**INFO-H420 - Management of Data Science and Business Workflows**

# Assignment 3:Workflows with Apache Airflow

Jintao Ma (ID 000586557)

Xianyun Zhuang (ID 000586733)

Professor

Dimitris Sacharidis

November 2023

# Contents

# Abstract

In this report, we have implemented the Assignment 3.

(i) In Exercise1, we finished the process_web_log data pipeline process with 4 tasks.

(ii) In Exercise2, we did test runs for each of the tasks we defined and a test run for the workflow. Additionally, we triggered the workflow and monitored a few runs. Document the test runs with our findings or observations from the runs.

(iii) In Exercise3, we added to the workflow a task that implemented sending a message to Slack after the last task.

At last, we uploaded the code to Github:

`https://github.com/woshimajintao/ETL-Pipeline-with-Apache-Airflow`

*Question: Define the DAG and workflow. Create 4 tasks to scan file, extract data, transform data, and load data*

The following is our implementation process and code snippet of process_web_log.

## 1.1 Environment Configuration:

1. **Docker:** Because both of us used the Windows system we had to install the Docker Desktop and Docker compose as the Figure1.1:

```
C:\Users\Jintao1999>docker-compose --version
Docker Compose version v2.23.0-desktop.1

C:\Users\Jintao1999>docker --version
Docker version 24.0.6, build ed223bc
```

Figure 1.1: Docker Set

2. **YAML File:** Docker-compose.yaml is a configuration file used by Docker Compose to define and run multi-container Docker applications. In this YAML file, we specified all the containers (services) required for an application, along with our configurations such as the image to be used, ports to be mapped, environment variables, volumes, and inter-container dependencies. Docker Compose utilizes this file to orchestrate the creation, configuration, and startup of all the specified services, making the management of complex applications with multiple components more manageable and efficient.

Here we used the airflow 2.7.3.[1]because it was a stable version.

We thought one noticed part of the docker-compose file was the volume. These configurations map the local directory to /opt/airflow/ inside the container. It was related to the path setting of the ETL pipeline.

```
1   volumes:
2     - ${AIRFLOW_PROJ_DIR:-.}/dags:/opt/airflow/dags
3     - ${AIRFLOW_PROJ_DIR:-.}/logs:/opt/airflow/logs
4     - ${AIRFLOW_PROJ_DIR:-.}/config:/opt/airflow/config
5     - ${AIRFLOW_PROJ_DIR:-.}/plugins:/opt/airflow/plugins
```

Listing 1.1: Volume Setting

3. **Make the Python File:** We found the Bashoperator was easy to use and implement 3 tasks of the whole ETL process.[2]And we used the PythonOperator to implement the message task. After we finish the Python file, we need to put it in the dag folder of Airflow.

4. **Run the Airflow:** We can use the Windows cmd command to open the airflow on Docker1.2:

```
1   D:\Apache\Airflow>docker compose up -d
```

Listing 1.2: Open Airflow



Figure 1.2: open Airflow

Then we clicked the `http://localhost:8080/` to access the Web UI of Airflow1.3.

Also, we can use the Windows cmd command to close the airflow on Docker1.4:

```
1   D:\Apache\Airflow>docker compose down
```

Listing 1.3: close Airflow

---

[1]https://airflow.apache.org/docs/apache-airflow/2.7.3/docker-compose.yaml
[2]https://airflow.apache.org/docs/apache-airflow/stable/howto/operator/bash.html

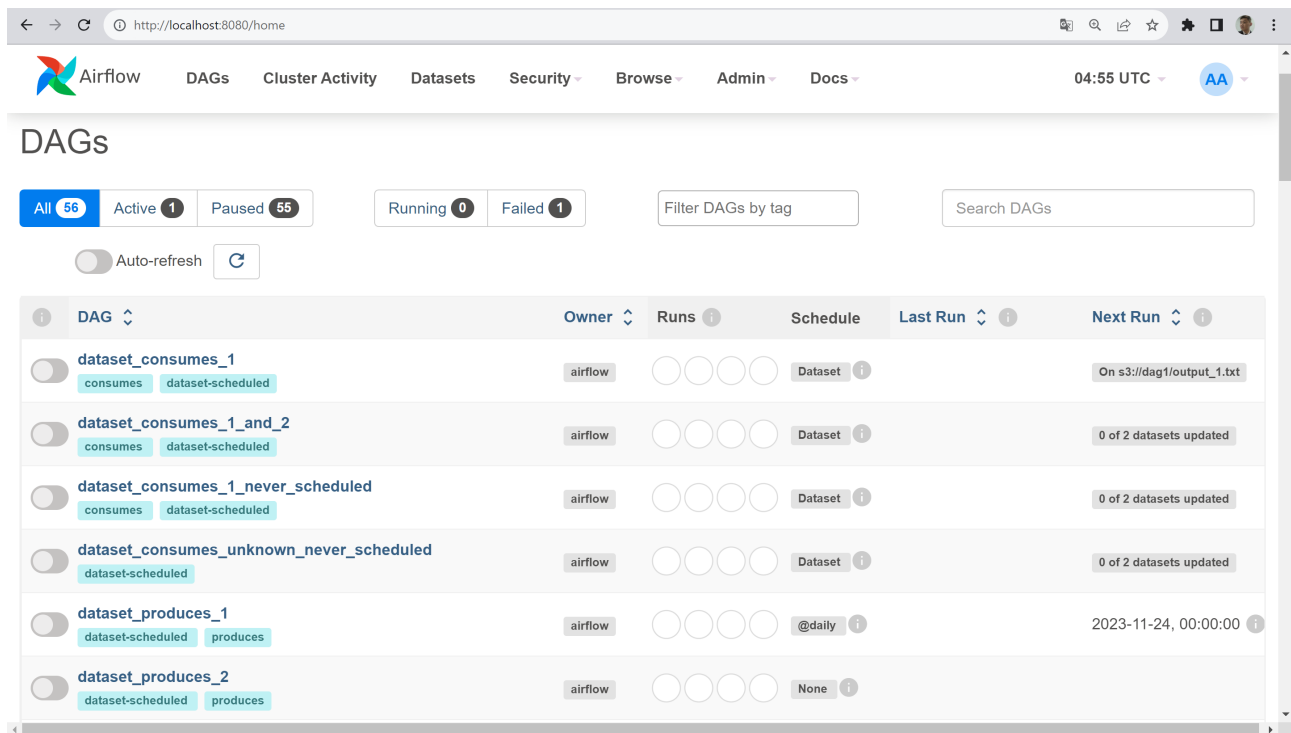Figure 1.3: Airflow Web UI



Figure 1.4: Close Airflow

## 1.2 Define the DAG:

The assignment is related to data science workflow monitor[Dum18]. All graphs were directed acyclic graphs (DAGs). This is our DAG definition block:

```
with DAG(
    dag_id='process_web_log',
    start_date=datetime(2023, 11, 18),
    default_args=default_args,
    schedule_interval=timedelta(days=1),
    catchup=False,
```

```
7 ) as dag:
```

Listing 1.4: DAG definition

The dag_id was the name of the DAG. The schedule_interval was to set the scheduling interval of the DAG to once a day.

## 1.3   Create a Task to Scan for a Log:

We used the FileSensor[3] to implement the scan_for_log task. The FileSensor is a built-in sensor in Apache Airflow that waits for a file or a group of files to be present. The sensor checks for the file's presence in the filesystem and ensures that the file is not being written to before returning.

```
1 scan_for_file = FileSensor(
2         task_id='scan_for_file',
3         filepath='/opt/airflow/dags/the_logs/log.txt',
4         fs_conn_id='my_file_system',  #connection ID of the filesystem
5         timeout=300,
6         mode='poke',
7         poke_interval=60,
8 )
```

Listing 1.5: scan_for_log

Please note that the fs_conn_id parameter[4] is the connection ID of the filesystem where the file will be located. We needed to set up this connection in Airflow's UI before running the DAG 1.5.

## 1.4   Create a Task to Extract Data:

We found the Bashoperator[5] was easy to use and implement the whole ETL process because it fits well with docker. So the three tasks were implemented by Bashoperaor.

Grep [6]is a useful command to search for matching patterns in a file. This command would extract all IP addresses from log.txt and write them into extracted_data.txt, separated by a newline.

```
1 extract_data = BashOperator(
2         task_id='extract_data',
```

---

[3]https://airflow.apache.org/docs/apache-airflow/stable/howto/operator/file.html
[4]https://www.restack.io/docs/airflow-faq-howto-operator-file-01
[5]https://airflow.apache.org/docs/apache-airflow/stable/howto/operator/bash.html
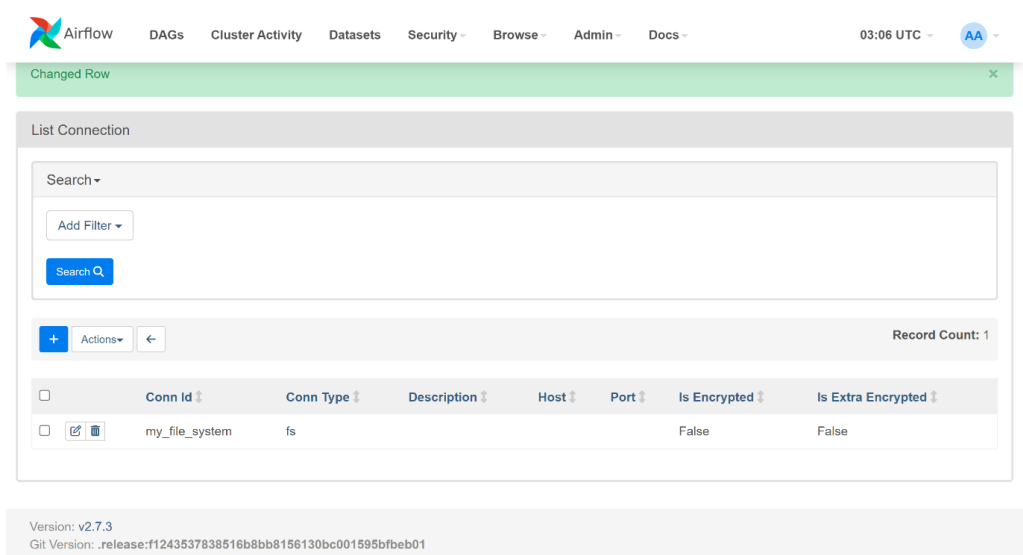[6]https://www.freecodecamp.org/news/grep-command-in-linux-usage-options-and-syntax-examples/

Figure 1.5: fs_conn_id setting

```
3        bash_command='grep -E -o "([0-9]{1,3}[\.]){3}[0-9]{1,3}" /opt/
    airflow/dags/the_logs/log.txt > /opt/airflow/dags/the_logs/extracted_data
    .txt',
4        dag=dag
5  )
```

Listing 1.6: extract_data

## 1.5   Create a Task to Transform Data:

Here we used two kinds of Command to implement the transform task[7]. Both of them could filter out all the occurrences of IP address 198.46.149.143 from "extracted_data.txt" and save the output to a file named "transformed_data.txt".

### 1.5.1   awk:

The awk is a widely used command for text processing.[8]  When using awk, we are able to select data – one or more pieces of individual text – based on a pattern we provide.

```
1  transform_data = BashOperator(
2        task_id='transform_data',
3        bash_command='awk \'$1 != "198.46.149.143"\' /opt/airflow/dags/
    the_logs/extracted_data.txt > /opt/airflow/dags/the_logs/transformed_data
    .txt',
```

---

[7]https://stackoverflow.com/questions/13579929/how-to-delete-lines-from-file-with-sed-awk
[8]https://www.freecodecamp.org/news/the-linux-awk-command-linux-and-unix-usage-syntax-examples/

```
4          dag=dag
5  )
```

Listing 1.7: transform_data

### 1.5.2   sed:

Sed stands for stream editor. It is a non-interactive command-line editor. It is primarily used for editing and filtering files. Sed is most often used in scripting, where it is used for editing either a single file or multiple files sequentially. Sed can be used to perform a sequence of editing actions on a single file or a group of files. It can be used to delete lines and words in a file, replace words, and even add text at specific locations in a file [9].

```
1  transform_data = BashOperator (
2         task_id='transform_data',
3         bash_command="sed '/198.46.149.143/d' /opt/airflow/dags/the_logs/
   extracted_data.txt> /opt/airflow/dags/the_logs/transformed_data.txt",
4         dag=dag
5  )
```

Listing 1.8: transform_data

## 1.6   Create a Task to Load the Data:

Here we used the Tar Command to implement the load task[10]. The 'tar' stands for tape archive, which is used to create Archive and extract the Archive files. tar command in Linux is one of the important commands that provides archiving functionality in Linux. We can use the Linux tar command to create compressed or uncompressed Archive files and also maintain and modify them.

The tar command uses the following flags to customize the command input[11]:

'-c': Creates a new archive.

'-z': Reduces the size of a given file.

'-v': Displays verbose output, showing the progress of the archiving process.

'-f': Specifies the filename of the archive.

```
1  load_data = BashOperator (
```

---

[9]https://tecadmin.net/sed-command-to-delete-line-in-file/google_vignette
[10]https://www.geeksforgeeks.org/tar-command-linux-examples/
[11]https://www.freecodecamp.org/news/how-to-compress-files-in-linux-with-tar-command/

```
2        task_id='load_data',
3        bash_command='tar -czvf weblog.tar /opt/airflow/dags/the_logs/
    transformed_data.txt',
4        dag=dag
5 )
```

Listing 1.9: load_data

## 1.7   Define the Workflow:

At last, we defined the workflow that executes the aforementioned tasks in sequence[12]. The load_data task depends on the transform_data task, the transform_data task depends on the extract_data task, and the extract_data task depends on the scan_for_file task.

```
1 scan_for_file >> extract_data >> transform_data >>load_data
```

Listing 1.10: Dependency Relationship

Putting all of the pieces together, we had our completed DAG.

---

[12]https://airflow.apache.org/docs/apache-airflow/stable/tutorial/pipeline.html

*Question: Do a test run for each of the tasks you defined. Once all work as expected, do a test run for the workflow. Finally, trigger/run the workflow and monitor a few runs. In the report, document the test runs, and include any findings or observations that you may have from the runs.*

## 2.1 Test run for each of the tasks:

First, comment out all tasks in the Python file except for the task currently needed for the test run. The Figure2.1 below shows the result for running only the 'scan' task.
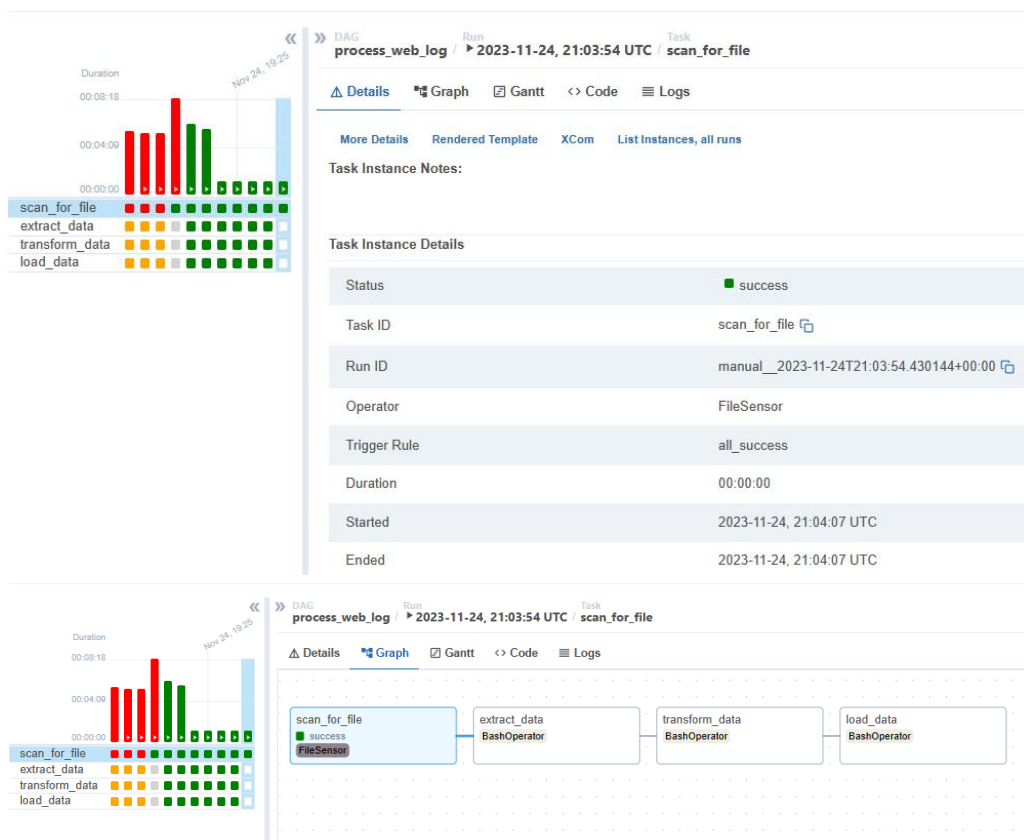


Figure 2.1: Test run result of 'scan' task

We can see from the Details and Graph that the 'scan' task ran successfully. Therefore, we performed the same operation on the other three tasks. The results of the test runs for the three tasks are as follows from the Figure2.2 to the Figure2.4:
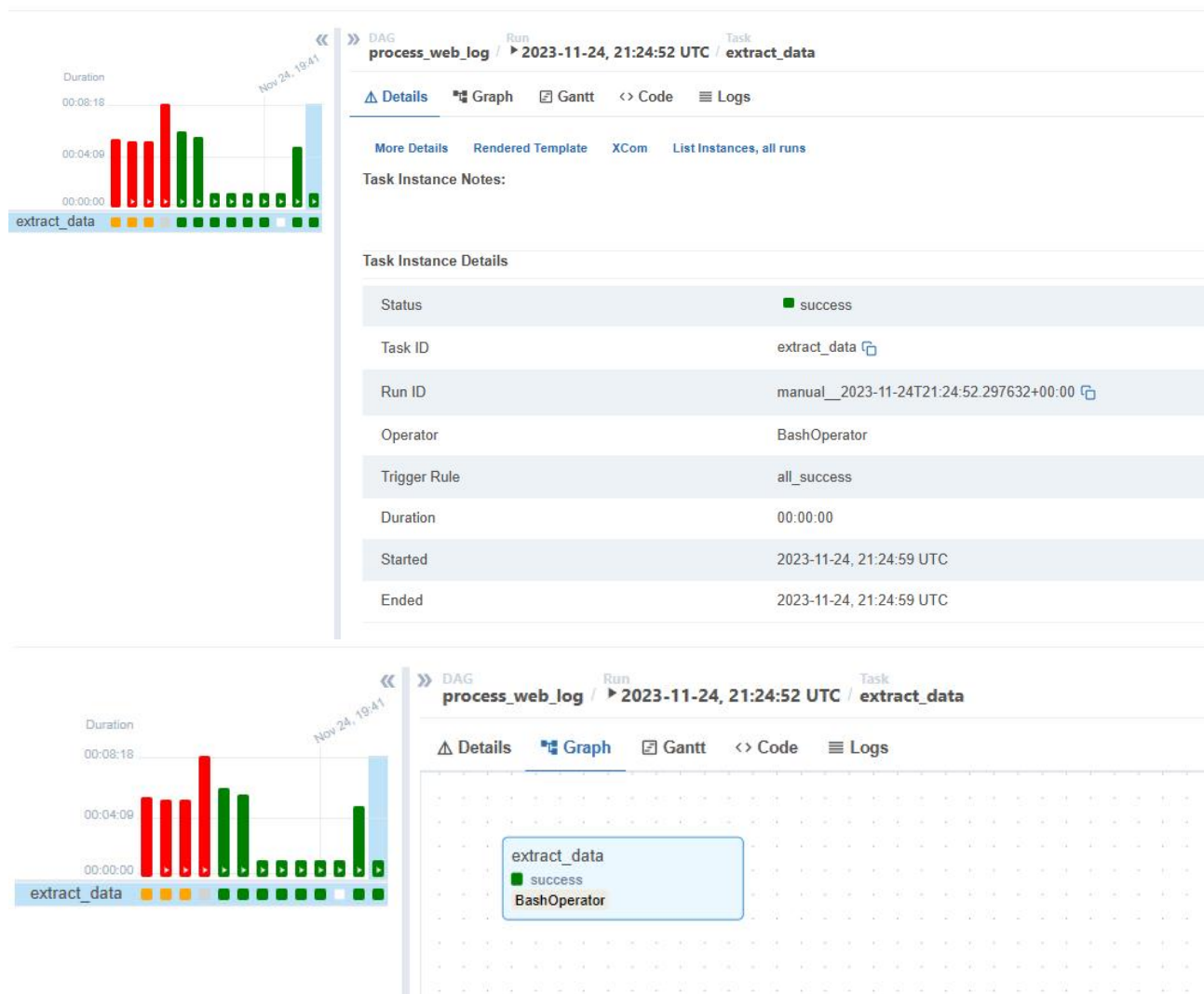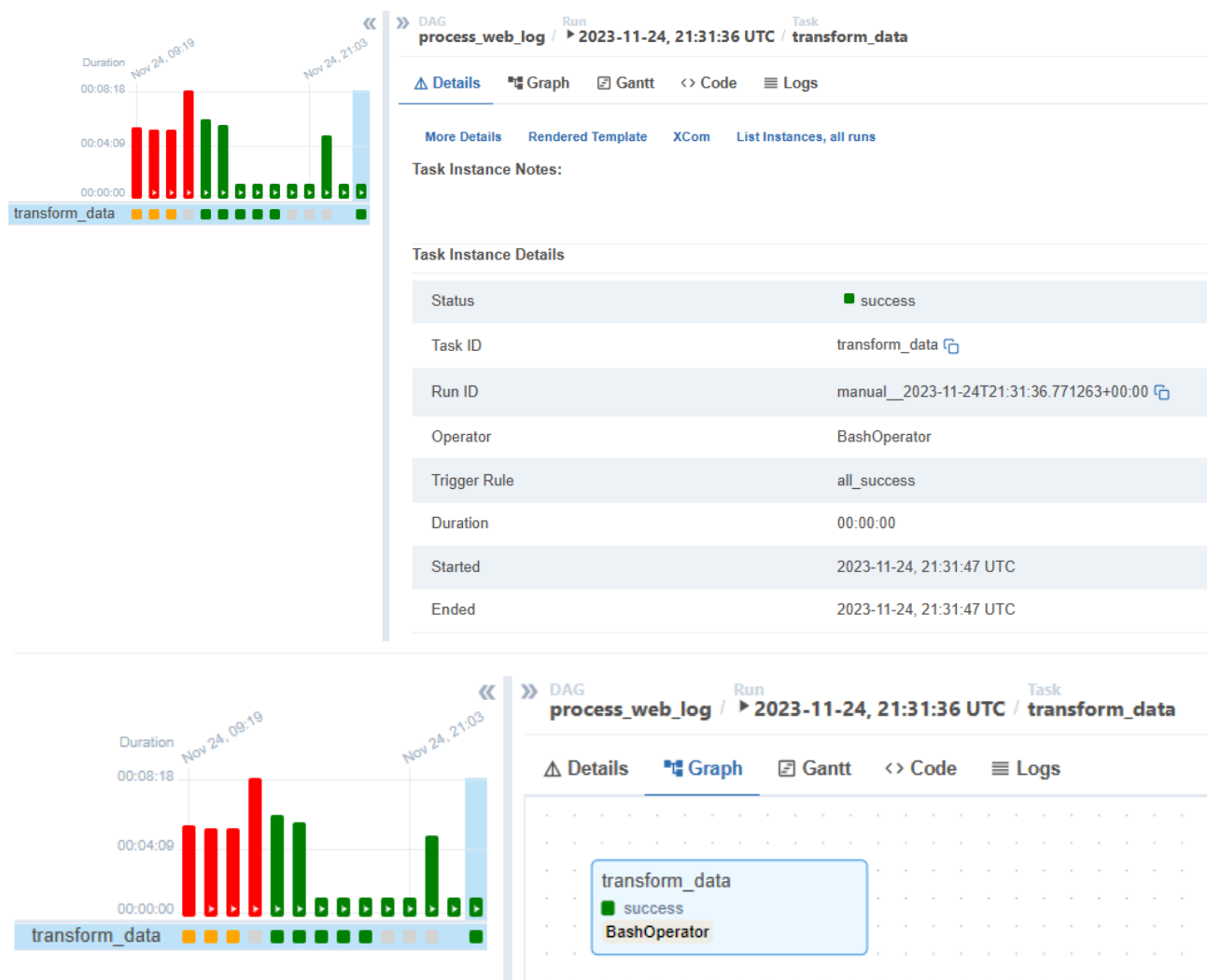


Figure 2.2: Test run result of 'extract' task

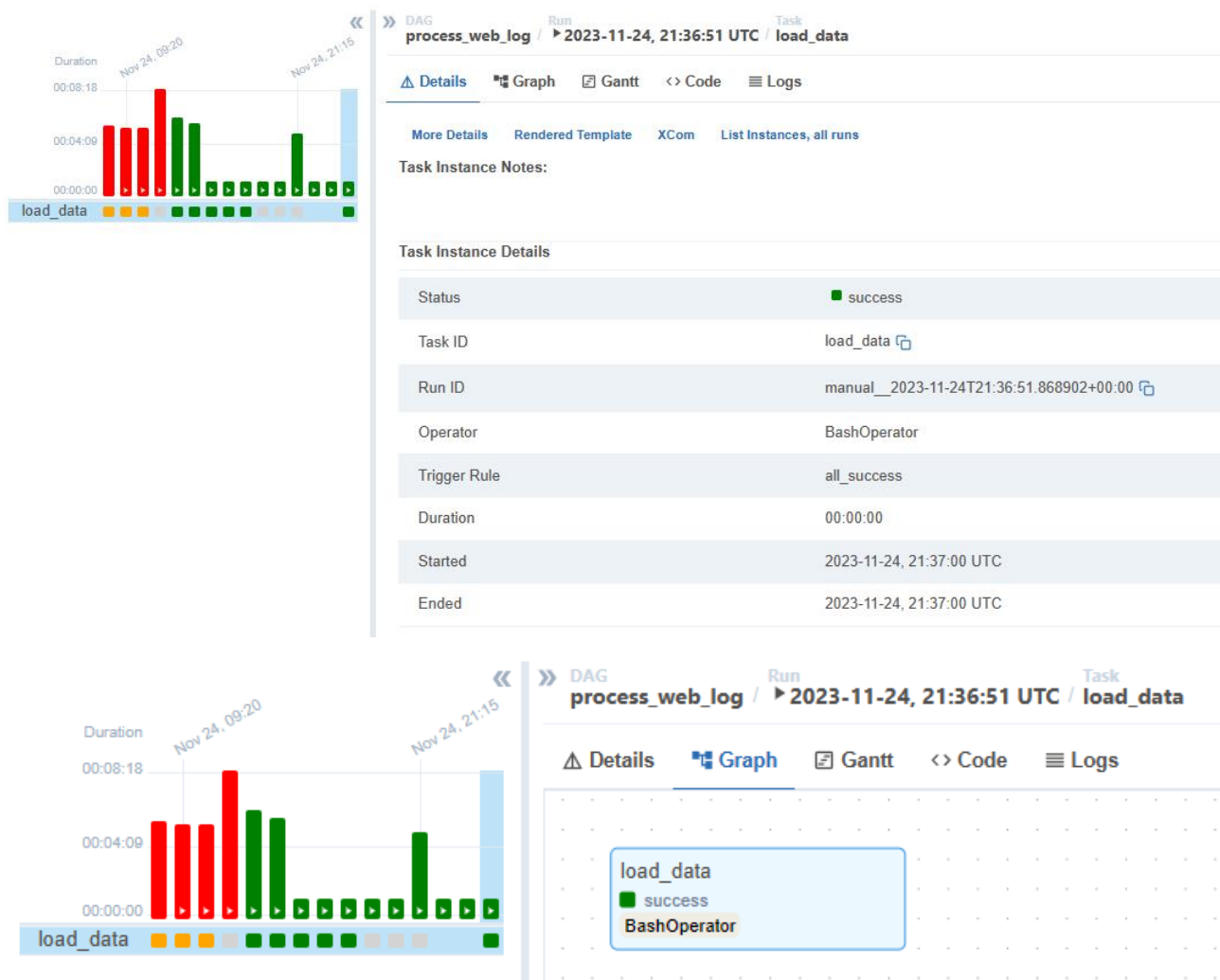Figure 2.3: Test run result of 'transform' task

Figure 2.4: Test run result of 'load data' task

## 2.2 Test run for the workflow:

We added all task codes to the DAG for a test run, which was executed successfully. The results are as shown in Figure2.5.
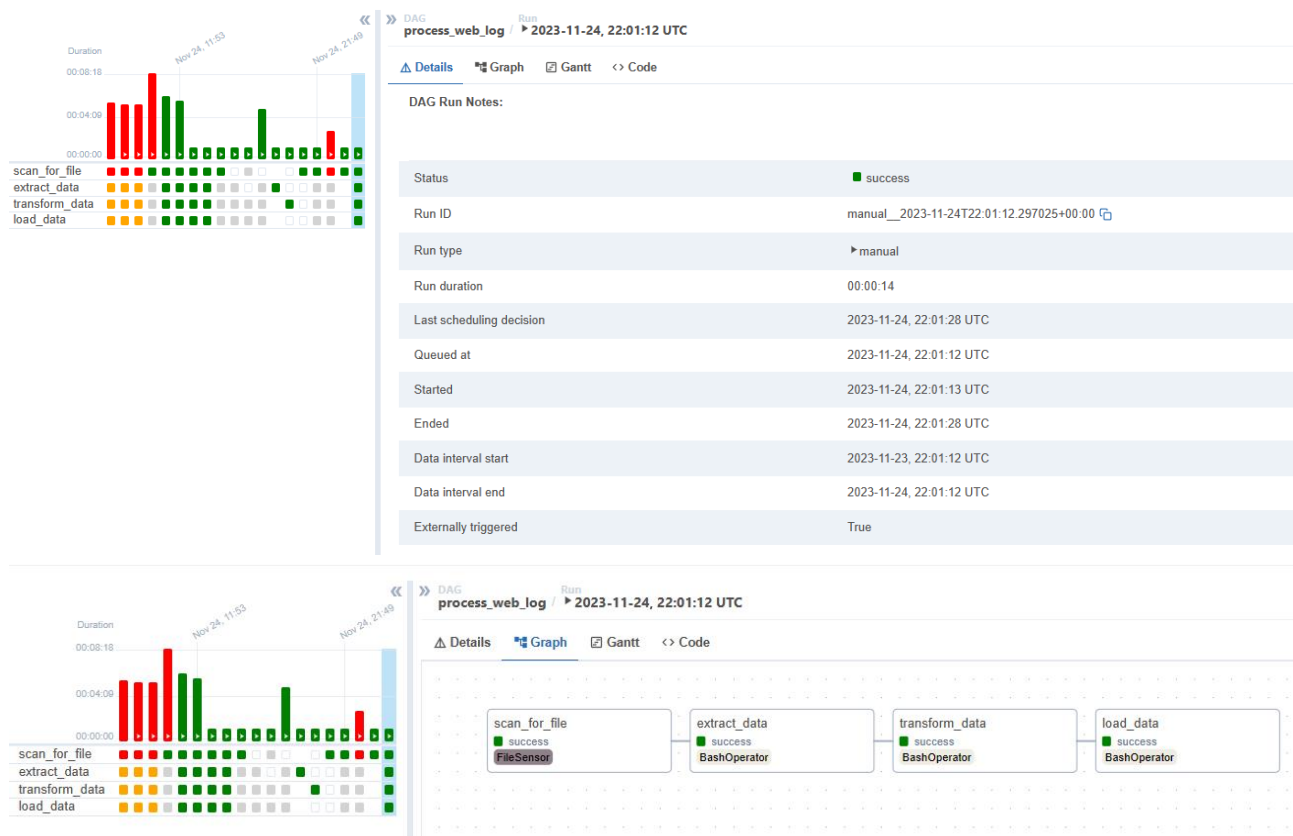
Figure 2.5: Test run result of the workflow

## 2.3 Trigger the workflow and monitor a few runs:

For testing purposes, we conducted two sets of tests by adjusting the execution intervals to 1 minute and 30 seconds, respectively.

When the interval was set to 1 minute, the DAG ran successfully but it took a long time to open the DAG and some of the tasks encountered 'up for retry' situations, resulting in increased processing time as shown in Figure2.6.

The scheduler does not appear to be running. Last heartbeat was received 1 minute ago.
The DAGs list may not update, and new tasks will not be scheduled.

DAG: process_web_log

Schedul

Grid | Graph | Calendar | Task Duration | Task Tries | Landing Times | Gantt | Details | <> Code | Audit Log

2023/11/24 22:30:02 | 25 | All Run Types | All Run States | Clear Filters

Press shift + / for Shortcuts

deferred failed queued removed restarting running scheduled skipped success up_for_re

DAG "process_web_log" seems to be missing from DagBag.

The scheduler does not appear to be running. Last heartbeat was received 2 minutes ago.
The DAGs list may not update, and new tasks will not be scheduled.
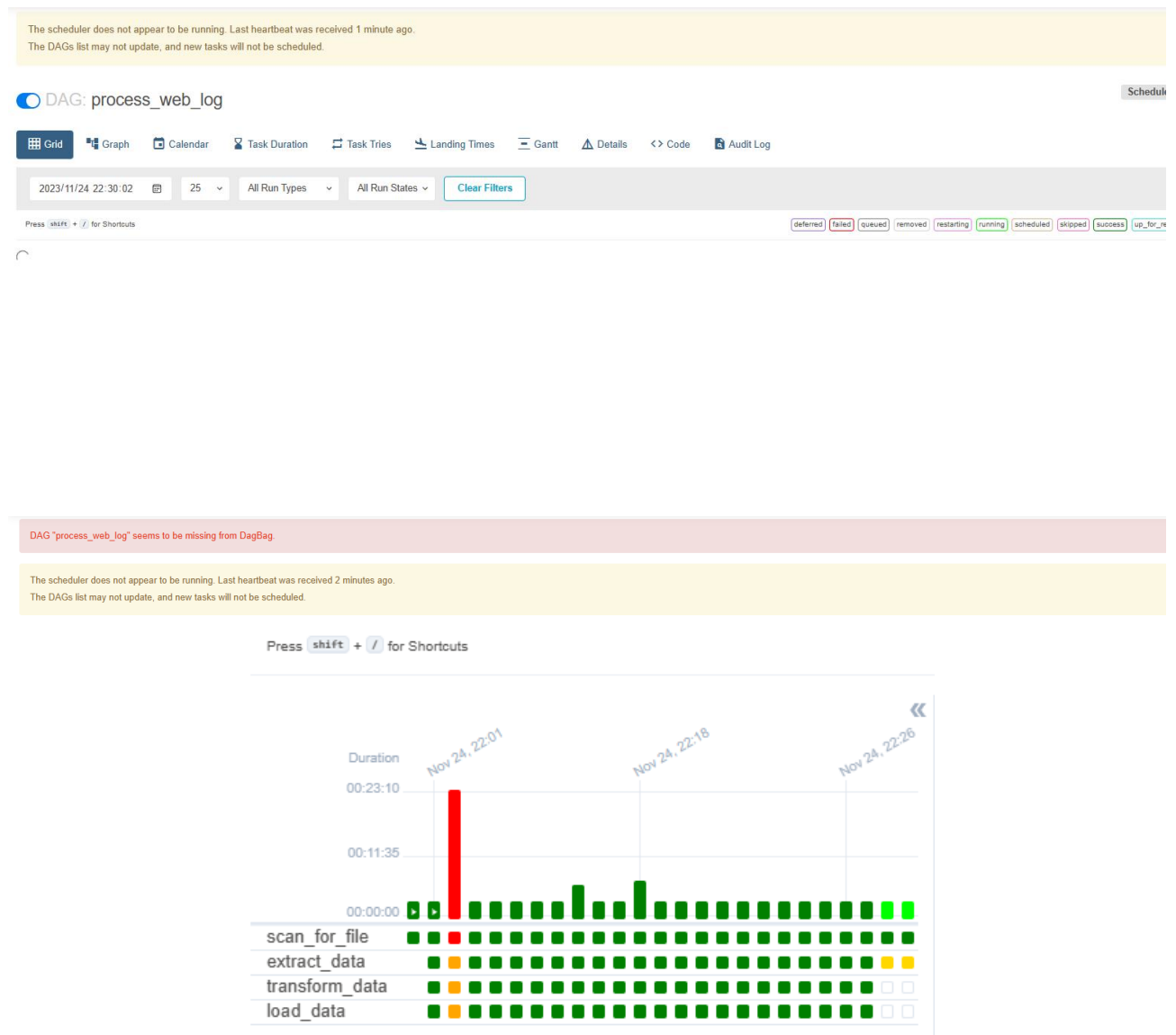
Press shift + / for Shortcuts

Figure 2.6: Test run result of triggering the workflow with 1-minute interval

When the interval was set to 30 seconds, the localhost page became inaccessible, and there were

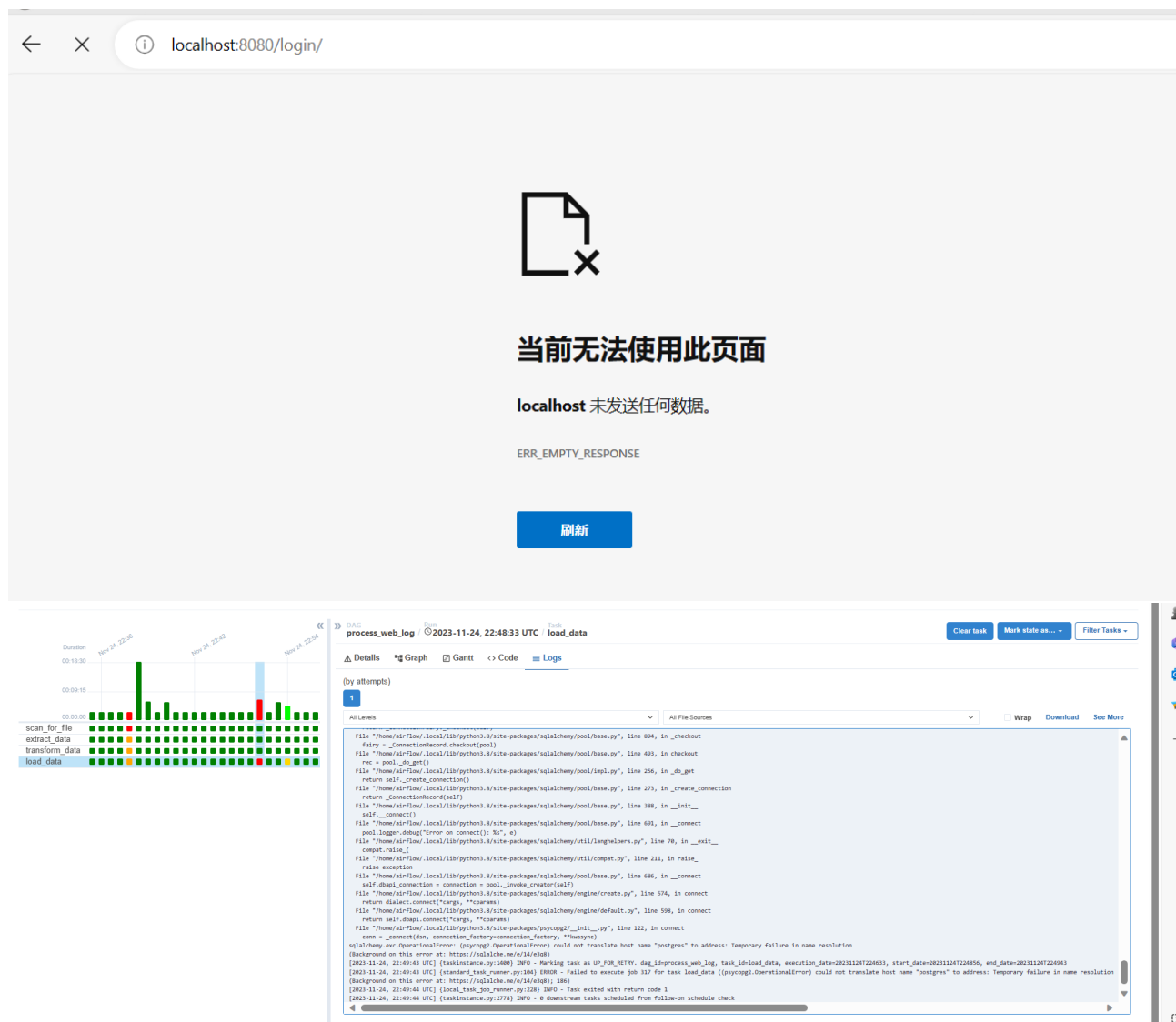also task failures as shown in Figure2.7.

Figure 2.7: Test run result of triggering the workflow with 30-seconds interval

## 2.4   Findings

During the test runs, we've noticed a significant impact on the success or failure of task execution based on computer performance and resource utilization. This could be attributed to several factors:

### 2.4.1   Computational Power and Speed:

Computational Power and Speed: The stronger the computer performance and faster the processor speed, the higher the efficiency in executing tasks. Some tasks require substantial computational power and speed, and inadequate computer performance might lead to slow or failed task execution.

### 2.4.2   Memory and Storage:

Certain tasks demand significant memory or disk space for data storage and processing. Insufficient memory or limited storage space on the computer can hinder task execution.

### 2.4.3   Resource Competition:

If other programs or processes are consuming a considerable amount of resources (e.g., large programs, background updates, or system processes), they might compete for the required resources, impacting the execution of tasks.

### 2.4.4   System Configuration and Environment:

System Configuration and Environment: Occasionally, task success also correlates with the operating system, software versions, and configurations of dependencies. Different system environments can influence task execution.

Therefore, optimizing computer performance and resource utilization by ensuring adequate computational power, memory, and storage, along with optimal system configurations and environments, can potentially enhance the success rate of task execution.

*3*

*Question: Add to the workflow an additional task after the last task. That task should send a message that the workflow was executed. You are free to define how and where this message is sent. For example, you can post to a Discord or Slack channel, write the message to a text file uploaded to Google Drive or Dropbox. In the report, document how you implemented the message task, and provide evidence that it executes.*

## 3.1    Create an Incoming Webhook in Slack

We started by setting up an Incoming Webhook in Slack.[1] To do this, create a new application on the Slack API website and enable the Incoming Webhook feature as shown in Figure3.1.



Figure 3.1: Test run result of triggering the workflow with 30-seconds interval

---

[1]https://www.restack.io/docs/airflow-knowledge-apache-providers-slack-webhook-http-pypi-operator

## 3.2 Create a connection to Slack

We used the Incoming Webhook URL obtained from Slack as the host while adding a new connection[2] in Airflow as shown in Figure3.2.



Figure 3.2: Test run result of triggering the workflow with 30-seconds interval

## 3.3 Create a new hook within Airflow

Firstly, we've defined a class called SlackHook, which inherits from Airflow's BaseHook class. This class includes methods for communicating with Slack [3].

### 3.3.1 init:

This is the initialization function that takes a parameter, slack_hook_conn_id, defaulted to 'slack_hook'. It calls the initialization function of the parent class and sets the passed parameter as the class attribute self.slack_hook_conn_id.

### 3.3.2 get_conn:

This method retrieves Slack's connection information. It uses the get_connection method, passing slack_hook_conn_id, and returns an object conn representing the Slack connection.

### 3.3.3 send_message:

This method is for sending messages to Slack. It first obtains the Slack connection's URL, constructs the data structure of the message to be sent, and uses requests. post method to send the message to Slack, and return the response result, response.

---

[2]https://airflow.apache.org/docs/apache-airflow-providers-slack/stable/connections/slack-incoming-webhook.html

[3]https://airflow.apache.org/docs/apacheairflow/stable/_api/airflow/hooks/base/index.html

```
1  class SlackHook(BaseHook):
2      def __init__(self, slack_hook_conn_id='slack_hook', *args, **kwargs):
3          super().__init__(*args, **kwargs)
4          self.slack_hook_conn_id = slack_hook_conn_id
5
6      def get_conn(self):
7          conn = self.get_connection(self.slack_hook_conn_id)
8          return conn
9
10     def send_message(self, message):
11         conn = self.get_conn()
12         slack_webhook_url = conn.host
13         slack_data = {'text': message}
14         response = requests.post(
15             slack_webhook_url, json=slack_data,
16             headers={'Content-Type': 'application/json'}
17         )
18         return response
```

Listing 3.1: SlackHook

The send_slack_message function stands alone; it instantiates the SlackHook class and sends a message to Slack using the send_message method. The message content is 'Workflow was executed successfully!'.

```
1  def send_slack_message():
2      slack_hook = SlackHook(slack_hook_conn_id='slack')
3      slack_hook.send_message('Workflow was executed successfully!')
```

Listing 3.2: send_slack_message

Finally, we defined a task within an Airflow DAG. It uses a PythonOperator named send_slack_message, setting the send_slack_message function as the python_callable to be executed when this task runs.

```
1  send_message_task = PythonOperator(
2          task_id='send_slack_message',
3          python_callable=send_slack_message,
4  )
```

Listing 3.3: send_message_task

## 3.4 Result of sending message to slack

After running the DAG, as shown in Figure3.3, we received the message 'Workflow was executed successfully!' in Slack, as shown in Figure3.4.
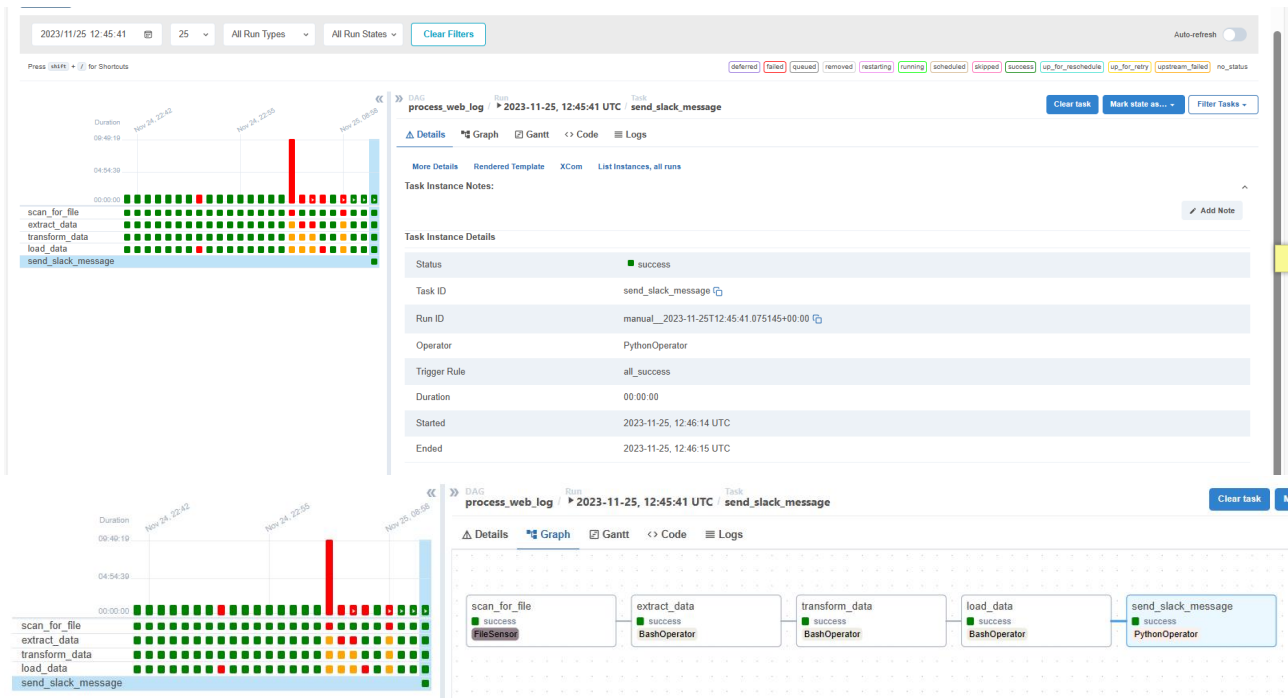


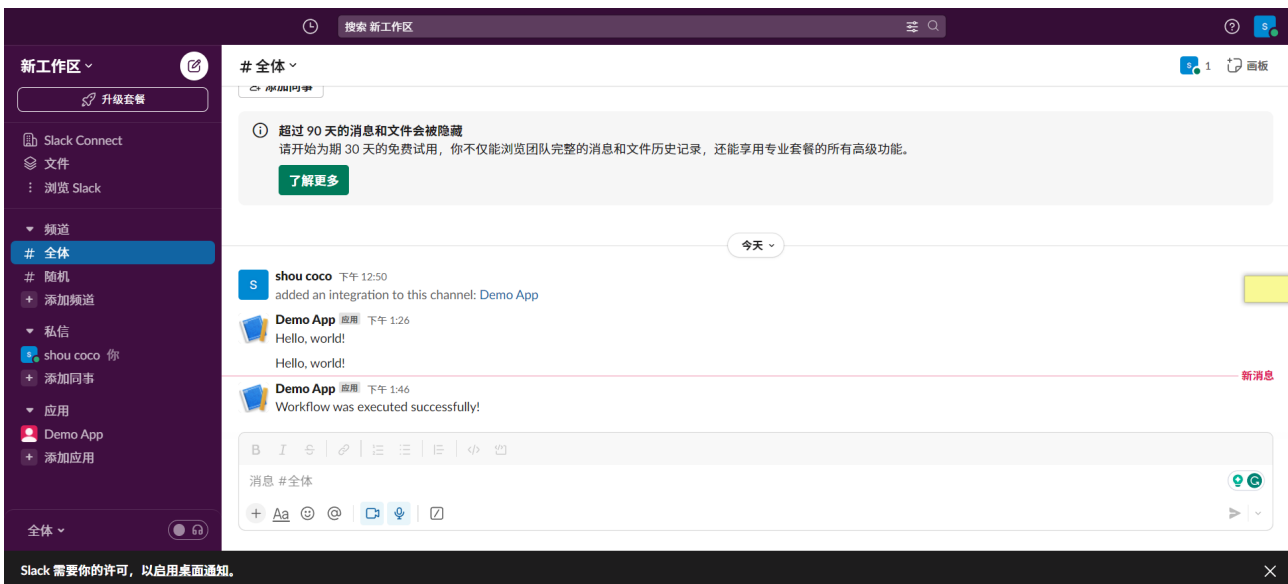Figure 3.3: running the DAG with send message to slack task



Figure 3.4: Test run result of triggering the workflow with 30-seconds interval

# Articles

[Dum18]    Dumas, Marlon (2018). "Fundamentals of Business Process Management". In: URL:
https://link.springer.com/book/10.1007/978-3-662-56509-4.