



# Manageable data pipelines with **Airflow** (and Kubernetes)



SQL

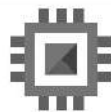
Google Cloud Platform



Bigtable



Spanner



Compute Engine

[AIRFLOW-3275] Add Google Cloud SQL Query operator (#4170)



potiuk authored and kaxil committed 17 days ago ✓



d7e80d6

[AIRFLOW-3345] Add Google Cloud Storage (GCS) operators for ACL (#4192) ...



sprzedwojski authored and kaxil committed 15 days ago ✓



0f7eca2





# Airflow

Airflow is a platform to programmatically author, schedule and monitor workflows.

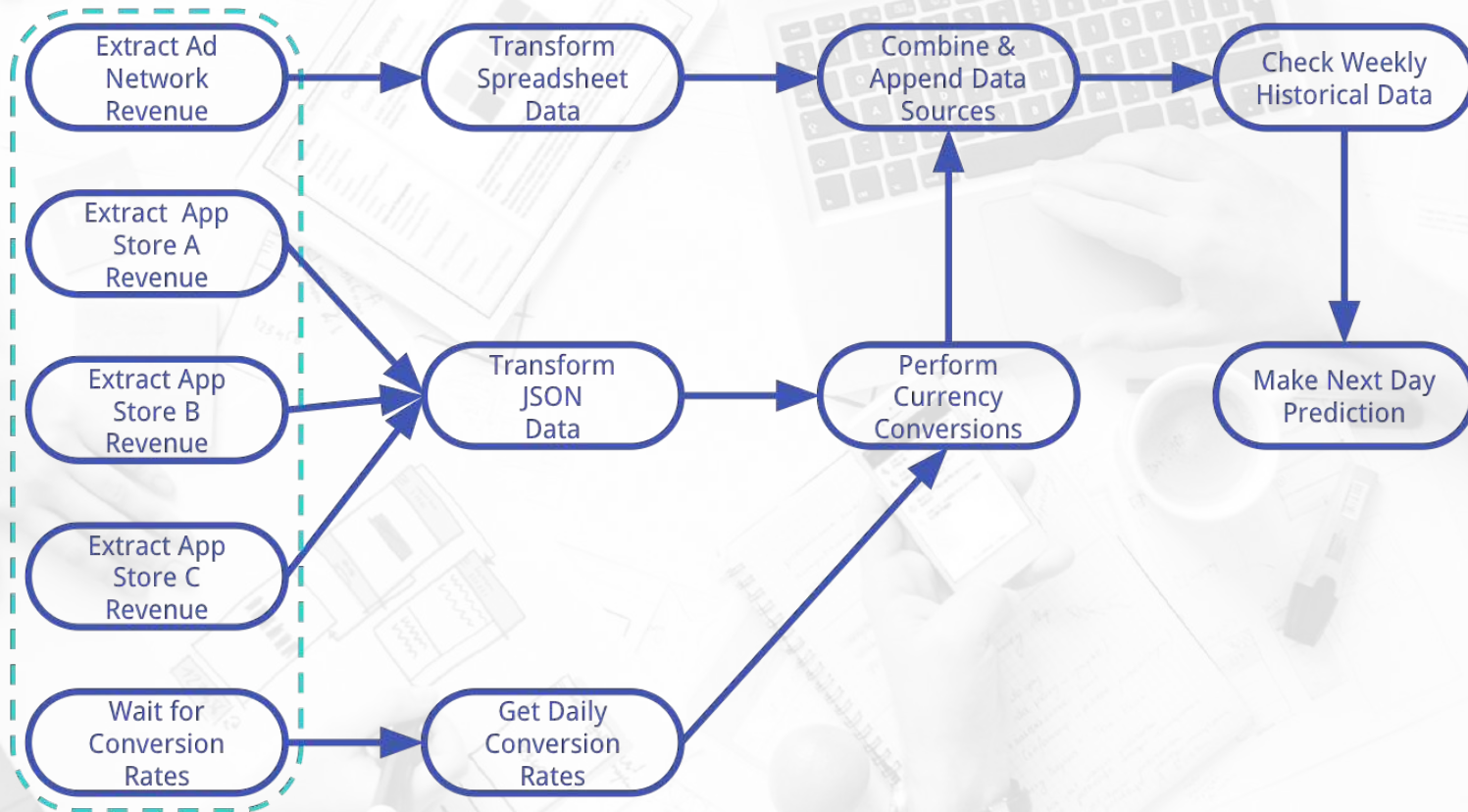
**Dynamic/Elegant**

**Extensible**

**Scalable**



# Workflows





# Companies using Airflow

(>200 officially)



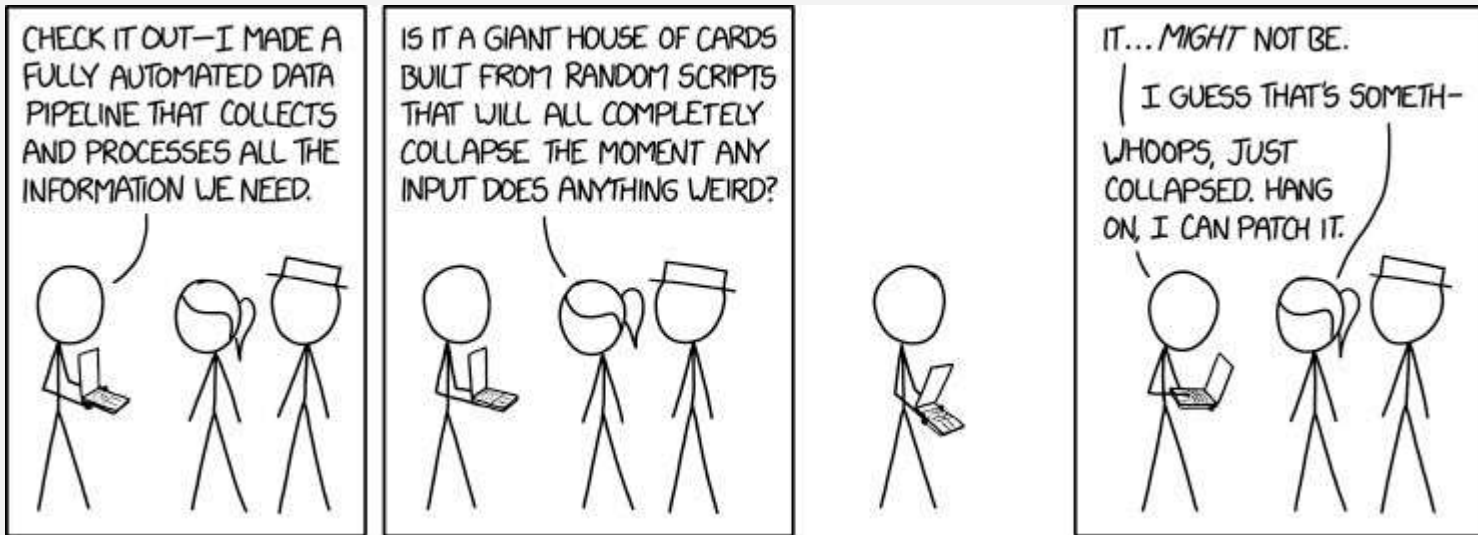
DigitalOcean





# Data Pipeline

<https://xkcd.com/2054/>



# Airflow vs. other workflow platforms

- Programming workflows
  - writing code not XML
  - versioning as usual
  - automated testing as usual
  - complex dependencies between tasks
- Managing workflows
  - aggregate logs in one UI
  - tracking execution
  - re-running, backfilling (run all missed runs)



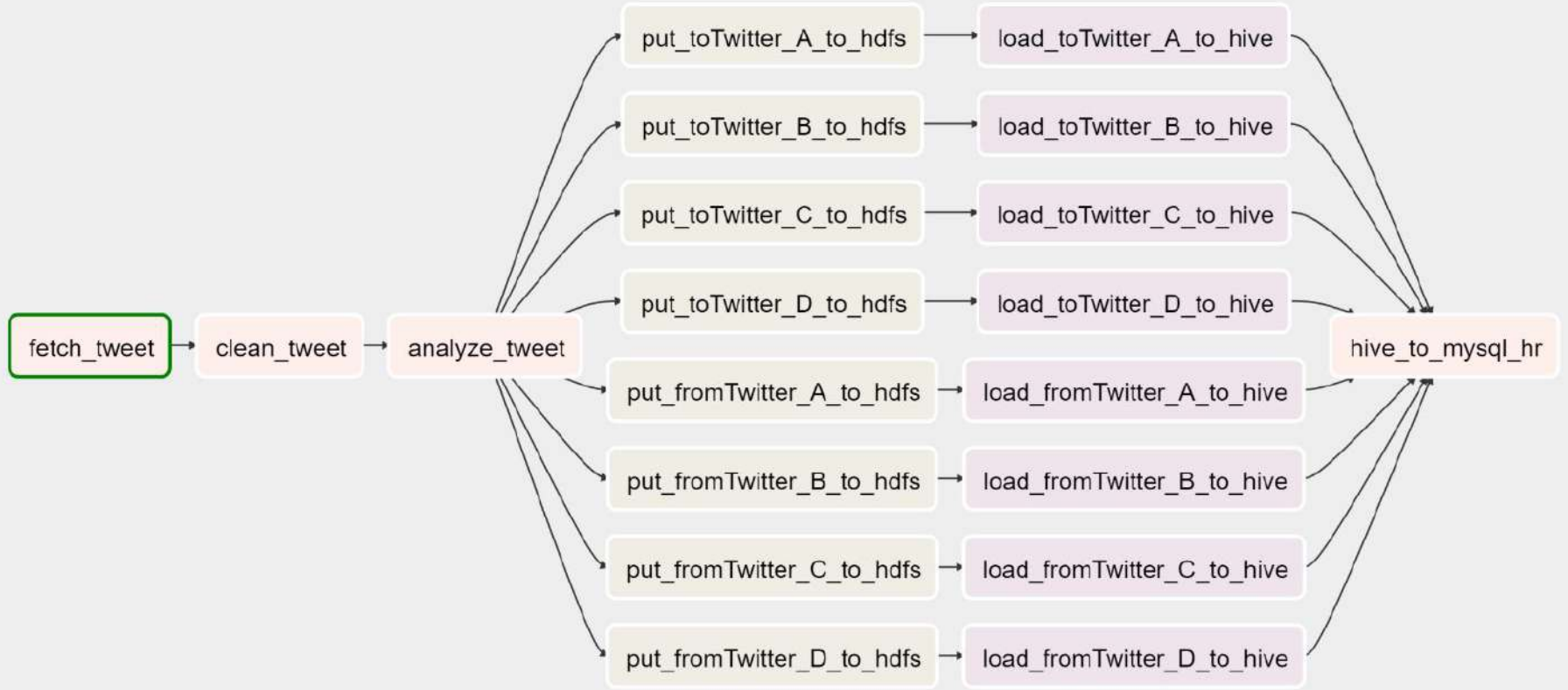
# Airflow use cases

- ETL jobs
- ML pipelines
- Regular operations:
  - Delivering data
  - Performing backups
- ...



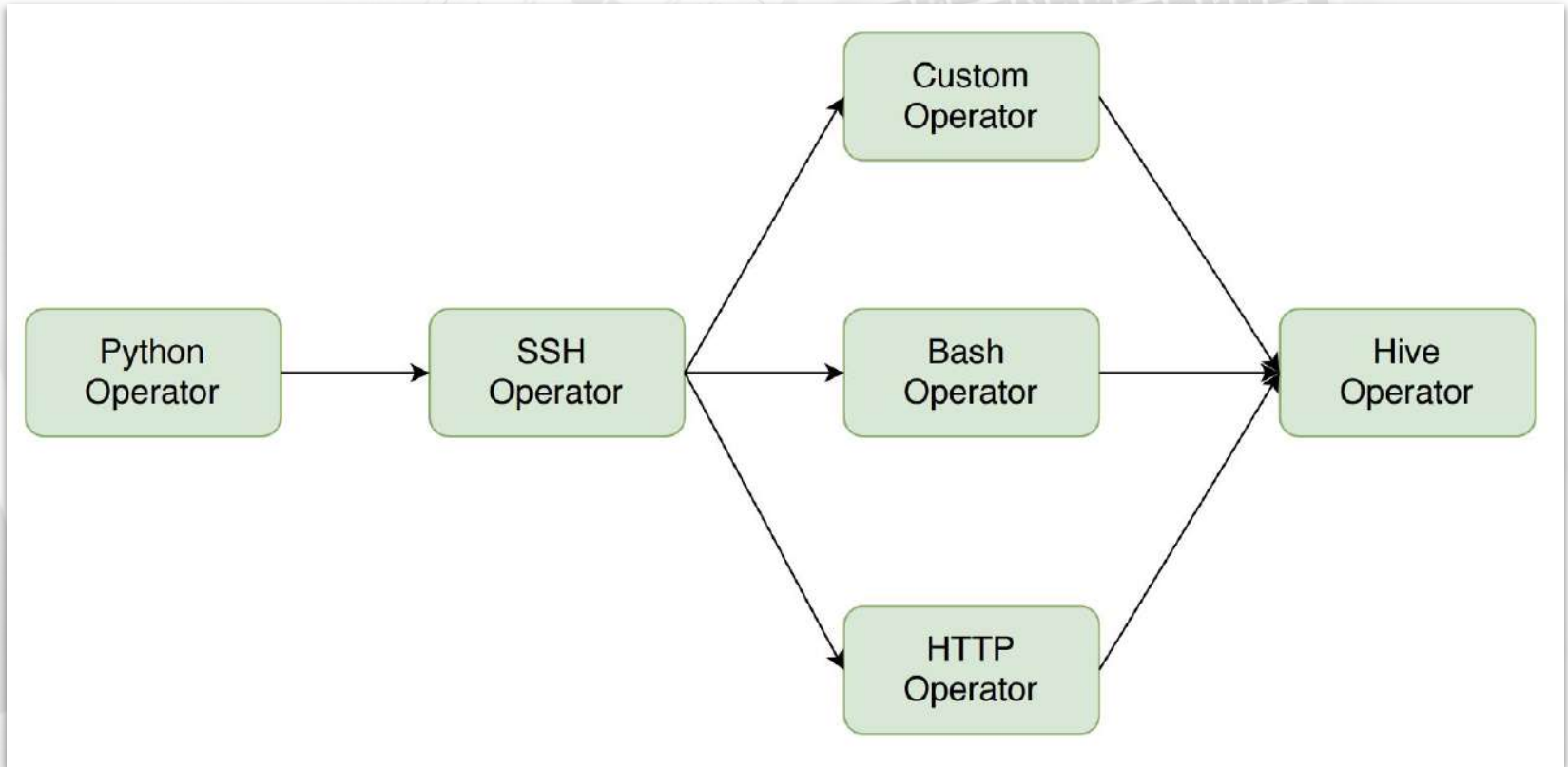


# Core concepts - Directed Acyclic Graph (DAG)



Source: [https://github.com/apache/incubator-airflow/blob/master/airflow/contrib/example\\_dags/example\\_twitter\\_README.md](https://github.com/apache/incubator-airflow/blob/master/airflow/contrib/example_dags/example_twitter_README.md)

# Core concepts - Operators



Source: <https://blog.usejournal.com/testing-in-airflow-part-1-dag-validation-tests-dag-definition-tests-and-unit-tests-2aa94970570c>



# Operator types

- Action Operators
  - Python, Bash, Docker, GCEInstanceStart, ...
- Sensor Operators
  - S3KeySensor, HivePartitionSensor, BigtableTableWaitForReplicationOperator , ...
- Transfer Operators
  - MsSqlToHiveTransfer, RedshiftToS3Transfer, ...



# Operator and Sensor

```
class ExampleOperator(BaseOperator):  
    def execute(self, context):  
        # Do something  
        pass
```



# Operator and Sensor

```
class ExampleOperator(BaseOperator):  
    def execute(self, context):  
        # Do something  
        pass  
  
class ExampleSensorOperator(BaseSensorOperator):  
    def poke(self, context):  
        # Check if the condition occurred  
        return True
```



# Operator good practices

- Idempotent
- Atomic
- No direct data sharing
  - Small portions of data between tasks: XCOMs
  - Large amounts of data: S3, GCS, etc.



# Core concepts - Tasks, TaskInstances, DagRuns

## UI DAG Tree View

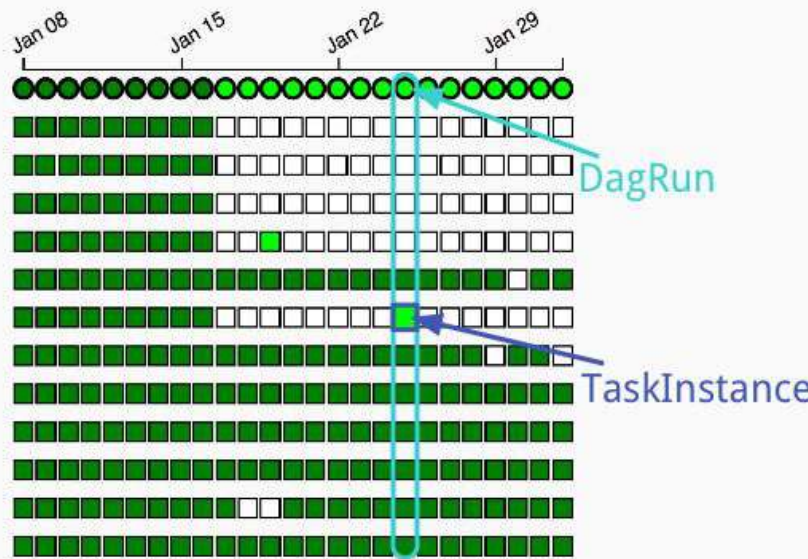
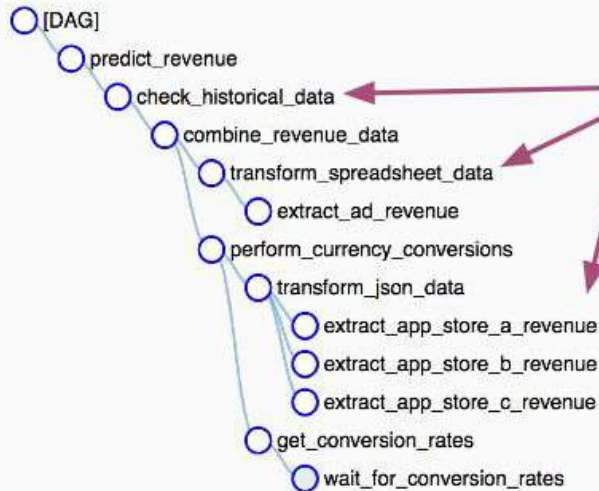
- CheckHistoricalDataOperator
- CombineDataRevenueDataOperator
- ConversionRatesSensor
- CurrencyConversionsOperator
- ExtractAdRevenueOperator
- ExtractAppStoreRevenueOperator
- ExtractConversionRatesOperator
- RevenuePredictionOperator
- TransformJSONDataOperator
- TransformSpreadsheetDataOperator

- success
- running
- failed
- skipped
- retry
- queued
- no status

Operators

Sensor

Tasks



DagRun

TaskInstance

Source: <https://medium.com/@dustinstansbury/understanding-apache-airflows-key-concepts-a96fed52b1a>



**Show me the code!**





Google Cloud Platform



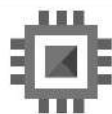
SQL



Bigtable



Spanner

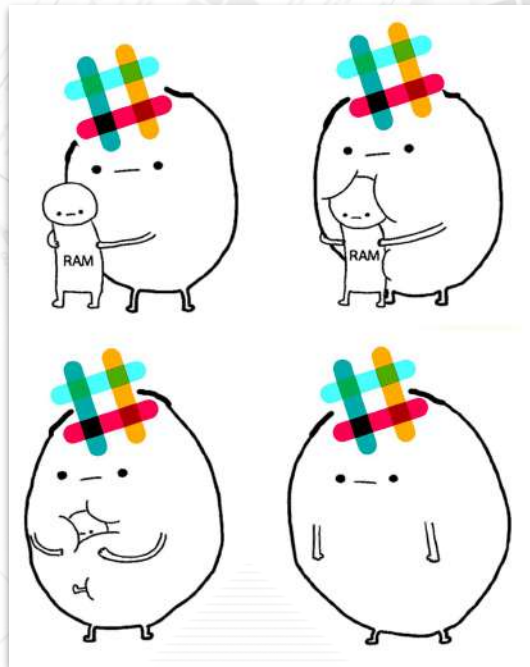


Compute Engine





# The solution



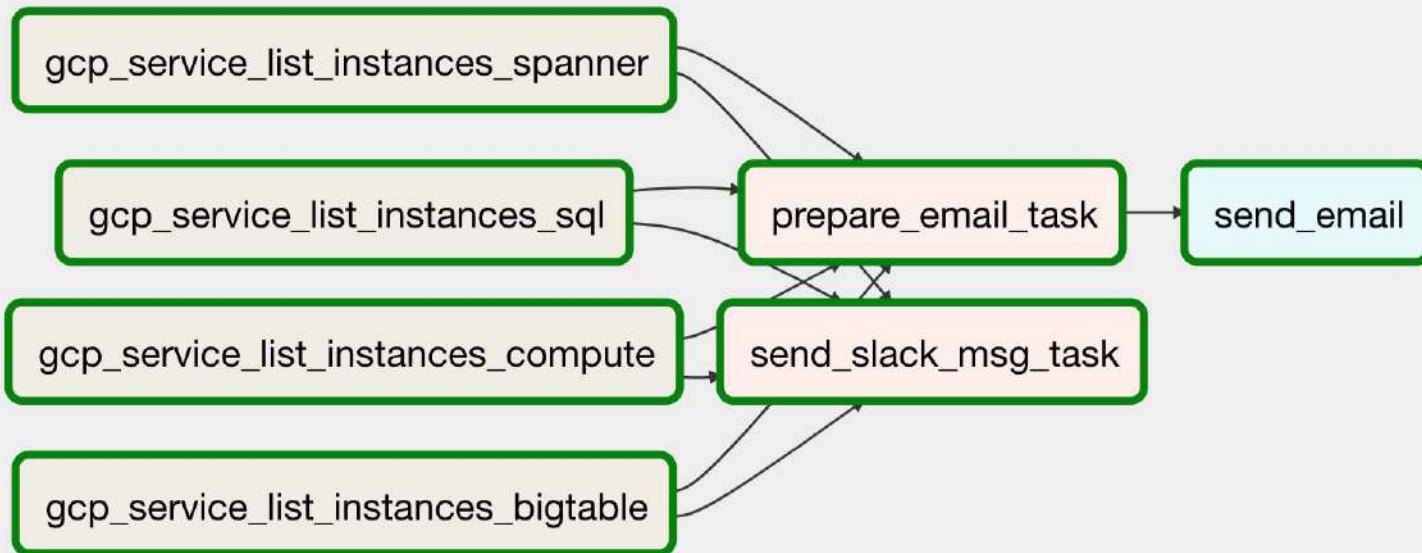


# Solution components

- Generic
  - BashOperator
  - PythonOperator
- Specific
  - EmailOperator



# The DAG





# Initialize DAG

```
dag = DAG(dag_id='gcp_spy',  
          ...  
          )
```



# Initialize DAG

```
dag = DAG(dag_id='gcp_spy',  
          default_args={  
            'start_date': utils.dates.days_ago(1),  
            'retries': 1  
          },  
          ...  
)
```



# Initialize DAG

```
dag = DAG(dag_id='gcp_spy',  
          default_args={  
              'start_date': utils.dates.days_ago(1),  
              'retries': 1  
          },  
          schedule_interval='0 16 * * *'  
)
```





# List of instances

```
bash_task = BashOperator(  
    task_id="gcp_service_list_instances_sql",  
    ...  
)
```



# List of instances

```
bash_task = BashOperator(  
    task_id="gcp_service_list_instances_sql",  
    bash_command=  
        "gcloud sql instances list | tail -n +2 | grep -oE '^[^ ]+' "  
        "| tr '\n' ' '",  
    ...  
)
```



# List of instances

```
bash_task = BashOperator(  
    task_id="gcp_service_list_instances_sql",  
    bash_command=  
        "gcloud sql instances list | tail -n +2 | grep -oE '^[^ ]+' "  
        "| tr '\n' ' '",  
    xcom_push=True,  
    ...  
)
```



# List of instances

```
bash_task = BashOperator(  
    task_id="gcp_service_list_instances_sql",  
    bash_command=  
        "gcloud sql instances list | tail -n +2 | grep -oE '^[^ ]+' "  
        "| tr '\n' ' '",  
    xcom_push=True,  
    dag=dag  
)
```



# All services

```
GCP_SERVICES = [  
    ('sql', 'Cloud SQL'),  
    ('spanner', 'Spanner'),  
    ('bigtable', 'BigTable'),  
    ('compute', 'Compute Engine'),  
]
```



# List of instances - all services

????

```
bash_task = BashOperator(  
    task_id="gcp_service_list_instances_sql",  
    bash_command=  
        "gcloud sql instances list | tail -n +2 | grep -oE '^[^ ]+' "  
        "| tr '\n' ' '",  
    xcom_push=True,  
    dag=dag  
)
```



# List of instances - all services

```
for gcp_service in GCP_SERVICES:
    bash_task = BashOperator(
        task_id="gcp_service_list_instances_{}".format(gcp_service[0]),
        bash_command=
            "gcloud {} instances list | tail -n +2 | grep -oE '^[^ ]+' "
            "| tr '\n' ' '".format(gcp_service[0]),
        xcom_push=True,
        dag=dag
    )
```

# Send Slack message

```
send_slack_msg_task = PythonOperator(  
    python_callable=send_slack_msg,  
    provide_context=True,  
    task_id='send_slack_msg_task',  
    dag=dag  
)
```





# Send Slack message

```
send_slack_msg_task = PythonOperator(  
    python_callable=send_slack_msg,  
    provide_context=True,  
    task_id='send_slack_msg_task',  
    dag=dag  
)
```

```
def send_slack_msg(**context):  
    for gcp_service in GCP_SERVICES:  
        result = context['task_instance'].\  
            xcom_pull(task_ids='gcp_service_list_instances_{}'  
                    .format(gcp_service[0]))  
  
    ...
```

```
def send_slack_msg(**context):  
    for gcp_service in GCP_SERVICES:  
        result = context['task_instance'].\  
            xcom_pull(task_ids='gcp_service_list_instances_{}'  
                    .format(gcp_service[0]))  
  
    ...
```

```
def send_slack_msg(**context):  
    for gcp_service in GCP_SERVICES:  
        result = context['task_instance'].\  
            xcom_pull(task_ids='gcp_service_list_instances_{}'  
                    .format(gcp_service[0]))  
  
    data = ...  
    ...
```

```
def send_slack_msg(**context):
    for gcp_service in GCP_SERVICES:
        result = context['task_instance'].\
            xcom_pull(task_ids='gcp_service_list_instances_{}'
                      .format(gcp_service[0]))

    data = ...
    requests.post(
        url=SLACK_WEBHOOK,
        data=json.dumps(data),
        headers={'Content-type': 'application/json'})
)
```

# Prepare email

```
prepare_email_task = PythonOperator(
    python_callable=prepare_email,
    provide_context=True,
    task_id='prepare_email_task',
    dag=dag
)
```

# Prepare email

```
prepare_email_task = PythonOperator(  
    python_callable=prepare_email,  
    provide_context=True,  
    task_id='prepare_email_task',  
    dag=dag  
)
```

```
def prepare_email(**context):
    for gcp_service in GCP_SERVICES:
        result = context['task_instance'].xcom_pull(
            task_ids='gcp_service_list_instances_{}'.format(
                gcp_service[0]))
    ...
    html_content = ...
    context['task_instance'].xcom_push(key='email', value=html_content)
```



```
def prepare_email(**context):  
    for gcp_service in GCP_SERVICES:  
        result = context['task_instance'].\  
            xcom_pull(task_ids='gcp_service_list_instances_{}'  
                    .format(gcp_service[0]))  
  
        ...  
        html_content = ...  
        context['task_instance'].xcom_push(key='email', value=html_content)
```



# Send email

```
send_email_task = EmailOperator(  
    task_id='send_email',  
    to='szymon.przedwojski@polidea.com',  
    subject=INSTANCES_IN_PROJECT_TITLE,  
    html_content=...,  
    dag=dag  
)
```



# Send email

```
send_email_task = EmailOperator(  
    task_id='send_email',  
    to='szymon.przedwojski@polidea.com',  
    subject=INSTANCES_IN_PROJECT_TITLE,  
    html_content=  
        "{{ task_instance.xcom_pull(task_ids='prepare_email_task', key='email') }}",  
    dag=dag  
)
```

# Dependencies

```

for gcp_service in GCP_SERVICES:
    bash_task = BashOperator(
        ...
    )
    bash_task >> send_slack_msg_task
    bash_task >> prepare_email_task

```



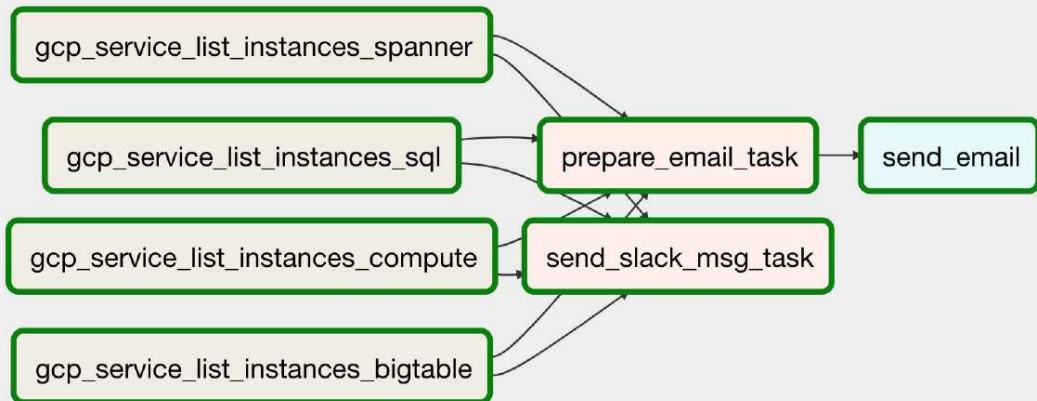
# Dependencies

```
for gcp_service in GCP_SERVICES:  
    bash_task = BashOperator(  
        ...  
    )  
    bash_task >> send_slack_msg_task  
    bash_task >> prepare_email_task  
  
prepare_email_task >> send_email_task
```



# Dependencies

```
for gcp_service in GCP_SERVICES:  
    bash_task = BashOperator(  
        ...  
    )  
    bash_task >> send_slack_msg_task  
    bash_task >> prepare_email_task  
  
prepare_email_task >> send_email_task
```





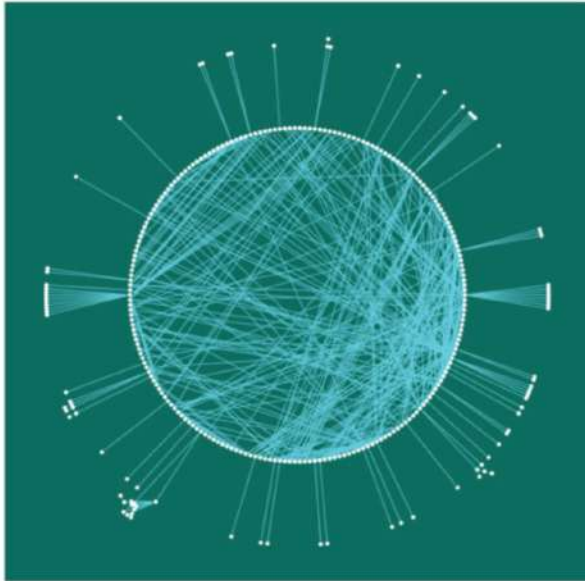
# Demo

<https://github.com/PolideaInternal/airflow-gcp-spy>



# Complex DAGs

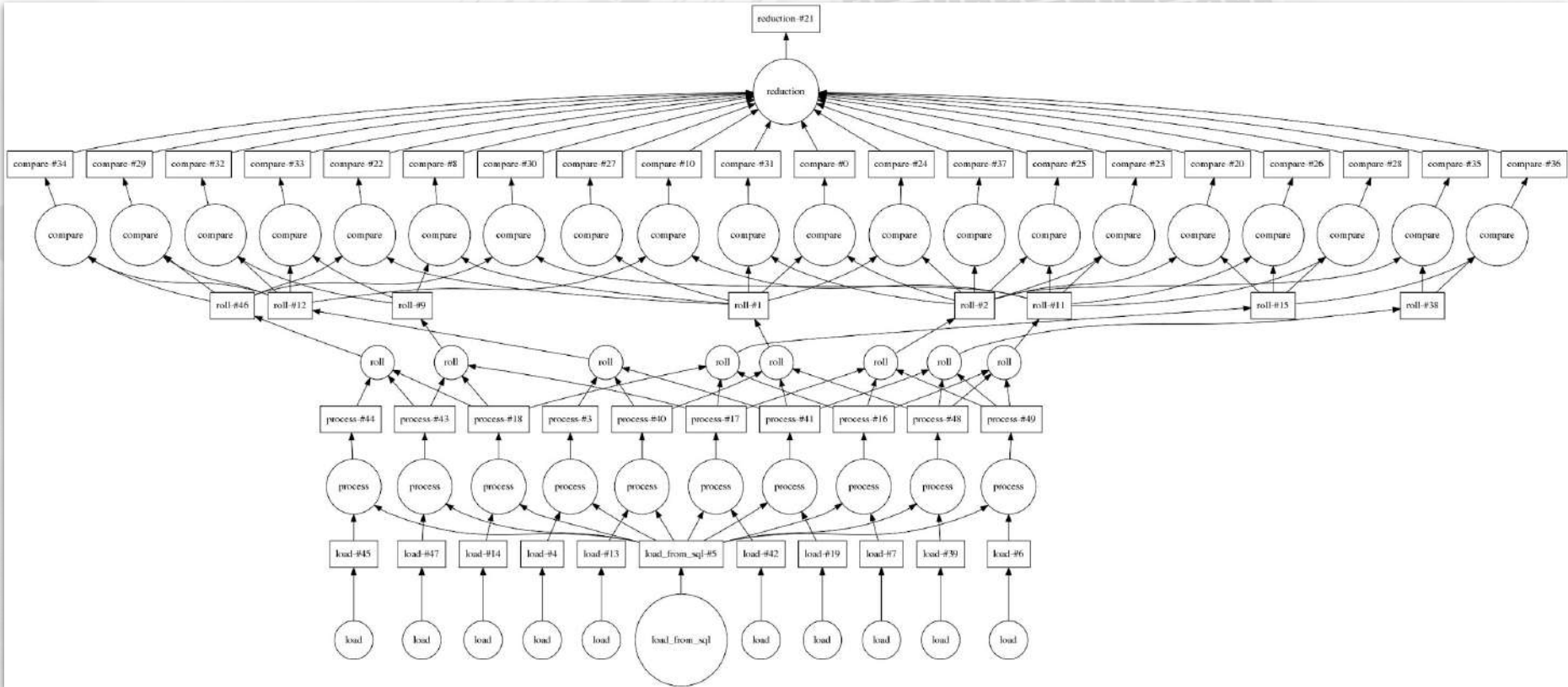
## Real Pipelines of Clover





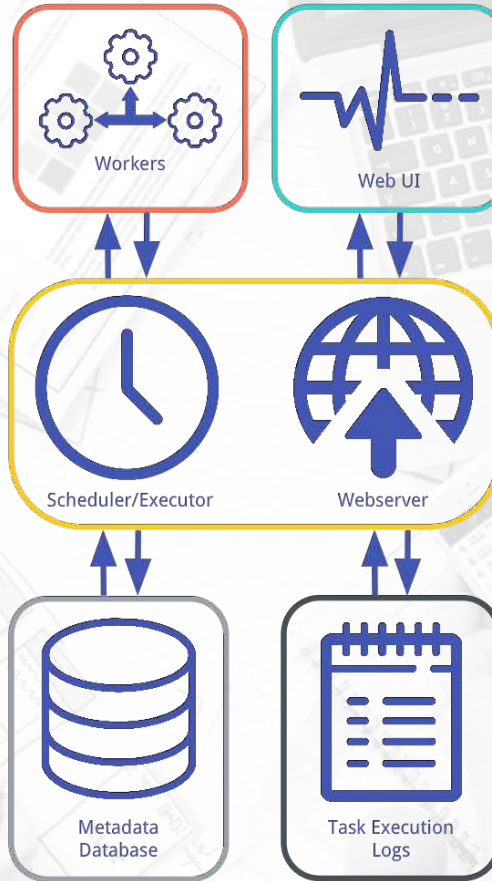


# Complex, Manageable, DAGs



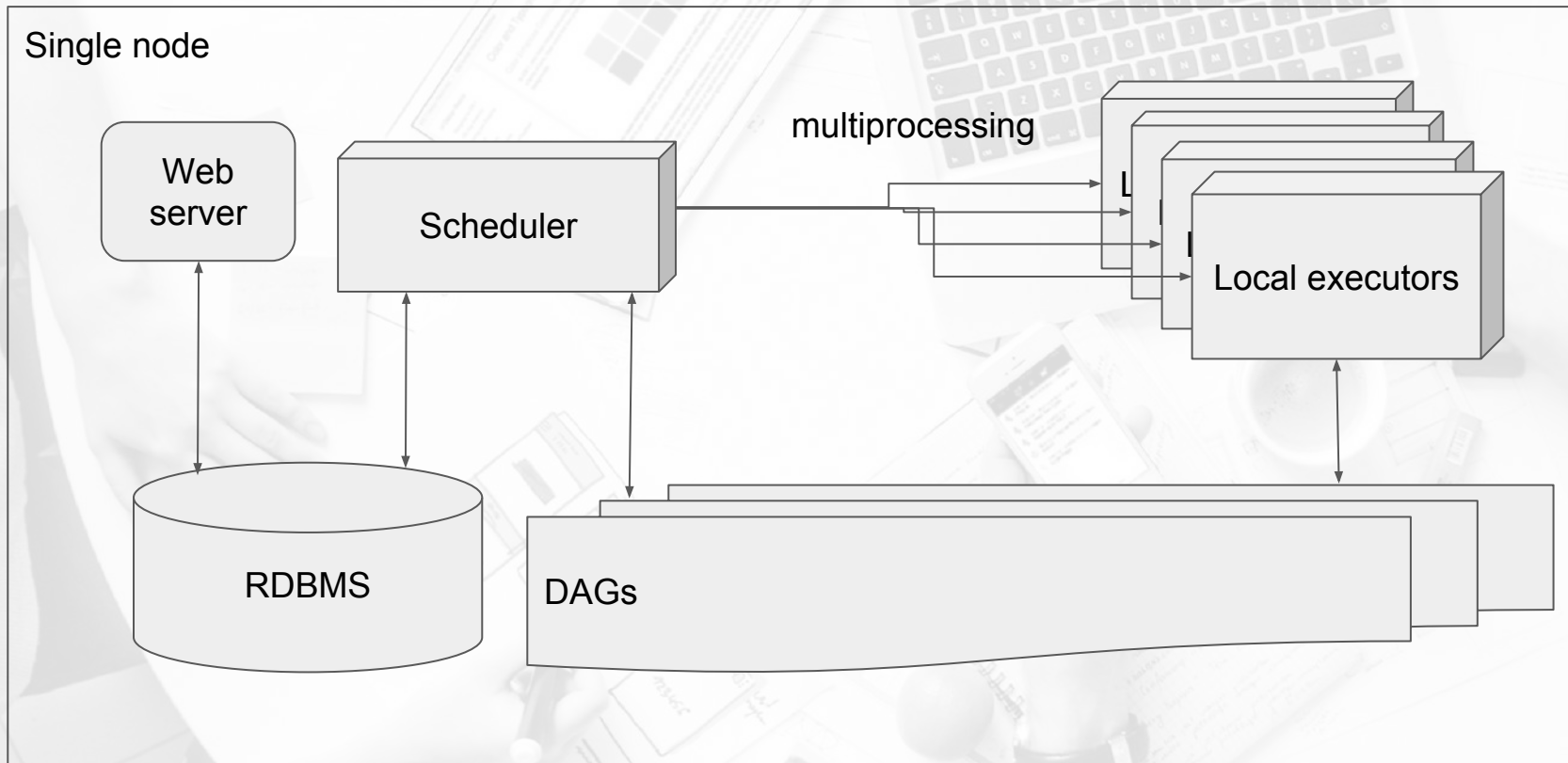


## Airflow's General Architecture



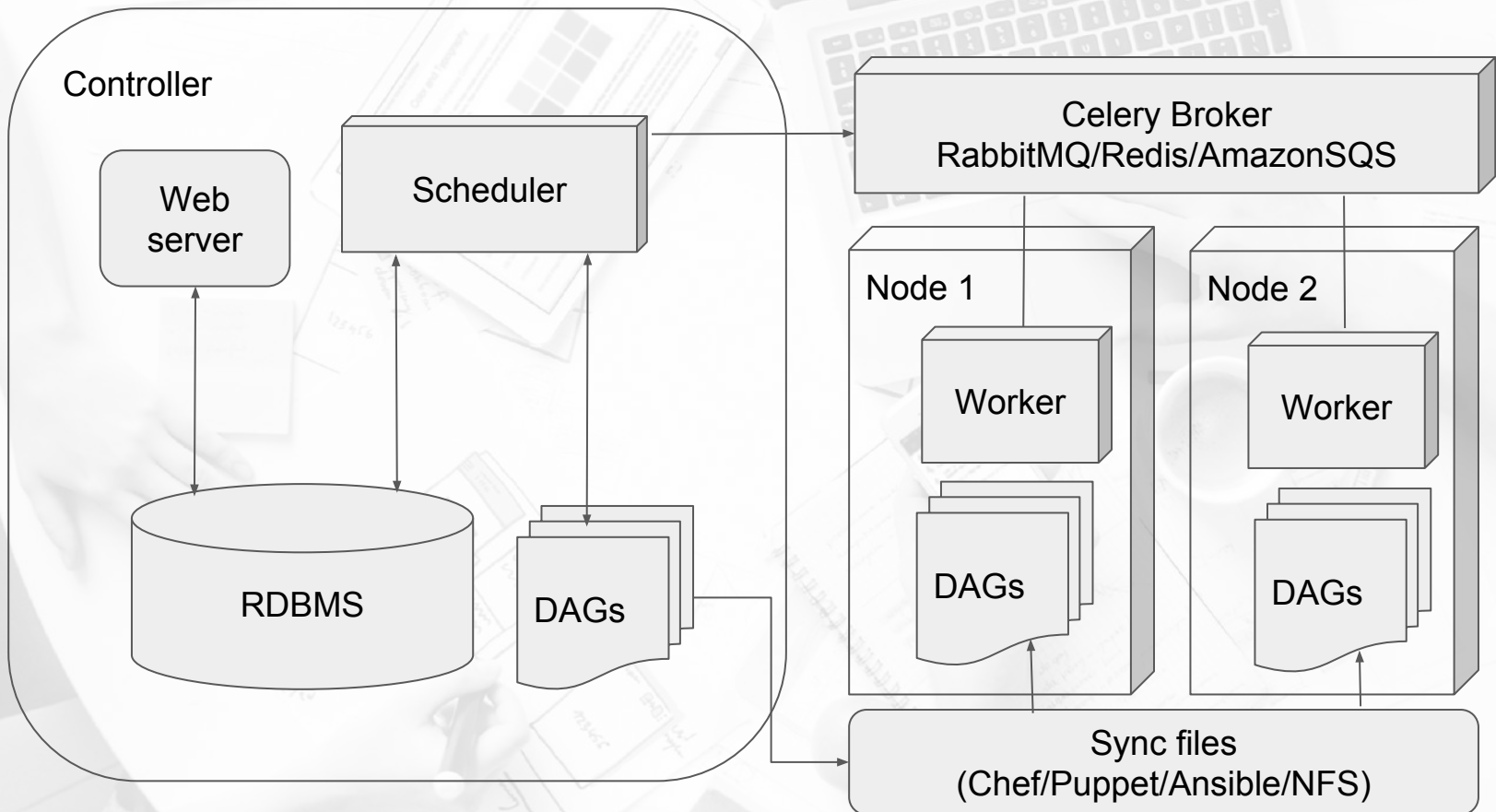


# Local Executor



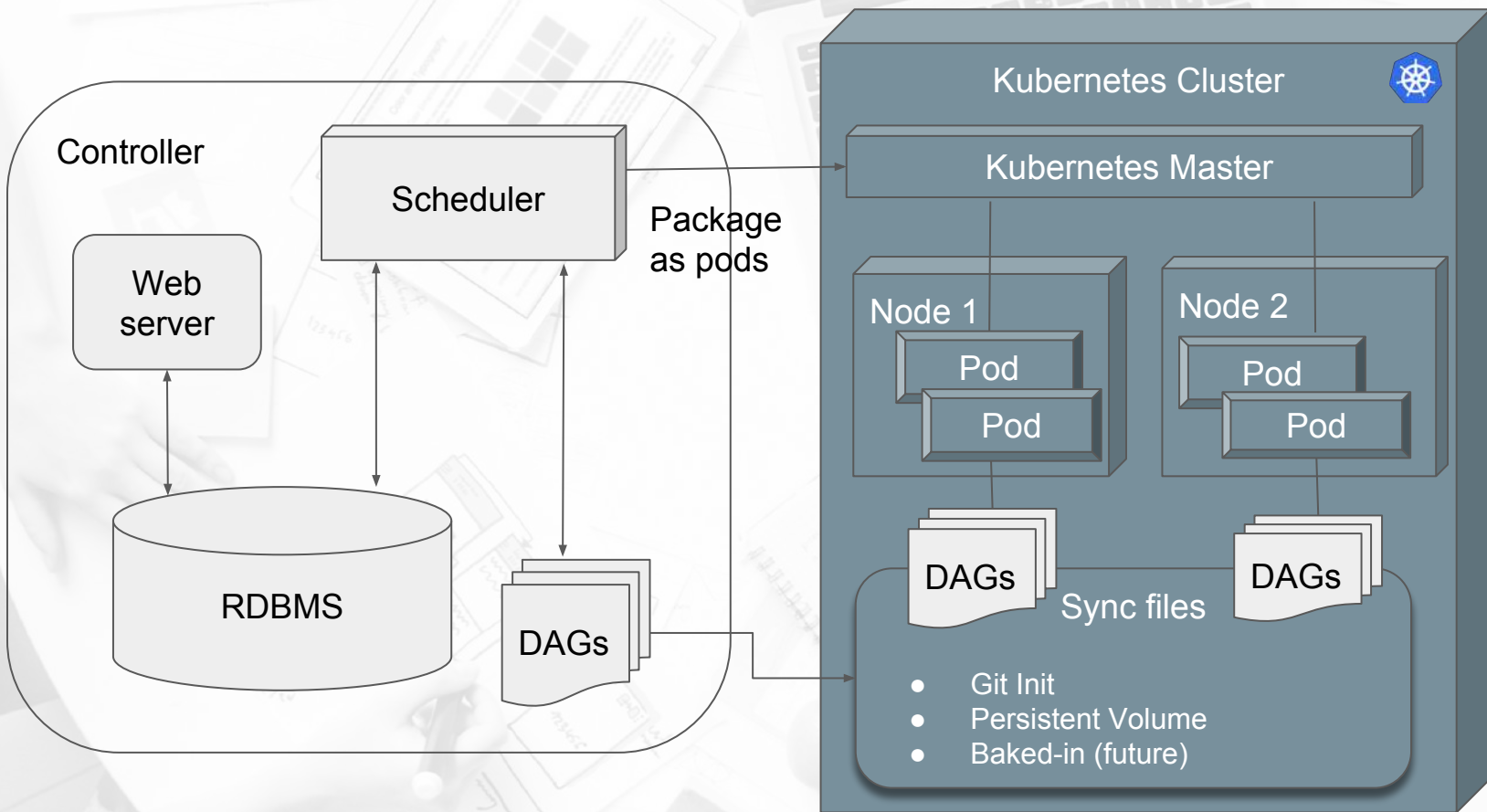


# Celery Executor





# (Beta): Kubernetes Executor





# GCP - Composer

Google Cloud Platform polidea-airflow

Composer Environment details REFRESH DELETE

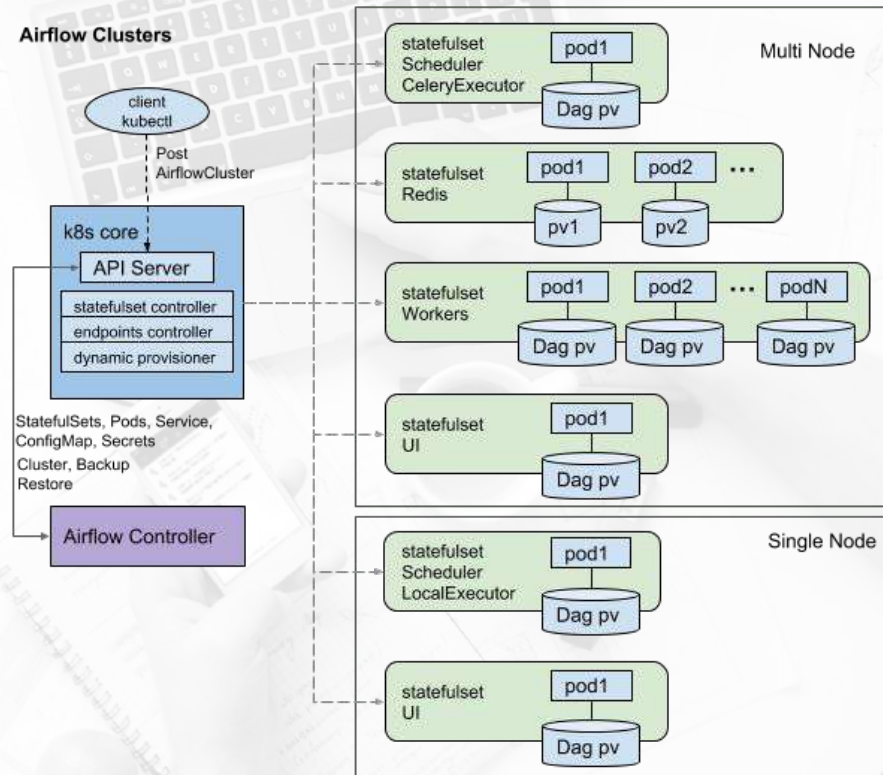
**gcp-spy**  
This environment is running.

ENVIRONMENT CONFIGURATION AIRFLOW CONFIGURATION OVERRIDES ENVIRONMENT VARIABLES

EDIT

Name	gcp-spy
Zone	europe-west1-d
Service account	198907790164-compute@developer.gserviceaccount.com
Google API scopes	https://www.googleapis.com/auth/cloud-platform
GKE cluster ID	projects/polidea-airflow/zones/europe-west1-d/clusters/europe-west1-gcp-spy-198cc3c1-gke
Image version	composer-1.3.0-airflow-1.9.0
Network tags	None
Worker nodes	
Disk size (GB)	20
Machine type	n1-standard-1
Node count	3
Network ID	default
Subnetwork ID	-
DAGs folder	gs://europe-west1-gcp-spy-198cc3c1-bucket/dags
Airflow web UI	https://o75ce47b5c088dc5a-tp.appspot.com
Stackdriver	view logs
Created	Mon Nov 19 2018 10:22:20 GMT+0100 (Central European Standard Time)
Updated	Mon Nov 19 2018 10:38:22 GMT+0100 (Central European Standard Time)

## Airflow Clusters



<https://github.com/GoogleCloudPlatform/airflow-operator>



**Thank You!**

Follow us @  
[polidea.com/blog](https://polidea.com/blog)





**We're hiring!**





Questions? :)

Follow us @  
[polidea.com/blog](https://polidea.com/blog)

