

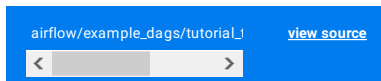
# Tutorial on the Taskflow API

This tutorial builds on the regular Airflow Tutorial and focuses specifically on writing data pipelines using the Taskflow API paradigm which is introduced as part of Airflow 2.0 and contrasts this with DAGs written using the traditional paradigm.

The data pipeline chosen here is a simple ETL pattern with three separate tasks for Extract, Transform, and Load.

## Example “Taskflow API” ETL Pipeline

Here is very simple ETL pipeline using the Taskflow API paradigm. A more detailed explanation is given below.



```

import json

from airflow.decorators import dag,
from airflow.utils.dates import days

# These args will get passed on to each task
# You can override them on a per-task basis
default_args = {
    'owner': 'airflow',
}

@dag(default_args=default_args, schedule_interval='@daily')
def tutorial_taskflow_api_etl():
    """
    ### TaskFlow API Tutorial Documentation
    This is a simple ETL data pipeline defined using
    the TaskFlow API using three simple tasks. For more
    Documentation that goes along with this tutorial, please see
    located
    [here](https://airflow.apache.org/docs/apache-airflow/2.0.0/tutorial.html)
    """

    @task()
    def extract():
        """
        ##### Extract task
        A simple Extract task to get data from a source into the
        pipeline. In this case, getting data from a hardcoded
        hardcoded JSON string.
        """
        data_string = '{"1001": 301.26, "1002": 42.99, "1003": 5.34}'

        order_data_dict = json.loads(data_string)
        return order_data_dict

    @task(multiple_outputs=True)
    def transform(order_data_dict: dict):
        """
        ##### Transform task
        A simple Transform task which takes the data from the
        computes the total order value.
        """
        total_order_value = 0

        for value in order_data_dict.values():
            total_order_value += value

        return {"total_order_value": total_order_value}

    @task()
    def load(total_order_value: float):
        """
        ##### Load task
        A simple Load task which takes the total order value
        instead of saving it to end storage.
        """

        print("Total order value is: {}".format(total_order_value))

        order_data = extract()
        order_summary = transform(order_data)
        load(order_summary["total_order_value"])

    tutorial_etl_dag = tutorial_taskflow_api_etl()

```

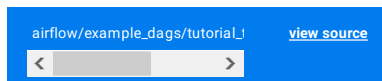
## It's a DAG definition file

If this is the first DAG file you are looking at, please note that this Python script is interpreted by Airflow and is a configuration file for your data pipeline. For a complete introduction to DAG files, please look at the core [Airflow tutorial](#) which covers DAG structure and definitions extensively.

## Instantiate a DAG

We are creating a DAG which is the collection of our tasks with dependencies between the tasks. This is a very simple definition, since we just want the DAG to be run when we set this up with Airflow, without any retries or complex scheduling. In this example, please notice that we are creating this DAG using the

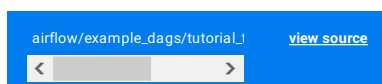
`@dag` decorator as shown below, with the python function name acting as the DAG identifier.



```
@dag(default_args=default_args, schedule_interval='@daily')
def tutorial_taskflow_api_etl():
    """
    ### TaskFlow API Tutorial Documentation
    This is a simple ETL data pipeline using the TaskFlow API using three simple tasks.
    The documentation that goes along with this tutorial is located
    [here](https://airflow.apache.org/docs/apache-airflow/2.0.0/tutorial_taskflow_api.html)
    """
```

## Tasks

In this data pipeline, tasks are created based on Python functions using the `@task` decorator as shown below. The function name acts as a unique identifier for the task.



```
@task()
def extract():
    """
    ##### Extract task
    A simple Extract task to get order data from the API. In this case, getting
    a hardcoded JSON string.
    """
    data_string = '{"1001": 301.28, "1002": 42.99, "1003": 5.04, "1004": 65.83, "1005": 21.30}'

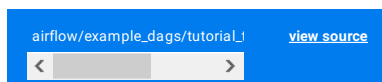
    order_data_dict = json.loads(data_string)
    return order_data_dict
```

The returned value, which in this case is a dictionary, will be made available for use in later tasks.

The Transform and Load tasks are created in the same manner as the Extract task shown above.

## Main flow of the DAG

Now that we have the Extract, Transform, and Load tasks defined based on the Python functions, we can move to the main part of the DAG.



```
order_data = extract()
order_summary = transform(order_data)
load(order_summary["total_order_value"])
```

That's it, we are done! We have invoked the Extract task, obtained the order data from there and sent it over to the Transform task for summarization, and then invoked the Load task with the summarized data. The dependencies between the tasks and the passing of data between these tasks which could be running on different workers on different nodes on the network is all handled by Airflow.

Now to actually enable this to be run as a DAG, we invoke the python function `tutorial_taskflow_api_etl` set up using the `@dag` decorator earlier, as shown below.

## But how?

For experienced Airflow DAG authors, this is startlingly simple! Let’s contrast this with how this DAG had to be written before Airflow 2.0 below:

```
import json

# The DAG object; we'll need this to instantiate a DAG
from airflow import DAG

# Operators; we need this to operate on XComs
from airflow.operators.python import PythonOperator
from airflow.utils.dates import days_ago

# These args will get passed on to each task
# You can override them on a per-task basis using args in the
# task dict, like:
# {'task_id': 'task_name', 'args': {'foo': 'bar', 'qux': 1999}}
default_args = {
    'owner': 'airflow',
}

with DAG(
    'tutorial_etl_dag',
    default_args=default_args,
    description='ETL DAG tutorial',
    schedule_interval=None,
    start_date=days_ago(2),
    tags=['example'],
) as dag:
    dag.doc_md = __doc__

    def extract(**kwargs):
        ti = kwargs['ti']
        data_string = '{"1001": 301.27, "1002": 42.19}'
        ti.xcom_push('order_data', data_string)

    def transform(**kwargs):
        ti = kwargs['ti']
        extract_data_string = ti.xcom_pull(task_id='extract')
        order_data = json.loads(extract_data_string)

        total_order_value = 0
        for value in order_data.values():
            total_order_value += value

        total_value = {"total_order_value": total_order_value}
        total_value_json_string = json.dumps(total_value)
        ti.xcom_push('total_order_value', total_value_json_string)

    def load(**kwargs):
        ti = kwargs['ti']
        total_value_string = ti.xcom_pull(task_id='transform')
        total_order_value = json.loads(total_value_string)

        print(total_order_value)
        extract_task = PythonOperator(
            task_id='extract',
            python_callable=extract,
        )

        extract_task.doc_md = """\
#### Extract task
A simple Extract task to get data ready for transform.
In this case, getting data is simulated by pushing a string to xcom.
This data is then put into xcom, so that the transform task can use it.
"""

        transform_task = PythonOperator(
            task_id='transform',
            python_callable=transform,
        )

        load_task = PythonOperator(
            task_id='load',
            python_callable=load,
        )

        extract_task.set_downstream(transform_task)
        transform_task.set_downstream(load_task)
```

```

        task_id='transform',
        python_callable=transform,
    )
    transform_task.doc_md = """\
#### Transform task
A simple Transform task which takes
and computes the total order value.
This computed value is then put into
"""

    load_task = PythonOperator(
        task_id='load',
        python_callable=load,
    )
    load_task.doc_md = """\
#### Load task
A simple Load task which takes in the
from xcom and instead of saving it to
"""

    extract_task >> transform_task >

```

All of the processing shown above is being done in the new Airflow 2.0 dag as well, but it is all abstracted from the DAG developer.

Let's examine this in detail by looking at the Transform task in isolation since it is in the middle of the data pipeline. In Airflow 1.x, this task is defined as shown below:

airflow/example\_dags/tutorial\_1 [view source](#)

```

def transform(**kwargs):
    ti = kwargs['ti']
    extract_data_string = ti.xcom_pull(
        task_id='extract_data',
    )
    order_data = json.loads(extract_data_string)

    total_order_value = 0
    for value in order_data.values():
        total_order_value += value

    total_value = {"total_order_value": total_order_value}
    total_value_json_string = json.dumps(total_value)
    ti.xcom_push('total_order_value', total_value_json_string)

```

As we see here, the data being processed in the Transform function is passed to it using Xcom variables. In turn, the summarized data from the Transform function is also placed into another Xcom variable which will then be used by the Load task.

Contrasting that with Taskflow API in Airflow 2.0 as shown below.

airflow/example\_dags/tutorial\_1 [view source](#)

```

@task(multiple_outputs=True)
def transform(order_data_dict: Dict[str, int]) -> Dict[str, int]:
    """
    #### Transform task
    A simple Transform task which
    computes the total order value
    """
    total_order_value = 0

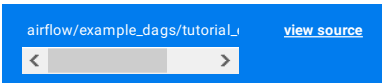
    for value in order_data_dict.values():
        total_order_value += value

    return {"total_order_value": total_order_value}

```

All of the Xcom usage for data passing between these tasks is abstracted away from the DAG author in Airflow 2.0. However, Xcom variables are used behind the scenes and can be viewed using the Airflow UI as necessary for debugging or DAG monitoring.

Similarly, task dependencies are automatically generated within TaskFlows based on the functional invocation of tasks. In Airflow 1.x, tasks had to be explicitly created and dependencies specified as shown below.



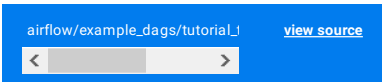
```
extract_task = PythonOperator(
    task_id='extract',
    python_callable=extract,
)
extract_task.doc_md = """\
#### Extract task
A simple Extract task to get data re
In this case, getting data is simula
This data is then put into xcom, so
"""

transform_task = PythonOperator(
    task_id='transform',
    python_callable=transform,
)
transform_task.doc_md = """\
#### Transform task
A simple Transform task which takes
and computes the total order value.
This computed value is then put into
"""

load_task = PythonOperator(
    task_id='load',
    python_callable=load,
)
load_task.doc_md = """\
#### Load task
A simple Load task which takes in th
from xcom and instead of saving it t
"""

extract_task >> transform_task >
```

In contrast, with the Taskflow API in Airflow 2.0, the invocation itself automatically generates the dependencies as shown below.



```
order_data = extract()
order_summary = transform(order_data)
load(order_summary["total_order_
```

## Multiple outputs inference

Tasks can also infer multiple outputs by using dict python typing.

```
@task
def identity_dict(x: int, y: int) -> Dict[str, int]:
    return {"x": x, "y": y}
```

By using the typing `Dict` for the function return type, the `multiple_outputs` parameter is automatically set to true.

Note, If you manually set the `multiple_outputs` parameter the inference is disabled and the parameter value is used.

## Adding dependencies to decorated tasks from regular tasks

The above tutorial shows how to create dependencies between python-based tasks. However, it is quite possible while writing a DAG to have some pre-existing tasks such as `BashOperator` or `FileSensor` based tasks which need to be run first before a python-based task is run.

Building this dependency is shown in the code below:

```
@task()
def extract_from_file():
    """
    ##### Extract from file task
    A simple Extract task to get data ready for the rest of the data
    pipeline, by reading the data from a file into a pandas dataframe
    """
    order_data_file = '/tmp/order_data.csv'
    order_data_df = pd.read_csv(order_data_file)

file_task = FileSensor(task_id='check_file', filepath='/tmp/order_data.csv')
order_data = extract_from_file()

file_task >> order_data
```

In the above code block, a new python-based task is defined as `extract_from_file` which reads the data from a known file location. In the main DAG, a new `FileSensor` task is defined to check for this file. Please note that this is a Sensor task which waits for the file. Finally, a dependency between this Sensor task and the python-based task is specified.

## What's Next?

You have seen how simple it is to write DAGs using the Taskflow API paradigm within Airflow 2.0. Please do read the [Concepts page](#) for detailed explanation of Airflow concepts such as DAGs, Tasks, Operators, etc, and the [Python task decorator](#) in particular.

More details about the Taskflow API, can be found as part of the Airflow Improvement Proposal [AIP-31: "Taskflow API" for clearer/simpler DAG definition](#) and specifically within the Concepts guide at [Concepts - Taskflow API](#).

[Previous](#)[Next](#)

Was this entry helpful?



Want to be a part of Apache Airflow?

[Join community](#)

License Donate Thanks  
Security  
© The Apache Software Foundation 2019

Apache Airflow, Apache, Airflow, the Airflow logo, and the Apache feather logo are either registered trademarks or trademarks of The Apache Software Foundation. All other products or name brands are trademarks of their respective holders, including The Apache Software Foundation.