
70       #the following replaces new line and carriage return char
71       html = html.replace('\n', ' ').replace('\r', '')
72       #the following create an html parser
73       soup = BeautifulSoup(html, "html.parser")
74       results = soup.find_all('h3')
75       titles = []
76
77       # tag inner text
78       for item in results:
79           titles.append(item.text)
80       store_json(titles, 'titles.json')
81
```

7. Provide a screenshot of the fully implemented `clean()` *method* with the correct code to remove all punctuation, numbers, and one-character words

from the `titles.json` file.

```python
81
82
83
84     # remove all punctuation, numbers, and one-character words from the titles.json file.
85     def clean():
86         def store_json(data,file):
87             with open(file, 'w', encoding='utf-8') as f:
88                 json.dump(data, f, ensure_ascii=False, indent=4)
89                 print('wrote file: ' + file)
90
91         with open('titles.json', 'r') as file:
92             titles = json.load(file)
93
94             # remove punctuation/numbers
95             for index, title in enumerate(titles):
96                 punctuation = '''!()-[]{};:'"\,<>./?@#$%^&*_~1234567890'''
97                 translationTable= str.maketrans("","",punctuation)
98                 clean = title.translate(translationTable)
99                 titles[index] = clean
100
101             # remove one character words
102             for index, title in enumerate(titles):
103                 clean = ' '.join( [word for word in title.split() if len(word)>1] )
104                 titles[index] = clean
105
106             store_json(titles, 'titles_clean.json')
107
```

8. Provide a screenshot of the completed `count_words()` *method* with the correct code to call the `store_json(data,file)` helper *function*.

```python
103                    clean = ' '.join( [word for word in title.split() if len(word)>1] )
104                    titles[index] = clean
105
106            store_json(titles, 'titles_clean.json')
107
108
109
110    def count_words():
111        from collections import Counter
112        def store_json(data, file):
113            with open(file, 'w', encoding='utf-8') as f:
114                json.dump(data, f, ensure_ascii=False, indent=4)
115                print('wrote file: ' + file)
116
117
118        with open('titles_clean.json', 'r') as file:
119            titles = json.load(file)
120            words = []
121            # extract words and flatten
122            for title in titles:
123                words.extend(title.split())
124
125            # count word frequency
126            counts = Counter(words)
127            store_json(counts, 'words.json')
128
129
130
```

9. Provide a screenshot of the DAG declaration with all six *tasks*, from `t0` to `t5`, correctly defined.

```python
124
125
126    with DAG(
127        "assignment",
128        start_date=days_ago(1),
129        schedule_interval="@daily",catchup=False,
130    ) as dag:
131
132    # INSTALL BS4 BY HAND THEN CALL FUNCTION
133
134        # ts are tasks
135        t0 = BashOperator(
136            task_id='task_zero',
137            bash_command='pip install beautifulsoup4',
138            retries=2
139        )
140        t1 = PythonOperator(
141            task_id='task_one',
142            depends_on_past=False,
143            python_callable=catalog
144        )
145        t2 = PythonOperator(
146            task_id='task_two',
147            depends_on_past=False,
148            python_callable=combine
149        )
150        t3 = PythonOperator(
151            task_id='task_three',
152            depends_on_past=False,
153            python_callable=titles
154        )
155        t4 =PythonOperator (
156            task_id='task_four',
157            depends_on_past=False,
158            python_callable=clean
159        )
160        t5 = PythonOperator(
161            task_id='task_five',
162            depends_on_past=False,
163            python_callable=count_words
164        )
165
166        t0>>t1>>t2>>t3>>t4>>t5
167
168    |
```
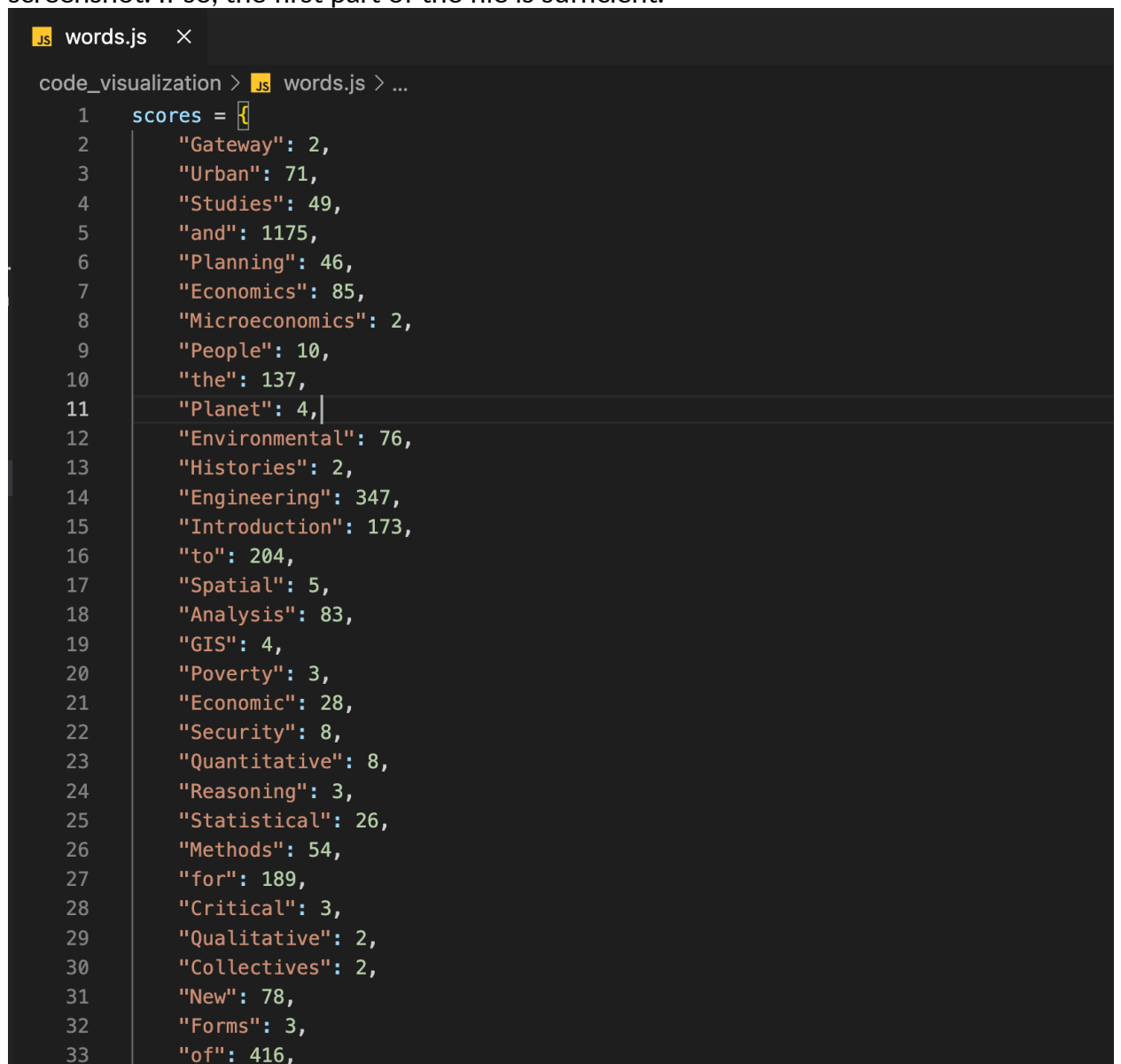
**Part 2: Code Execution**

1.  Provide a screenshot of your Docker application that shows that your Airflow
    Docker *container* has initiated.



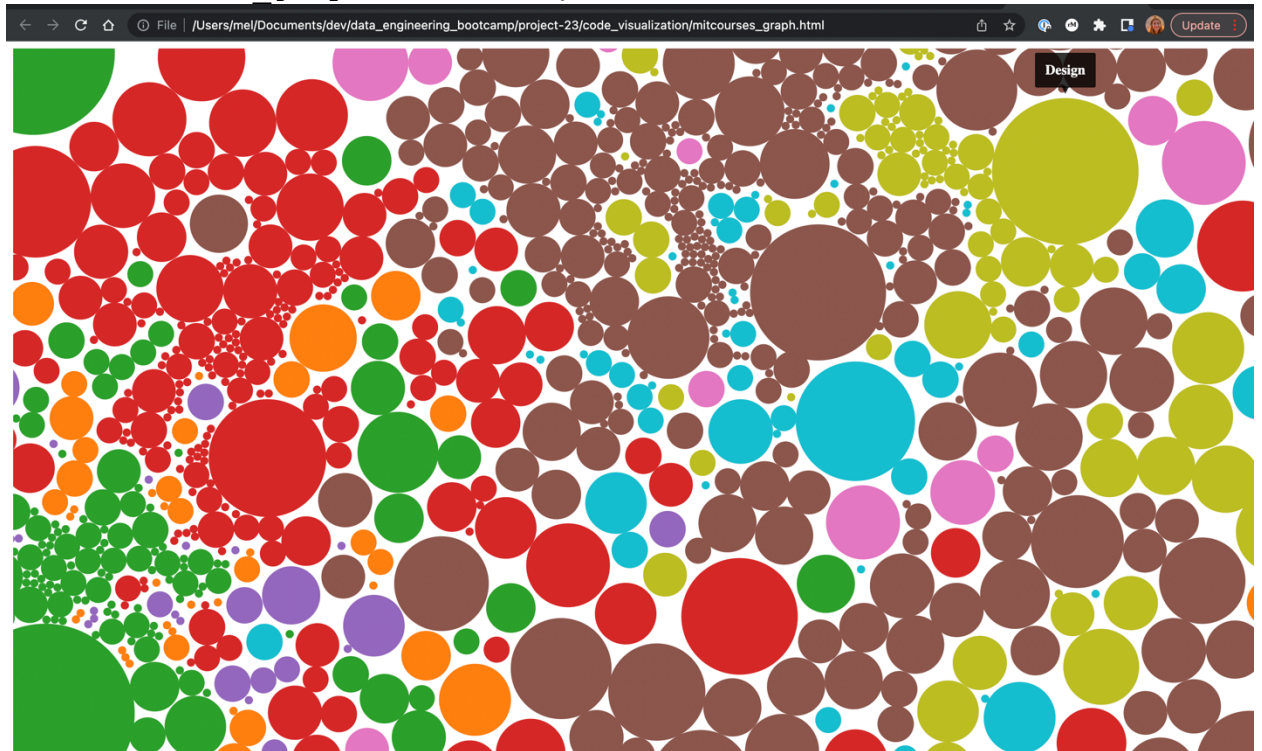2.  Provide a screenshot of the *task* boxes to show that the DAG ran successfully.

3. Provide a screenshot of the `words.js` file with the data from the `words.json` file. You may not be able to fit all the data in one screenshot. If so, the first part of the file is sufficient.

```js
scores = {
    "Gateway": 2,
    "Urban": 71,
    "Studies": 49,
    "and": 1175,
    "Planning": 46,
    "Economics": 85,
    "Microeconomics": 2,
    "People": 10,
    "the": 137,
    "Planet": 4,
    "Environmental": 76,
    "Histories": 2,
    "Engineering": 347,
    "Introduction": 173,
    "to": 204,
    "Spatial": 5,
    "Analysis": 83,
    "GIS": 4,
    "Poverty": 3,
    "Economic": 28,
    "Security": 8,
    "Quantitative": 8,
    "Reasoning": 3,
    "Statistical": 26,
    "Methods": 54,
    "for": 189,
    "Critical": 3,
    "Qualitative": 2,
    "Collectives": 2,
    "New": 78,
    "Forms": 3,
    "of": 416,
```

4. Provide a screenshot of the visualization produced with the `mitcourses_graph.html` file in your web browser.



5. Provide a screenshot of your enhanced visualization created with the D3 *library* and the modified example code.