

DA5402: Machine Learning Operations Laboratory

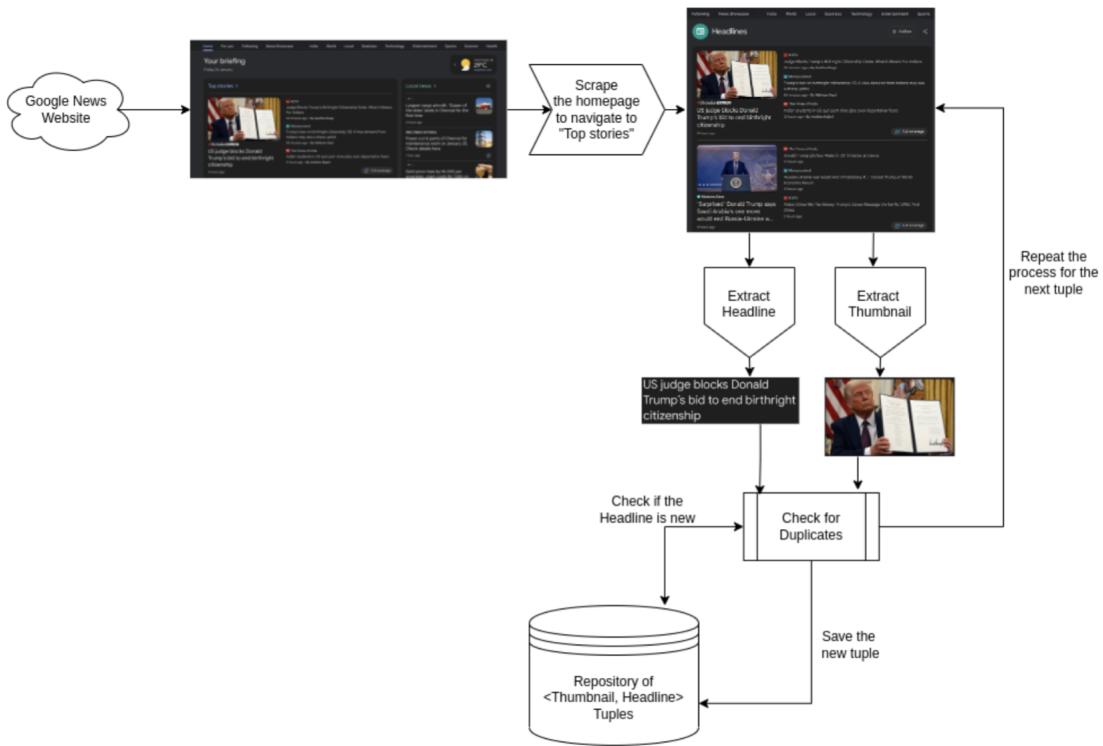
Assignment 2

Nandhakishore C S (DA24M011)

February 18, 2025

1 Problem Statement

We got introduced to Apache Airflow today, which happens to be one of most popular open-source orchestration tool amongst the industry practitioners. Let's rebuild our solution for Assignment 1 using Airflow this time. This exercise would teach you the power of using the right tool to the right problem. The workflow that we implemented in A1 is shown again here as a ready reckoner.



1.1 Module 1 [5 points]

Create a suitable Airflow Operator that would either import your Python function or directly execute your Python script to scrape the home page of Google News. The URLs to scrape may change over time, so you should keep them as configurable parameters or operator properties.

The code for module 1 is re-purposed from assignment1. This returns the HTML data from the webpage which is scrapped. The functions from airflow to pull and push data into Xcom are added.

1.2 Module 2 [5 points]

Create a suitable Airflow Operator that would either import your Python function or directly executing your Python script to scrape the “Top stories” from the home page of Google News. The URLs to scrape may change over time, so you have to keep them as configurable parameters or operator properties.

The code for module 2 is re-purposed from assignment1. This returns the URL of the Top stories section from the Google News webpage.

```

1 # Module 1 - get HTML content
2 def module1(**kwargs):
3     response = requests.get(url)
4     if response.status_code == 200:
5         kwargs['ti'].xcom_push(key='home_page', value=response.text)
6     else:
7         raise Exception(f"\nFailed to scrape:\t {url}\n")

```

Figure 1: Snap shot of the module1 from DAG1.

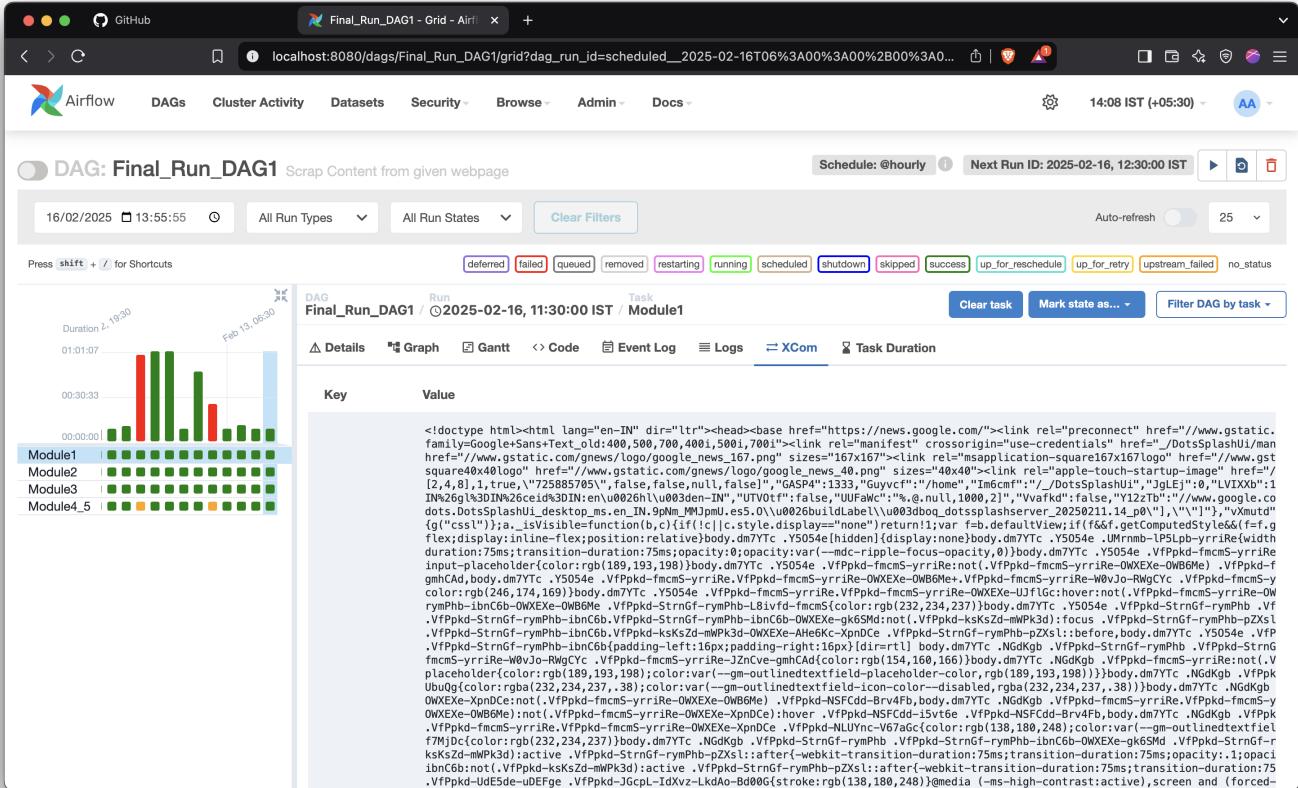


Figure 2: HTML Tags scrapped from news.google.com

```

1 # Module 2 - get URL for top stories sub section
2 def module2(**kwargs):
3     ti = kwargs['ti']
4     html = ti.xcom_pull(task_ids='sub_section_url', key='home_page')
5     soup = BeautifulSoup(html, 'html.parser')
6     # top_story = soup.find('a', {'aria-label': 'Top stories'})
7     top_story = soup.find_all('a')
8     pattern = sub_section_str
9
10    for link in top_story:
11        link_text = link.get_text().strip()
12        if(pattern.lower() in link_text.lower()):
13            href = link.get('href')
14            if(href is not None):
15                # href = f'{url}{href[1 : -1]}'
16                top_stories_url = f'https://news.google.com{href[1:-1]}'
17                ti.xcom_push(key='top_stories_url', value=top_stories_url)

```

Figure 3: Snap shot of the module2 from DAG1.

1.3 Module 3 [10 points]

Create a suitable Airflow Operator that would either import your Python function or directly executing your Python script to extract the thumbnail and the headline of every story from that page. Remember that the page is set up for lazy loading. Your operator should factor lazy loading. The layout of the “Top stories” page may change over time, so ensure that you create your Operator in a way for easy updates later.

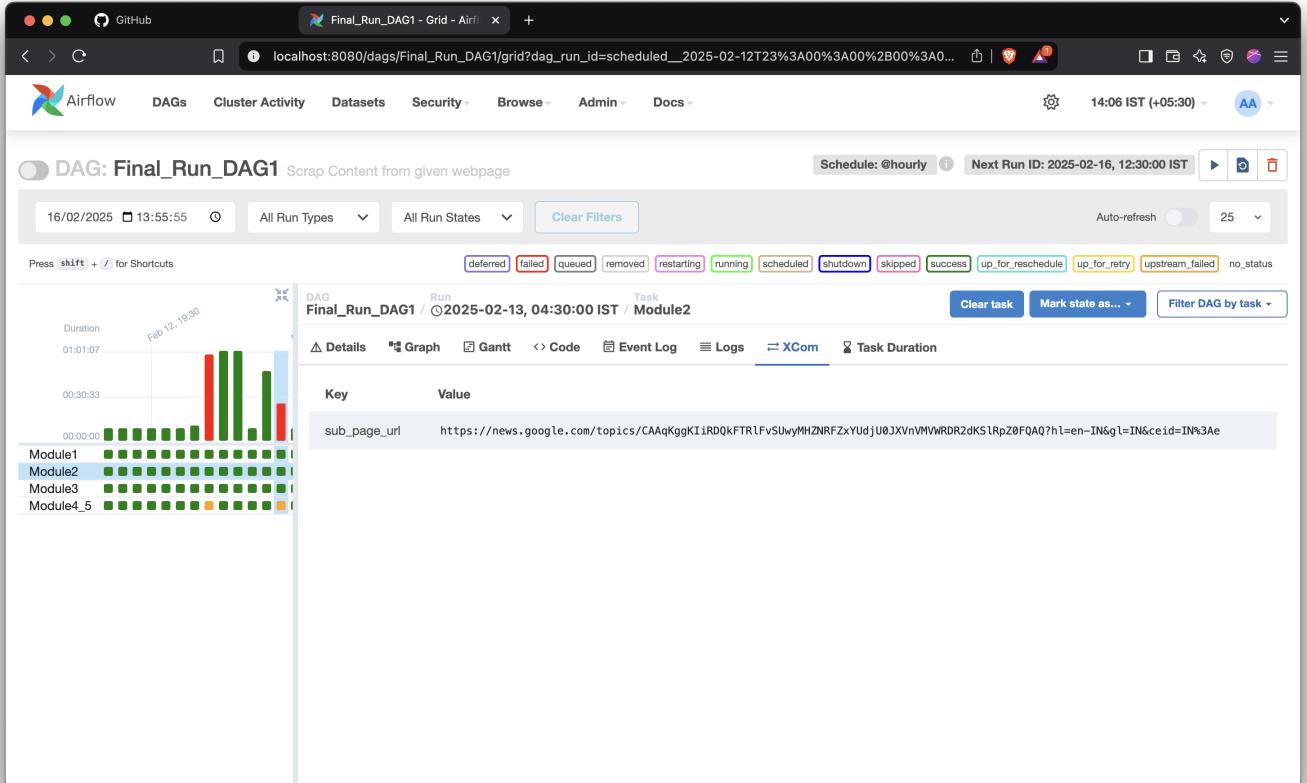


Figure 4: Top Stories URL from module2

The code for module 3 is re-purposed from assignment1. The headline, the thumbnail, date and article URL are dumped in the Airflow's Xcom as a dictionary, such that it can be accessed by a database operator if needed.

The module3 has two helper functions which are again repurposed from assignment1. The functions get the time (for getting article's publication date / save date) and to get the content (English text and URL) from the scrapped HTML tags. To handle lazy loading, a while loop with time.sleep function is utilised to wait for the page to load. Snap shot of the module3 from DAG1.

1.4 Module 4 [10 points]

Create a Postgres operator to setup your database tables with the necessary unique/primary keys. Your database should have one table for storing the image data (base64 encoded) and the other table for storing the headlines and other meta information such as URLs, scrape timestamp, article date, etc. You can run a join query with the image and headlines tables to fetch the {title, image} pair. After the “insert” operation, you should create a status file at “dags/run/status”, which should contain just one number representing the number of successful inserts. If all the records are duplicates, you will have a ‘0’ in the status file.

1.5 Module 5 [10 points]

Setup a Airflow DAG to orchestrate the workflow pipeline with necessary properties and load them up in the Airflow console. You should setup crontab, to make it runnable every hour. As long as the Airflow container is running in the background, you will keep scrapping the articles for <image, headline> tuple

The code for module4 and module5 are combined together for the ease of coding it without errors. The function is repurposed from assignment1. The module creates a status file in the local system named 'status' and connects the DAG to a postgres database using the **psycopg2** library. This is a Airflow's Python operator which acts as a database connector and inserts data on the database. This can be viewed using the DBeaver application in the local system.

The module uses Jaccard Similarity as a measure to find the similarity between the files.

```
1 def _jaccard_similarity(document1: str, document2: str) -> float:
2     s1 = set(document1.lower().split())
```

```

1 # Helper function for module3
2 def _extract_story_data(article: bs4.element.Tag) -> dict:
3     try:
4         # Extract headline and URL first
5         headline = None
6         for a_tag in article.find_all('a'):
7             if(a_tag.get_text().strip() is not None):
8                 headline = a_tag.get_text().strip()
9                 article_url = a_tag.get('href')
10                if((article_url) and (not article_url.startswith('http'))):
11                    # article_url = f"https://news.google.com{article_url[1:-1]}"
12                    article_url = "https://news.google.com" + article_url[1:-1]
13                else:
14                    break
15
16            if(not headline):
17                return None
18
19        # Extract thumbnail
20        thumbnail_url = None
21
22        # Try figure tag first
23        figure = article.find('figure')
24        if(figure is not None):
25            img = figure.find('img')
26            if(img is not None):
27                thumbnail_url = img.get('src') or img.get('data-src')
28                # thumbnail_url = f"https://news.google.com{thumbnail_url}"
29                thumbnail_url = "https://news.google.com" + thumbnail_url
30
31        # If no image found, skip this story
32        if(not thumbnail_url): return
33
34        # Extract and convert publication date
35        time_elem = article.find('time')
36        relative_time = time_elem.get_text().strip() if time_elem else None
37        pub_date = _get_time(relative_time) if relative_time else datetime.now().strftime("%Y-%m-%d %H
38 :%M:%S")
39
40        return {
41            'Headline': headline,
42            'Thumbnail_url': thumbnail_url,
43            'Article_url': article_url,
44            'Date': pub_date
45        }
46
47    except Exception as error:
48        print(f"Error:\t {error}")
49        return
50
51 # Helper function for module 3
52 def _get_time(relative_time: str) -> str:
53     try:
54         now = datetime.now()
55
56         if 'minute' in relative_time or 'minutes' in relative_time:
57             minutes = int(relative_time.split()[0])
58             return (now - timedelta(minutes=minutes)).strftime("%Y-%m-%d %H:%M:%S")
59
60         elif 'hour' in relative_time or 'hours' in relative_time:
61             hours = int(relative_time.split()[0])
62             return (now - timedelta(hours=hours)).strftime("%Y-%m-%d %H:%M:%S")
63
64         elif 'day' in relative_time or 'days' in relative_time:
65             days = int(relative_time.split()[0])
66             return (now - timedelta(days=days)).strftime("%Y-%m-%d %H:%M:%S")
67
68         elif 'Yesterday' in relative_time:
69             return (now - timedelta(days=1)).strftime("%Y-%m-%d %H:%M:%S")
70
71     else:
72         return now.strftime("%Y-%m-%d %H:%M:%S")
73
74 except Exception:
75     return now.strftime("%Y-%m-%d %H:%M:%S")

```

Figure 5: Helpers functions for module3 from DAG1.

```

1 # Module 3 - get thumbnail, headlines and META data
2 def module3(**kwargs):
3     ti = kwargs['ti']
4     top_stories_url = ti.xcom_pull(task_ids='extract_top_stories', key='top_stories_url')
5
6     stories = []
7     page = 1
8     while(page <= 5):
9         try:
10             response = requests.get(top_stories_url)
11             response.raise_for_status()
12             soup = BeautifulSoup(response.text, 'html.parser')
13             articles = soup.find_all('article')
14
15             for article in articles:
16                 story = _extract_story_data(article)
17                 if story and story not in stories:
18                     stories.append(story)
19
20             if ((not articles) or (len(articles) == 0)):
21                 break
22
23             page += 1
24             time.sleep(5)
25
26         except requests.RequestException as error:
27             print(f'Error fetching {page}: {error}')
28             break
29
30     ti.xcom_push(key='news_data', value=stories)

```

Figure 6: Snap shot of the module3 from DAG1.

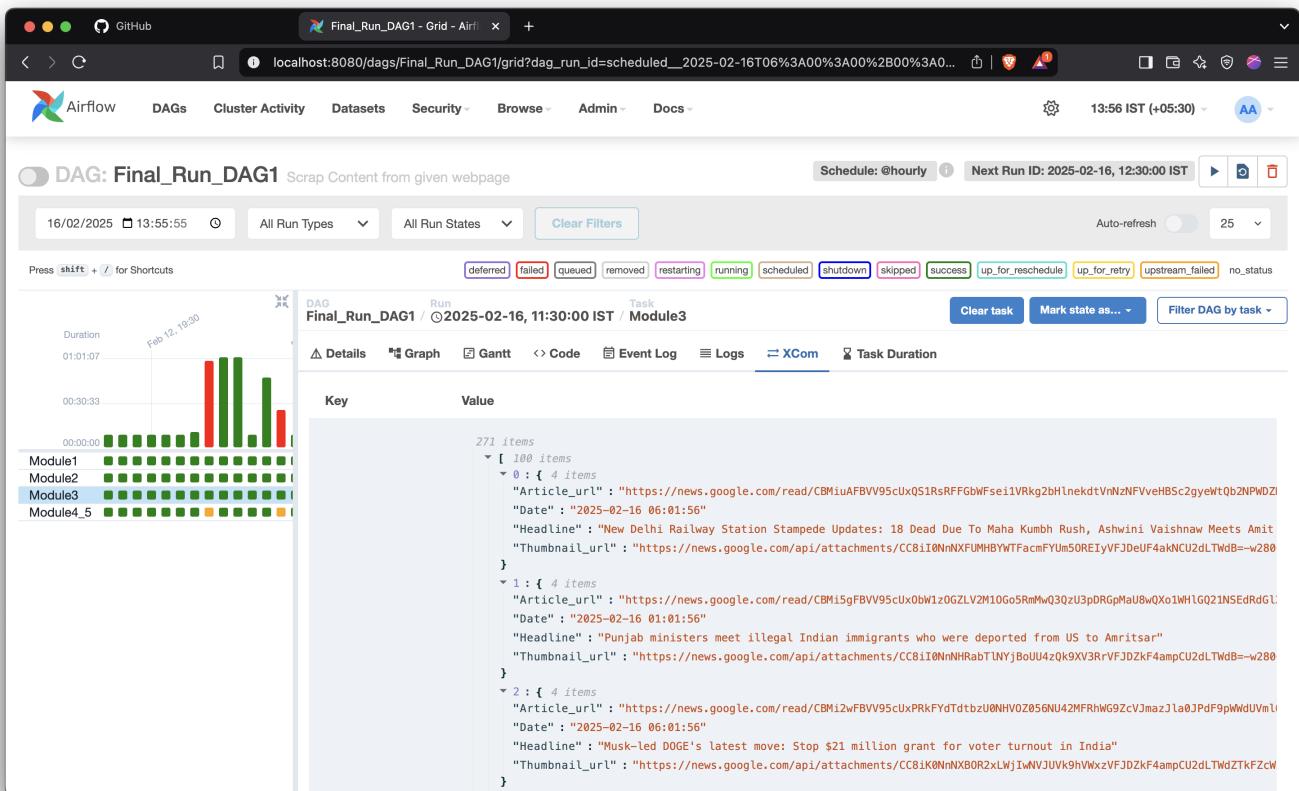


Figure 7: Data scrapped from top stories page

```

3 s2 = set(document2.lower().split())
4
5 intersection = len(s1.intersection(s2))
6 union = len(s1.union(s2))
7
8 return intersection / union if union > 0 else 0

```

The screenshot shows the DBeaver interface with the 'news_headlines' table selected. The table structure is as follows:

headline_id	image_id	headline	thumbnail_url	article_url	publication_date
1	1	Second FIR filed against Ranveer Alla	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 15:56
2	2	Maharashtra police files FIR against F	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 15:01
3	3	Veer Savarkar is the backdrop to PM	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 12:01
4	4	AAP chief Kejriwal meets Punjab CM	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 03:01
5	5	2 Soldiers Killed In Action As Terroris	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 13:01
6	6	JEE Main 2025 Session 1 Result Decl	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 03:01
7	7	Delhi new CM announcement live: BJ	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 05:01
8	8	Stuck in Bengaluru traffic, German pi	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 08:01
9	9	'No Vehicle Zone' enforced in Maha	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 07:01
10	10	Out on parole, Engineer Rashid raise	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 10:01
11	11	10 countries, including India, most af	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 04:01
12	12	Hamas lashes out at Trump's warning	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 08:01
13	13	Delhi, UP, Haryana declare February	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 14:01
14	14	How a con artist, a biz man and a 'ha	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 07:01
15	15	Rahul, Kharge skip Invest K'taka Sum	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 08:01
16	16	BMRC fare hike: BMRC had sought a	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-10 21:01
17	17	"Don't Delete Data": Supreme Court	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 12:01
18	18	Donald Trump halts US enforcement	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 08:01
19	19	Income Tax Bill 2025 Live Updates: S	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 09:01
20	20	Pariksha Pe Charcha: Radhika Gupta	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 17:01
21	21	How did Tamil Nadu Governor send n	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 06:01
22	22	4 months after J&K govt formed, rule	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 17:01
23	23	Meet the mom who flies to work ever	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 10:01
24	24	Prashant Kishor gives reasons behin	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 17:01
25	25	Unnatural sex act by husband withou	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 13:01
26	26	1,000-pt Sensex crash: Has Trump s	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 14:01
27	27	Six U.S. Congressmen write to New A	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 09:01
28	28	Trump top economic adviser Kevin H	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 17:01
29	29	Vodafone Idea Q3 results: Net loss a	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 14:01
30	30	Palestinians won't have right to retur	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-10 18:01
31	31	Rishabh Pant & Arshdeep Singh IN,	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 16:50
32	32	Jasprit Bumrah's Champions Trophy	https://news.google.com/api/attach	https://news.google.com/read/CBN	2025-02-11 02:01

Figure 8: Table in database with entries of headlines, article URL , thumbnail URL and meta data. The column headline ID is used as a foreign key to infer the other table to save the thumbnail image

```

9
10 def is_similar(new_headline, existing_headlines):
11     for headline in existing_headlines:
12         if SequenceMatcher(None, new_headline, headline).ratio() > 0.85:
13             return True
14     return False

```

Helper function to check similarity between headlines for module 4 and 5

```

1 def module4(**kwargs):
2     try:
3         ti = kwargs['ti']
4         news_data = ti.xcom_pull(task_ids='extract_news_data', key='news_data')
5         pg_hook = PostgresHook(postgres_conn_id = postgres_connection_ID)
6         connection = pg_hook.get_conn()
7         cursor = connection.cursor()
8
9         # counter variable for successfull inserts
10        counter = 0
11
12        for story in news_data:
13            cursor.execute(
14                "SELECT Headline FROM news_headlines"
15            )
16            current_headlines = cursor.fetchall()
17
18            duplicate = False
19            for (current_headline, ) in current_headlines:
20                # similarity = is_similar(story['Headline'], current_headline)
21                # if(similarity is not True):
22                #     print("\nDuplicate values of headline found!\n")
23                #     duplicate = True
24                #     break
25                similarity = _jaccard_similarity(story['Headline'], current_headline)
26            if(similarity >= .75):
27                print("\nDuplicate values of headline found!\n")
28                duplicate = True
29                break
30

```

```

31     if(duplicate == True):
32         continue
33
34     thumbnail_url = story['Thumbnail_url']
35     thumbnail_data = requests.get(thumbnail_url)
36     thumbnail_data.raise_for_status()
37     thumbnail_image = base64.b64encode(thumbnail_data.content).decode('utf-8')
38
39     cursor.execute(
40         "INSERT INTO news_images (image_data) VALUES (%s) RETURNING image_id",
41         (thumbnail_image,))
42
43     thumbnail_ID = cursor.fetchone()[0]
44
45     cursor.execute(
46         """
47             INSERT INTO news_headlines
48                 (image_id, headline, thumbnail_url, article_url, publication_date,
49                  scrape_timestamp)
49                 VALUES (%s, %s, %s, %s, %s, %s)
50                 ON CONFLICT (headline, article_url) DO NOTHING
51                 RETURNING headline_id
52             """,
53             (
54                 thumbnail_ID,
55                 story['Headline'],
56                 story['Thumbnail_url'],
57                 story['Article_url'],
58                 story['Date'],
59                 datetime.now()
60             )
61         )
62
63
64     if (cursor.fetchone() is not None):
65         counter += 1
66     else:
67         cursor.execute(
68             "DELETE FROM news_images WHERE image_id = %s", (thumbnail_ID, ))
69
70     # connection commit
71     connection.commit()
72
73     # Status check
74     dir_path = '/opt/airflow/dags'
75     os.makedirs(dir_path, exist_ok = True)
76
77     with open(os.path.join(dir_path, 'status'), 'w') as file:
78         file.write(str(counter))
79
80 except Exception as error:
81     print(error)
82     raise

```

1.6 Module 6[10 points]

Setup another Airflow DAG that would send an email to yourself, whenever a new `(image, headline)` tuple is added to the database. This DAG has to remember the previous state of the database tables to detect new rows. When news rows are detected, it should remember its knowledge of the previous state and also send an email to your inbox. To send emails, you may have to learn about SMTP + TLS/SSL setup with credentials. The set up is the same process when you configure your email client (Thunderbird, Outlook, etc) to send emails from your gmail or institute accounts. The only difference is you are going to do it programmatically.

The module to check for new entries and send email requires the two factor authentication code from the Gmail account. Insert your code and email details at the start of the DAG.

```

1 # Configuration to send notification email using DAGS
2 SMTP_SERVER = "smtp.gmail.com"
3 SMTP_PORT = 587
4 SENDER_EMAIL = "" # sender email
5 PASSWORD = "" # send 2FA pass key
6 RECIPIENT_EMAIL = "" # recipient email
7
8 # ----- Helper Function to do Module 6 -----
9 def _send_email(message):
10     # Send email using SMTP with App Password
11     try:
12         with smtplib.SMTP(SMTP_SERVER, SMTP_PORT) as server:

```

DBeaver 24.3.4 - news_images

Database Navigator X news_headlines news_images X Properties Data ER Diagram

Enter a part of object name here

airflow localhost:54320 Databases airflow Schemas public Tables Foreign Tables Views Materialized Views Indexes Functions Sequences Data types Aggregate functions Event Triggers Extensions Storage System Info Roles postgres Administer System Info DBeaver Sample Database (SQLite) nandhakishorecs localhost:54320

news_images Grid 123 image_id image_data created_at

Record	image_id	image_data	created_at
1	1	[base64 encoded image data]	2025-02-11 16:56:41.588
2	2	[base64 encoded image data]	2025-02-11 17:01:55.131
3	3	[base64 encoded image data]	2025-02-11 17:01:55.131
4	4	[base64 encoded image data]	2025-02-11 17:01:55.131
5	5	[base64 encoded image data]	2025-02-11 17:01:55.131
6	6	[base64 encoded image data]	2025-02-11 17:01:55.131
7	7	[base64 encoded image data]	2025-02-11 17:01:55.131
8	8	[base64 encoded image data]	2025-02-11 17:01:55.131
9	9	[base64 encoded image data]	2025-02-11 17:01:55.131
10	10	[base64 encoded image data]	2025-02-11 17:01:55.131
11	11	[base64 encoded image data]	2025-02-11 17:01:55.131
12	12	[base64 encoded image data]	2025-02-11 17:01:55.131
13	13	[base64 encoded image data]	2025-02-11 17:01:55.131
14	14	[base64 encoded image data]	2025-02-11 17:01:55.131
15	15	[base64 encoded image data]	2025-02-11 17:01:55.131
16	16	[base64 encoded image data]	2025-02-11 17:01:55.131
17	17	[base64 encoded image data]	2025-02-11 17:01:55.131
18	18	[base64 encoded image data]	2025-02-11 17:01:55.131
19	19	[base64 encoded image data]	2025-02-11 17:01:55.131
20	20	[base64 encoded image data]	2025-02-11 17:01:55.131
21	21	[base64 encoded image data]	2025-02-11 17:01:55.131
22	22	[base64 encoded image data]	2025-02-11 17:01:55.131
23	23	[base64 encoded image data]	2025-02-11 17:01:55.131
24	24	[base64 encoded image data]	2025-02-11 17:01:55.131
25	25	[base64 encoded image data]	2025-02-11 17:01:55.131
26	26	[base64 encoded image data]	2025-02-11 17:01:55.131
27	27	[base64 encoded image data]	2025-02-11 17:01:55.131
28	28	[base64 encoded image data]	2025-02-11 17:01:55.131
29	29	[base64 encoded image data]	2025-02-11 17:01:55.131
30	30	[base64 encoded image data]	2025-02-11 17:01:55.131
31	31	[base64 encoded image data]	2025-02-11 17:01:55.131
32	32	[base64 encoded image data]	2025-02-11 17:01:55.131

Apply resultset changes: Ch...connection auto-commit state IST en_IN Refresh Save Cancel Export data 200 200+ Inserted: 0 / Deleted: 0 / Updated: 1

Figure 9: Table with thumbnail data in base64 encoded form

```

13     server.starttls()
14 # Using App Password here - get token from GMail 2FA
15 server.login(SENDER_EMAIL, PASSWORD)
16 server.send_message(message)
17     print("Email sent successfully")
18 except Exception as error:
19     print(f"Failed to send email: {error}")
20     raise
21
22 # ----- Helper Function to do Module 6 -----
23 def _load_previous_state():
24     # Load previous database content - the penultimate state from a JSON file
25     try:
26         with open('/opt/airflow/dags/previous_state.json', 'r') as file:
27             return json.load(file)
28     except FileNotFoundError:
29         return {'last_headline_id': 0}
30
31 # ----- Helper Function to do Module 6 -----
32 def save_current_state(state):
33     # Save current database state to a JSON file
34     with open('/opt/airflow/dags/previous_state.json', 'w') as file:
35         json.dump(state, file)
36
37 # ----- Helper Function to do Module 6 -----
38 def _compose_email(new_entries):
39     # To email the new changes compose an email with content from new entries
40     message = MIMEMultipart()
41     message["Subject"] = f"New Articles Alert: {len(new_entries)} new entries"
42     message["From"] = SENDER_EMAIL
43     message["To"] = RECIPIENT_EMAIL
44
45     if(len(new_entries) > 0):
46         content = "New articles have been added!:\n"
47         for entry in new_entries:
48             content += f"    {entry['headline']}\\n"
49             if entry.get('publication_date'):
50                 content += f"Published on\\t: {entry['publication_date']}\\n"
51             content += f"    URL: {entry['article_url']}\\n\\n"
52
53     message.attach(MIMEText(content, "plain"))

```

```

54     else:
55         content = 'No new top stories!\n'
56         message.attach(MIMEText(content, "plain"))
57     return message
58
59 # ----- Main Function for Module 6 -----
60 def module6(**context):
61     try:
62         # Get previous state
63         prev_state = _load_previous_state()
64         last_headline_id = prev_state['last_headline_id']
65
66         # Connect to database
67         pg_hook = PostgresHook(postgres_conn_id = postgres_connection_ID)
68         connection = pg_hook.get_conn()
69         cursor = connection.cursor()
70
71         # Get new entries - SQL Query
72         cursor.execute(
73             """
74                 SELECT h.headline_id, h.Headline, h.Article_url, h.publication_date
75                 FROM news_headlines h
76                 WHERE h.headline_id > %s
77                 ORDER BY h.headline_id
78             """,
79             (last_headline_id, )
80         )
81
82         new_entries = []
83         max_headline_id = last_headline_id
84
85         for row in cursor.fetchall():
86             headline_id, headline, article_url, pub_date = row
87             max_headline_id = max(max_headline_id, headline_id)
88             new_entries.append({
89                 'headline': headline,
90                 'article_url': article_url,
91                 'publication_date': pub_date.isoformat() if pub_date else None
92             })
93
94         if (new_entries is not None):
95             # Prepare and send email
96             message = _compose_email(new_entries)
97             _send_email(message)
98
99             # Update state
100            save_current_state({'last_headline_id': max_headline_id})
101
102        # Clean up status file
103        os.remove('/opt/airflow/dags/status')
104
105    except Exception as error:
106        print(f"Error in sending email: {error}")
107        raise
108
109    finally:
110        if 'cursor' in locals():
111            cursor.close()
112        if 'connection' in locals():
113            connection.close()

```

Function for module6 from DAG2.

```

1 # Configuration to send notification email using DAGs
2 SMTP_SERVER = "smtp.gmail.com"
3 SMTP_PORT = 587
4 SENDER_EMAIL = "" # sender email
5 PASSWORD = "" # send 2FA pass key
6 RECIPIENT_EMAIL = "" # recipient email
7
8 # ----- Helper Function to do Module 6 -----
9 def _send_email(message):
10     # Send email using SMTP with App Password
11     try:
12         with smtplib.SMTP(SMTP_SERVER, SMTP_PORT) as server:
13             server.starttls()
14             # Using App Password here - get token from GMail 2FA
15             server.login(SENDER_EMAIL, PASSWORD)
16             server.send_message(message)
17             print("Email sent successfully")

```

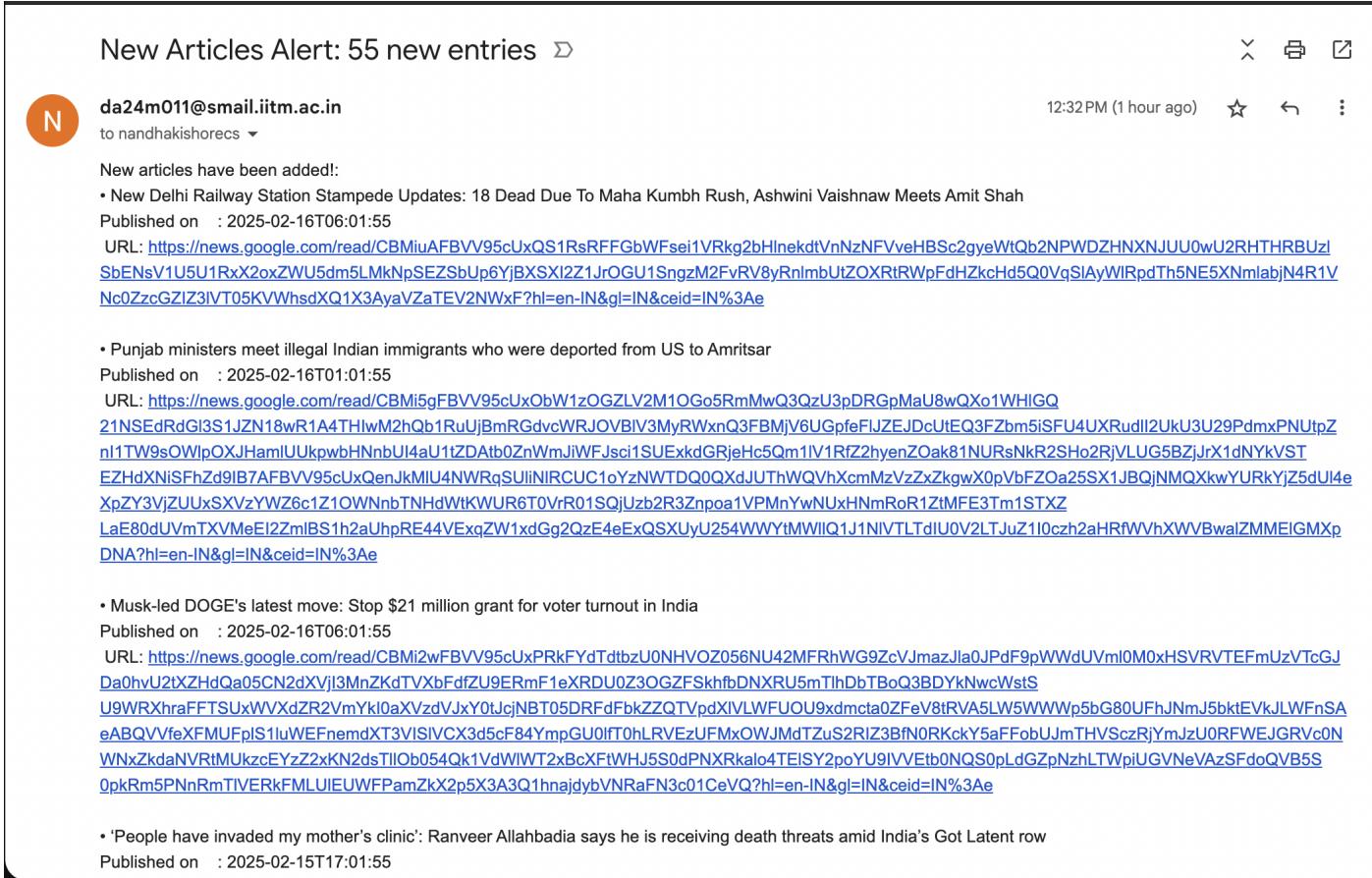


Figure 10: Emails being send from DAG2 to about new additions of news articles

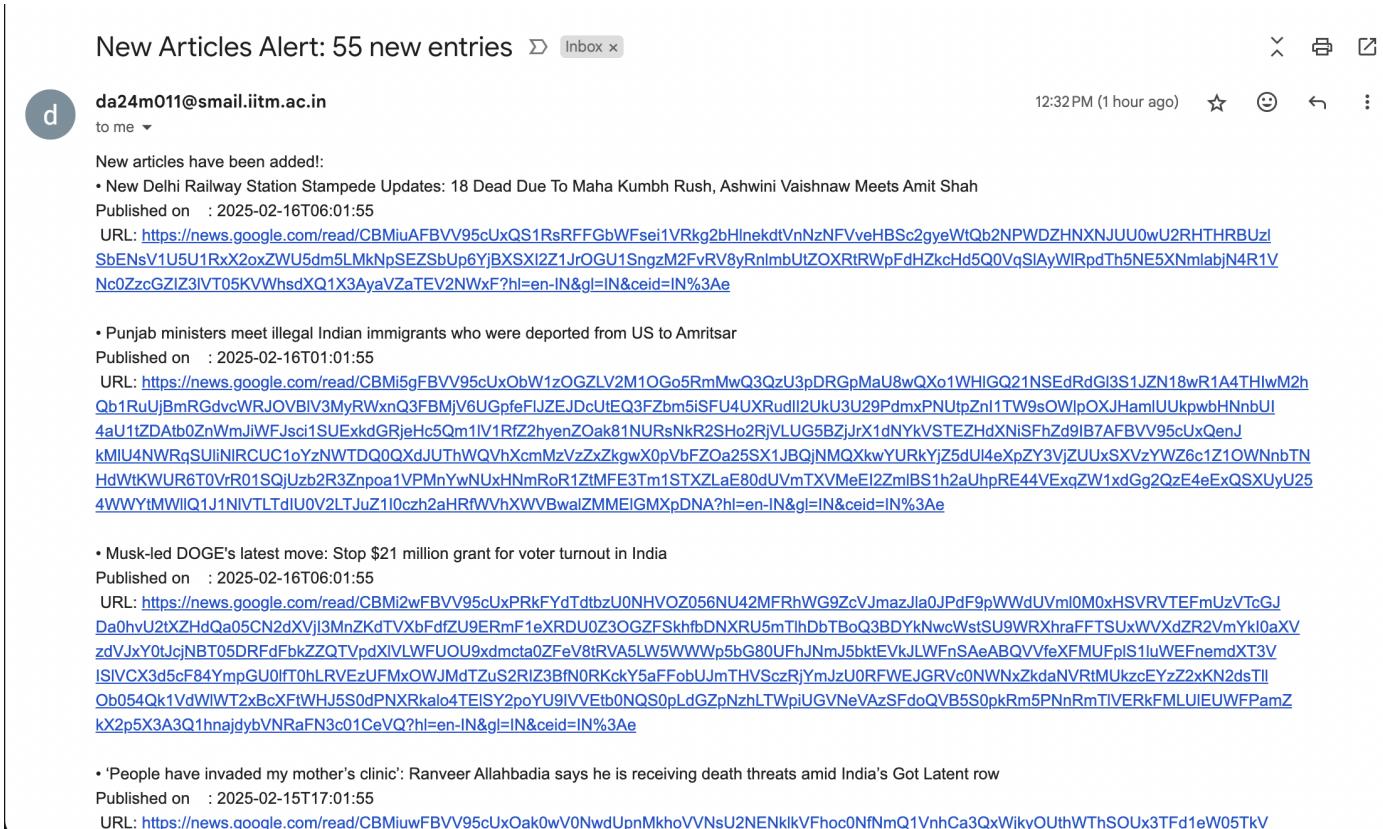


Figure 11: Emails received from DAG2 about new headlines

```

18     except Exception as error:
19         print(f"Failed to send email: {error}")
20         raise
21
22 # ----- Helper Function to do Module 6 -----
23 def _load_previous_state():
24     # Load previous database content - the pultimate state from a JSON file
25     try:
26         with open('/opt/airflow/dags/previous_state.json', 'r') as file:
27             return json.load(file)
28     except FileNotFoundError:
29         return {'last_headline_id': 0}
30
31 # ----- Helper Function to do Module 6 -----
32 def save_current_state(state):
33     # Save current database state to a JSON file
34     with open('/opt/airflow/dags/previous_state.json', 'w') as file:
35         json.dump(state, file)
36
37 # ----- Helper Function to do Module 6 -----
38 def _compose_email(new_entries):
39     # To email the new changes compose an email with content from new entries
40     message = MIMEMultipart()
41     message["Subject"] = f"New Articles Alert: {len(new_entries)} new entries"
42     message["From"] = SENDER_EMAIL
43     message["To"] = RECIPIENT_EMAIL
44
45     if(len(new_entries) > 0):
46         content = "New articles have been added!:\n"
47         for entry in new_entries:
48             content += f"    {entry['headline']}\\n"
49             if entry.get('publication_date'):
50                 content += f"Published on\\t: {entry['publication_date']}\\n"
51             content += f"    URL: {entry['article_url']}\\n\\n"
52
53         message.attach(MIMEText(content, "plain"))
54     else:
55         content = 'No new top stories!\\n'
56         message.attach(MIMEText(content, "plain"))
57     return message
58
59 # ----- Main Function for Module 6 -----
60 def module6(**context):
61     try:
62         # Get previous state
63         prev_state = _load_previous_state()
64         last_headline_id = prev_state['last_headline_id']
65
66         # Connect to database
67         Pg_hook = PostgresHook(postgres_conn_id = postgres_connection_ID)
68         connection = pg_hook.get_conn()
69         cursor = connection.cursor()
70
71         # Get new entries - SQL Query
72         cursor.execute(
73             """
74                 SELECT h.headline_id, h.Headline, h.Article_url, h.publication_date
75                 FROM news_headlines h
76                 WHERE h.headline_id > %s
77                 ORDER BY h.headline_id
78             """,
79             (last_headline_id, )
80         )
81
82         new_entries = []
83         max_headline_id = last_headline_id
84
85         for row in cursor.fetchall():
86             headline_id, headline, article_url, pub_date = row
87             max_headline_id = max(max_headline_id, headline_id)
88             new_entries.append({
89                 'headline': headline,
90                 'article_url': article_url,
91                 'publication_date': pub_date.isoformat() if pub_date else None
92             })
93
94         if (new_entries is not None):
95             # Prepare and send email
96             message = _compose_email(new_entries)
97             _send_email(message)

```

```

99     # Update state
100    save_current_state({'last_headline_id': max_headline_id})
101
102    # Clean up status file
103    os.remove('/opt/airflow/dags/status')
104
105 except Exception as error:
106     print(f"Error in sending email: {error}")
107     raise
108
109 finally:
110     if 'cursor' in locals():
111         cursor.close()
112     if 'connection' in locals():
113         connection.close()
114
115 # ----- Define default_args for DAGs -----
116 default_args = {
117     'owner': 'airflow',
118     'depends_on_past': False,
119     'start_date': datetime(2024, 2, 18),
120     'retries': 2,
121     'retry_delay': timedelta(minutes=5),
122 }
123
124 # ----- DAG Configuration -----
125 with DAG('DAG_1_FINAL', default_args=default_args, schedule_interval='@hourly', catchup=False) as dag1:
126     :
127     scrape_task = PythonOperator(task_id='sub_section_url', python_callable=module1, provide_context=True)
128     extract_task = PythonOperator(task_id='extract_top_stories', python_callable=module2, provide_context=True)
129     get_news_task = PythonOperator(task_id='extract_news_data', python_callable=module3, provide_context=True)
130     create_table_task = PostgresOperator(
131         task_id = 'create_tables',
132         postgres_conn_id = postgres_connection_ID,
133         sql= """
134             DO $$$
135             BEGIN
136                 -- Create news_images table if it doesn't exist
137                 IF NOT EXISTS (SELECT FROM pg_tables WHERE schemaname = 'public' AND tablename = 'news_images') THEN
138                     CREATE TABLE news_images (
139                         image_id SERIAL PRIMARY KEY,
140                         image_data BYTEA NOT NULL,
141                         created_at TIMESTAMP DEFAULT CURRENT_TIMESTAMP
142                     );
143                 END IF;
144
145                 -- Create news_headlines table if it doesn't exist
146                 IF NOT EXISTS (SELECT FROM pg_tables WHERE schemaname = 'public' AND tablename = 'news_headlines') THEN
147                     CREATE TABLE news_headlines (
148                         headline_id SERIAL PRIMARY KEY,
149                         image_id INTEGER REFERENCES news_images(image_id),
150                         headline TEXT NOT NULL,
151                         thumbnail_url TEXT,
152                         article_url TEXT,
153                         publication_date TIMESTAMP,
154                         scrape_timestamp TIMESTAMP DEFAULT CURRENT_TIMESTAMP,
155                         UNIQUE(headline, article_url)
156                     );
157                 END IF;
158             END $$;
159         """
160     )
161     store_data_task = PythonOperator(task_id='store_news_data', python_callable=module4, provide_context=True)
162     create_table_task >> scrape_task >> extract_task >> get_news_task >> store_data_task
163
164 with DAG('DAG_2_FINAL', default_args=default_args, schedule_interval='@hourly', catchup=False) as dag2:
165     :
166     wait_for_dag1 = ExternalTaskSensor(
167         task_id = 'wait_for_scraping',
168         # The external tag ID should be same as DAG1
169         external_dag_id = 'DAG_1_FINAL',
170         # Wait for the last task of DAG1
171         external_task_id = 'store_news_data',
172         timeout = 600,
173

```

Figure 12: DAG 1 and 2 with arguments taken from the user

```
58     poke_interval = 60
59 )
60
61 # Wait for status file
62 wait_for_status = FileSensor(
63     task_id = 'wait_for_status_file',
64     # Use 'file' as connection ID - cope this from DAG Dashboard - admin - connections
65     fs_conn_id = postgres_connection_ID,
66     # For macbook, the opt argument is required.
67     filepath = '/opt/airflow/dags/status',
68     poke_interval = 30,
69     timeout = 600
70 )
71
72 # Process new entries and send email
73 notify_task = PythonOperator(
74     task_id = 'Module6',
75     python_callable = module6,
76     provide_context = True
77 )
78
79 wait_for_dag1 >> wait_for_status >> notify_task
```