



Jagoan
Hosting_

<#YangPentingYakin/>

WORKFLOW USING AIRFLOW_

29 / 02

2020



ABOUT ME

1. **Raja Fathurrahman**
2. I am **Physicist** (Computational and Instrumentation physics)
3. Data Analyst at **JagoanHosting**

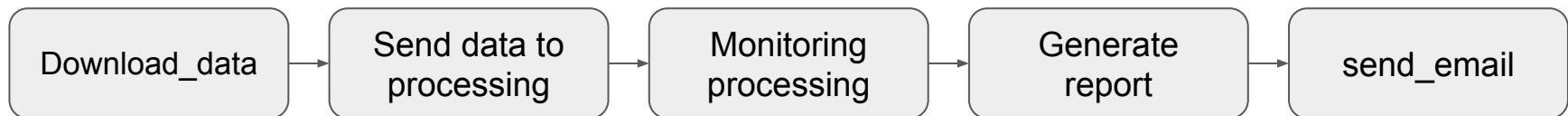


What is workflow ?

Workflows are the way process get work done, and can be illustrated as series of steps that need to be completed sequentially in a diagram or checklist.

- Sequence of **Tasks**
- Started on **schedule** or **trigger** by event
- Frequently used to handle data **pipelines**

Typical and Example what we can do with workflow



- ETL (Extract, Transform, Load)
- Data Warehousing
- A/B Testing
- Anomaly Detection
- Training Recommendation testing
- Orchestrating automated process



Workflow Manager



argo



Azkaban



Luigi



Jagoan
Hosting_

Open-source workflow orchestration tools and comparison

Tool	Paradigm	Workflow syntax/DSL	UI	Origin	Platform(s)	Popularity
Airflow	DAG configuration as code	Python	Mature UI	AirBnB, Apache	*nix	13481 Stars
Argo Workflows	Tasks described as containers	YAML-based templates	Mature UI	Intuit	Kubernetes	3521 Stars
Azkaban	Declarative with DSL	Flow 2.0 (YAML-like)	Mature UI	LinkedIn (Hadoop)	*nix	2727 Stars
Luigi	Object oriented (Task, Target)	Python	Mature UI	Spotify	*nix, Windows	12063 Stars
Oozie	Declarative with DSL	XML-based	Mature UI	Hadoop-only, Apache	*nix, Java	500 Stars



Why Apache airflow ?

We've also observed a general shift away from drag-and-drop ETL (Extract Transform and Load) tools towards a more programmatic approach. Product know-how on platforms like Informatica, IBM Datastage, Cognos, AbInitio or Microsoft SSIS isn't common amongst modern data engineers, and being replaced by more generic software engineering skills along with understanding of programmatic or configuration driven platforms like Airflow, Oozie, Azkabhan or Luigi. It's also fairly common for engineers to develop and manage their own job orchestrator/scheduler.

- ‘Maxime Beauchemin’

Airflow is a workflow scheduler to help with scheduling complex workflows and provide an easy way to maintain them



Introduction apache airflow

- **Programmatically** author
- **Schedule** and **Monitor** workflows
- More **Maintainable**
- **Versionable**
- **Testable**
- **Collaborative**
- **Open source**
- Rich **Web UI** for **Monitoring** and **Log**

How to install apache airflow

```
$ python3 --version  
Python 3.6.0  
$ virtualenv --vers  
15.1.0
```

```
$ cd /path/to/my/airflow/workspace  
$ virtualenv -p `which python3` venv  
$ source venv/bin/activate  
(venv) $
```

```
(venv) $ pip install airflow==1.8.0
```

1

2

3

```
(venv) $ airflow version
```

v1.8.0rc5+apache.incubating

```
(venv) $ cd /path/to/my/airflow/workspace  
(venv) $ mkdir airflow_home  
(venv) $ export AIRFLOW_HOME='pwd'/airflow_home
```

```
(venv) $ airflow initdb
```

```
(venv) $ airflow webserver
```

6

7



Airflow UI

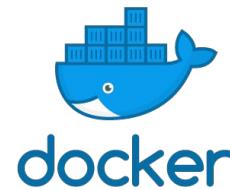
Airflow DAGs Data Profiling ▾ Browse ▾ Admin ▾ Docs ▾ About ▾ 2018-09-22 10:48:00 UTC ⌂

DAGs

Search:

		DAG	Schedule	Owner	Recent Tasks	Last Run	DAG Runs	Links
		example_bash_operator	0 0 * * *	airflow	(9)	(9)	(9)	
		example_branch_dop_operator_v3	/1 * * * *	airflow	(9)	(9)	(9)	
		example_branch_operator	@daily	airflow	(9)	(9)	(9)	
		example_http_operator	1 day, 0:00:00	airflow	(9)	(9)	(9)	
		example_kubernetes_executor	None	airflow	(9)	(9)	(9)	
		example_passing_params_via_test_command	/1 * * * *	airflow	(9)	(9)	(9)	
		example_python_operator	None	airflow	(9)	(9)	(9)	
		example_short_circuit_operator	1 day, 0:00:00	airflow	(9)	(9)	(9)	
		example_skip_dag	1 day, 0:00:00	airflow	(9)	(9)	(9)	

Alternative installing apache airflow

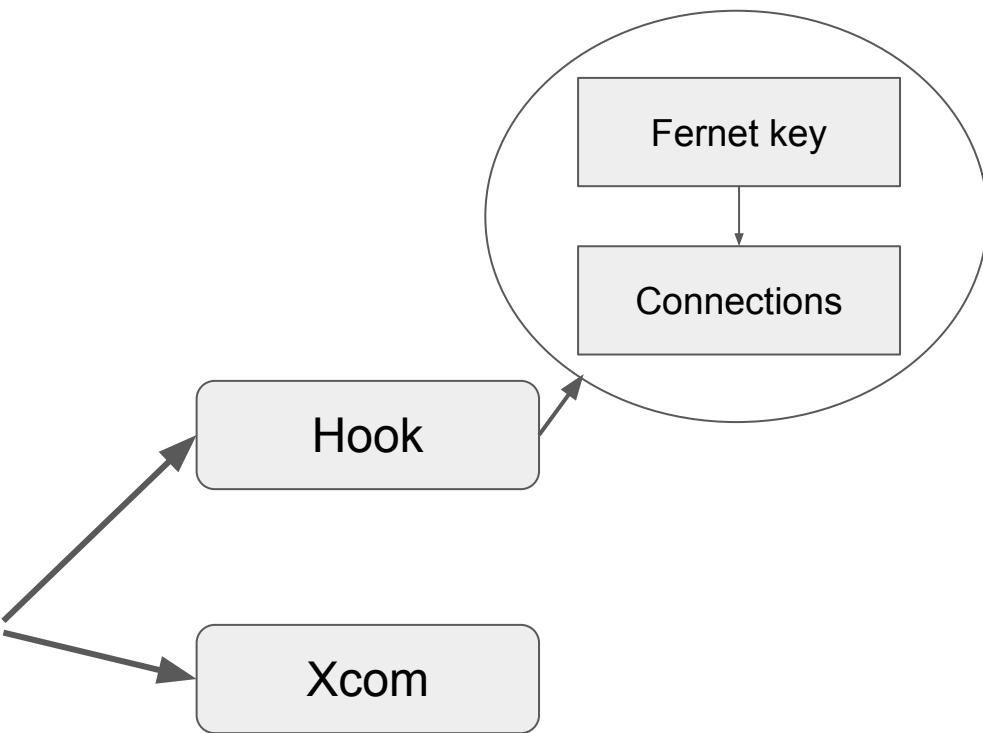


More easy installing airflow using docker just pull the images and run the container on docker and no need more configuration

```
docker pull puckel/docker-airflow
```

```
docker run -d -p 8080:8080 puckel/docker-airflow webserver
```

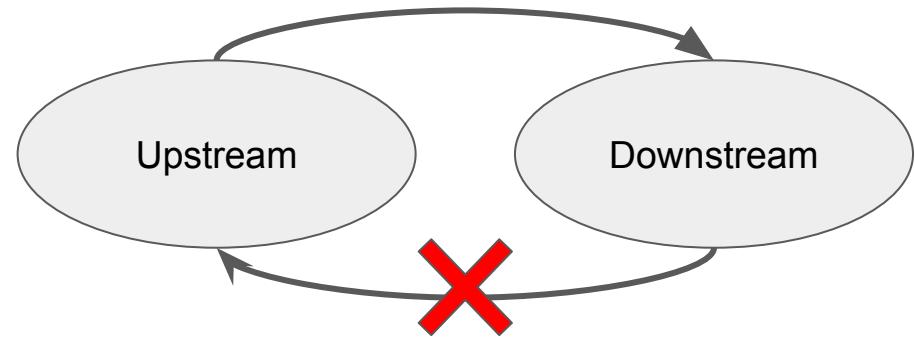
Workflow component in Airflow



Jagoan
Hosting_

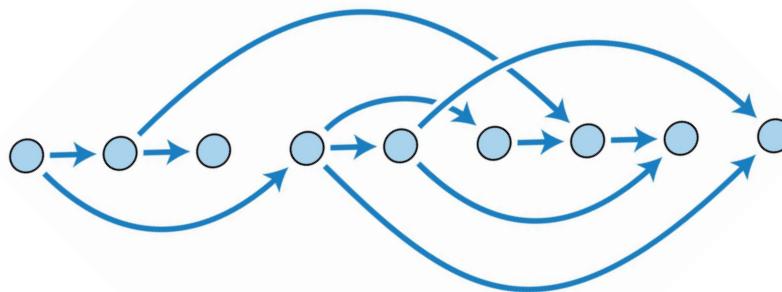
What workflow in Airflow?

This flow like a river flow from upstream to downstream



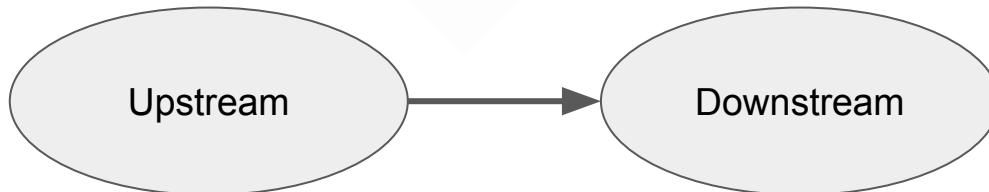
Workflow with Airflow

DAG concept using river flow analogy that is directed and not in cycle



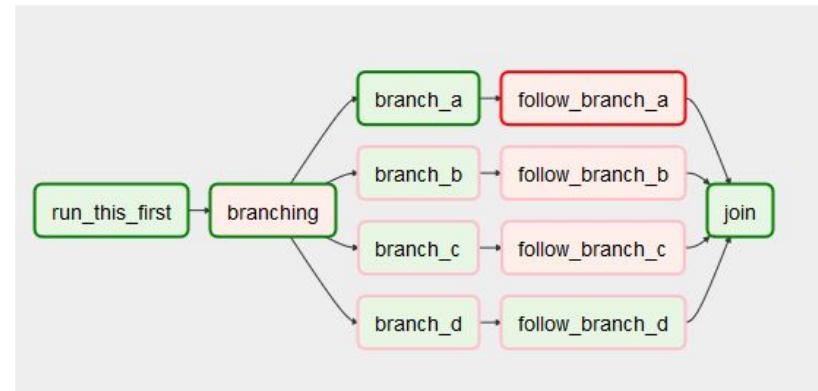
This is call a DAG (Directed Acyclic Graph)

DAG is the workflow in airflow



Workflow with Airflow

Like a river, they have a tributaries



In Airflow it call branch



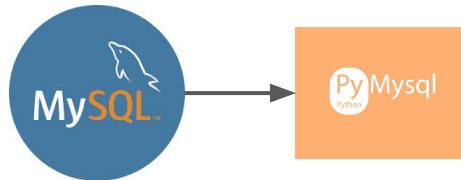
Jagoan
Hosting_

Workflow with apache Airflow

Hooks



Operators



Tasks

Process how to
pymysql connect
to mysql
database to get
data with query

Airflow Operators

- PythonOperator
 - BashOperator
 - EmailOperator
 - S3ToRedshiftTransfer
 - RedshiftToS3Transfer
 - BigQueryToCloudStorageOperator
 - GoogleCloudStorageToBigQueryOperator
 - SlackAPIPostOperator
 - MySqlToHiveTransfer
 - SimpleHttpOperator
 - DataprocOperationBaseOperator
 - JiraOperator
 - DockerOperator
 - PostgresOperator
- and many more...

Airflow Hooks

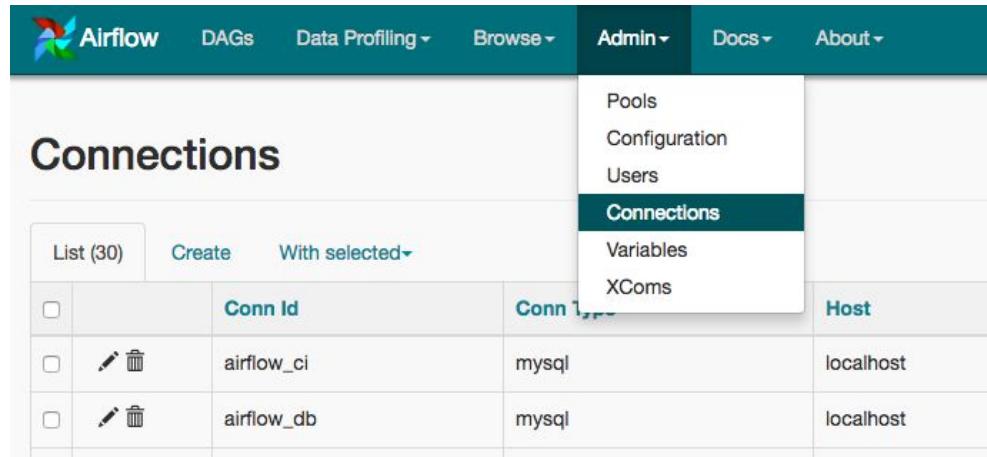
- AwsHook
 - JdbcHook
 - GoogleCloudStorageHook
 - CassandraHook
 - BigQueryHook
 - RedisHook
 - DatadogHook
 - AzureFileShareHook
 - GrpcHook
 - JiraHook
 - SSHHook
 - HiveCliHook
 - SlackHook
 - DockerHook
 - S3Hook
 - MySqlHook
 - PrestoHook
- and many more...

Connections (Hook) on Airflow

reduce connection
initialization in code

Connection Types

- Amazon Web Services Connection
- Google Cloud Platform Connection
- Google Cloud SQL Connection
- gRPC
- MySQL Connection
- Oracle Connection
- PostgresSQL Connection
- SSH Connection



The screenshot shows the Airflow web interface with the 'Connections' page selected. The top navigation bar includes links for Airflow, DAGs, Data Profiling, Browse, Admin (which is currently active), Docs, and About. The 'Connections' page has a sub-navigation menu with options: Pools, Configuration, Users, Connections (which is highlighted in blue), Variables, and XComs. Below the menu is a table with columns: Conn Id, Conn Type, and Host. There are three entries: airflow_ci (mysql, localhost), airflow_db (mysql, localhost), and another entry partially visible.

	Conn Id	Conn Type	Host
<input type="checkbox"/>	airflow_ci	mysql	localhost
<input type="checkbox"/>	airflow_db	mysql	localhost



Fernet key is important on Airflow connection

Cryptography concept on fernet key to secure your connection



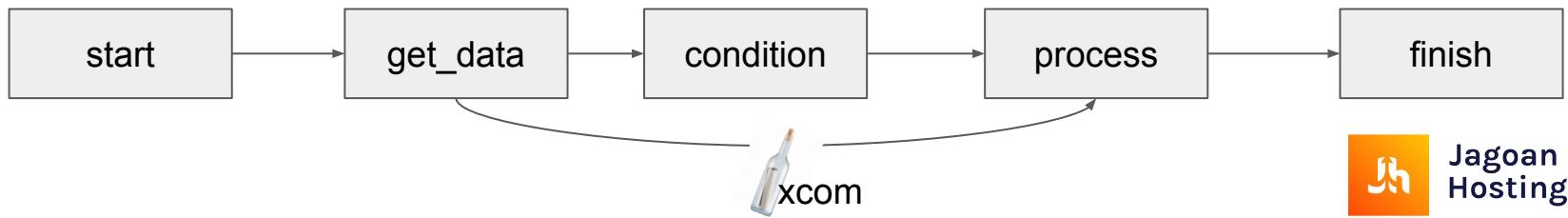
Static or Dynamic fernet key depend on yours

Workflow with apache Airflow

Xcom communication between operator



Xcom concept like a message in the bottle, it can be pull and push between tasks instance, and save the data as a pickle object and best suited for small pieces of data

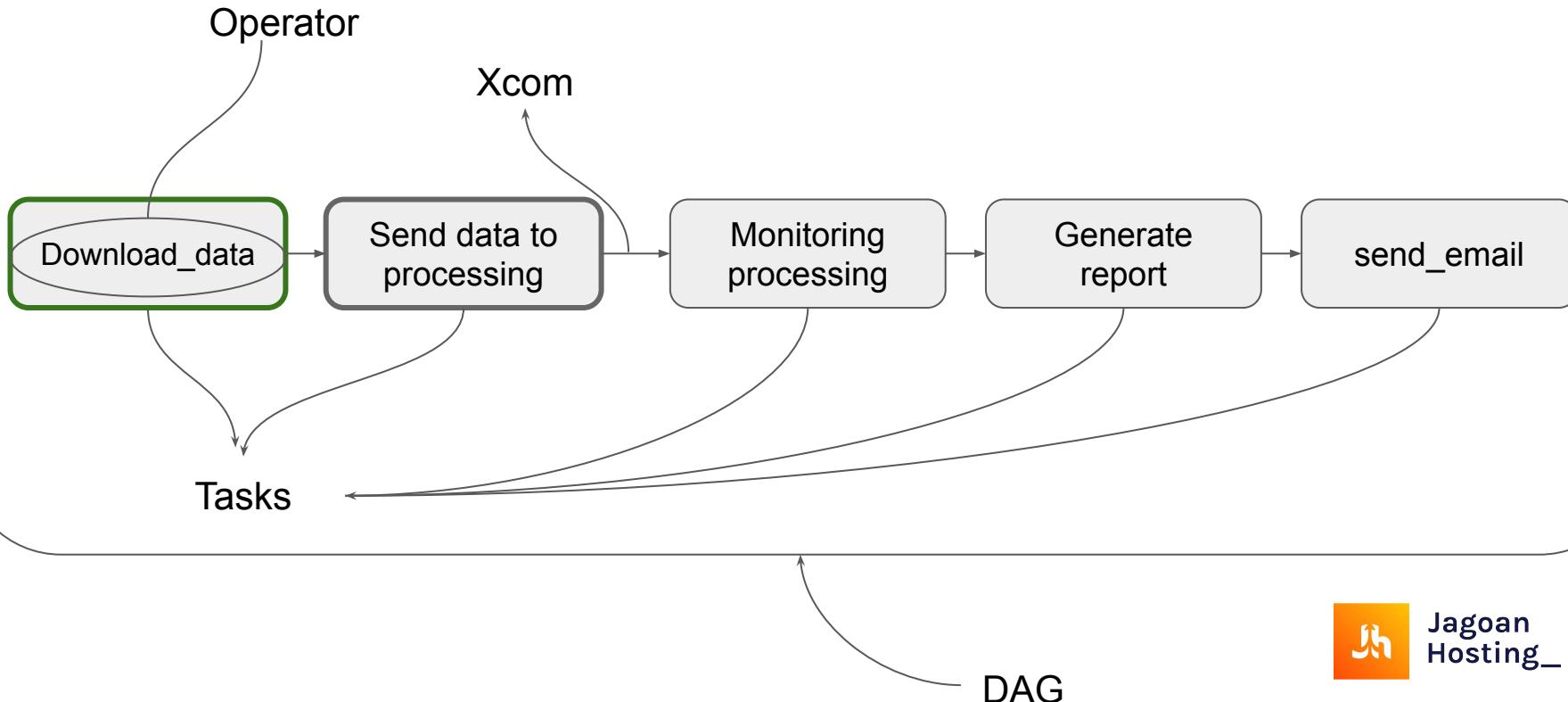


Workflow with apache Airflow

preset	meaning	cron
None	Don't schedule, use for exclusively "externally triggered" DAGs	
@once	Schedule once and only once	
@hourly	Run once an hour at the beginning of the hour	0 * * * *
@daily	Run once a day at midnight	0 0 * * *
@weekly	Run once a week at midnight on Sunday morning	0 0 * * 0
@monthly	Run once a month at midnight of the first day of the month	0 0 1 * *
@yearly	Run once a year at midnight of January 1	0 0 1 1 *

Single DAG Run will be created, with an execution_date of 2016-01-01, and the next one will be created just after midnight on the morning of 2016-01-03 with an execution date of 2016-01-02.

Workflow with apache Airflow



Jagoan
Hosting_

What value does Airflow add?

- Retries tasks elegantly
- Alerts on failure
- Re-run specific tasks in the DAG
- Supports distributed execution

<https://airflow.apache.org/docs/stable/concepts.html>



Demo

```
❸ printHelloworld.py /home/adib/jagoanhosting/airflow/dags/printHelloworld.py  finishProcess
#Basic Airflow Library
from airflow import DAG
from airflow.operators.python_operator import PythonOperator

from datetime import datetime, timedelta, date, timezone #library for time

default_args = {
    'owner': 'Hello',
    'depends_on_past': False,
    'start_date': datetime(2020, 2, 28),
    'retries': 1,
    'retry_delay': timedelta(minutes=1) }

dag = DAG('HelloWorld_01', default_args=default_args, schedule_interval=None)

def startProcess():
    print('FINISH')

def finishProcess():
    print('FINISH')

def printHelloworld(**kwargs):
    date = kwargs['ds']
    dates={}.format(date)
    execute = datetime.strptime(dates, '%Y-%m-%d')
    date_now = (execute).strftime("%Y-%m-%d")

    print("Hello world!!, This is", date_now)

t1 = PythonOperator(
    task_id= 'startProcess',
    python_callable=startProcess,
    dag=dag)

t2 = PythonOperator(
    task_id= 'printHelloworld',
    provide_context=True,
    python_callable=printHelloworld,
    dag=dag)

t3 = PythonOperator(
    task_id= 'finishProcess',
    python_callable=finishProcess,
    dag=dag)

t1 >> t2 >> t3
```

} Import module

} DAG

} Operator

} Tasks

Dependencies

```
xcom_example.py /home/adib/jagoanhosting/airflow/dags/xcom_example.py...
from airflow import DAG
from airflow.operators.python_operator import PythonOperator
from airflow.utils.dates import days_ago

args = {
    'owner': 'airflow',
    'start_date': days_ago(2)}

dag = DAG('example_xcom', default_args=args, schedule_interval="@once")

value_2 = {'a': 'b'}

def push(**kwargs):
    ti = kwargs['task_instance']
    value = [1, 2, 3]
    ti.xcom_push(key='value_1', value=value) → Xcom push

def push_by_returning(**kwargs):
    return value_2

def puller(**kwargs):
    ti = kwargs['ti']

    pulled_value_1 = ti.xcom_pull(key='value_1', task_ids='push')
    print('Xcom get value_1 : ', pulled_value_1)
    pulled_value_2 = ti.xcom_pull(task_ids='push_by_returning')
    print('Xcom get values_2', pulled_value_2)

    }
```

Xcom pull



Jagoan
Hosting_

Xcom Demo

```
def push(**kwargs):
    """Pushes an XCom without a specific target"""
    ti = kwargs['task_instance']
    value = [1, 2, 3]
    ti.xcom_push(key='value_1', value=value)
```

```
def puller(**kwargs):
    """Pull all previously pushed XComs and check if the pushed values match the pulled values."""
    ti = kwargs['ti']
    # get value_1 with key value pair
    pulled_value_1 = ti.xcom_pull(key='value_1', task_ids='push')
    print('Xcom get value_1 : ', pulled_value_1)
    # get value_2 without key value
    pulled_value_2 = ti.xcom_pull(task_ids='push_by_returning')
    print('Xcom get values_2', pulled_value_2)
```



Web UI (User Interface)

**JADI TUA ITU PASTI,
JADI KEREN ITU PILIHAN**

Mulai langkah Keren-mu dari Sekarang di :

bit.ly/skuymagang



Jagoan
Hosting -

Raja Fathurrahman
Data Analyst

Contact :
083129031420
raja@jagoanhosting.com