

1.

```
nifi@a848917905fd:~/ingest$ ls
ingest.sh
nifi@a848917905fd:~/ingest$ cat ingest.sh
#!/bin/bash

# Directorio de destino
DEST_DIR="/home/nifi/ingest"

# Crear el directorio si no existe
if [ ! -d "$DEST_DIR" ]; then
    echo "Creando el directorio $DEST_DIR"
    mkdir -p "$DEST_DIR"
fi

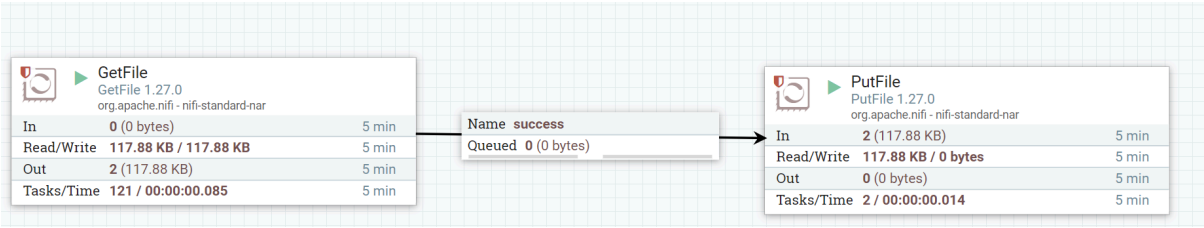
# Descargar el archivo titanic.csv
FILE_URL="https://dataengineerpublic.blob.core.windows.net/data-engineer/titanic.csv"
DEST_FILE="$DEST_DIR/titanic.csv"

echo "Descargando el archivo titanic.csv a $DEST_FILE"
curl -o "$DEST_FILE" "$FILE_URL"

# Verificar si la descarga fue exitosa
if [ $? -eq 0 ]; then
    echo "Descarga completada exitosamente."
else
    echo "Error al descargar el archivo."
    exit 1
fi
nifi@a848917905fd:~/ingest$
```

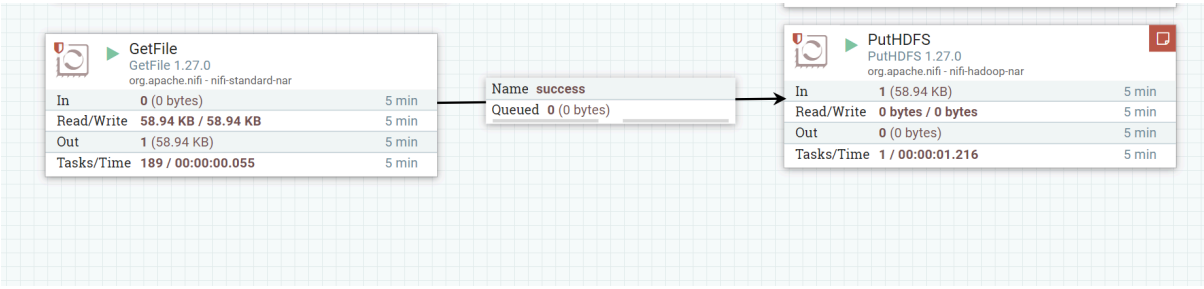
```
nifi@a848917905fd:~/ingest$ ./ingest.sh
Descargando el archivo titanic.csv a /home/nifi/ingest/titanic.csv
  % Total    % Received % Xferd  Average Speed   Time    Time     Time  Current
                                 Dload  Upload   Total   Spent    Left   Speed
100 60353  100 60353    0     0  67225      0  --:--:-- --:--:-- --:--:-- 67208
Descarga completada exitosamente.
nifi@a848917905fd:~/ingest$ ls
ingest.sh titanic.csv
nifi@a848917905fd:~/ingest$
```

2, 3 y 4



```
nifi@a848917905fd:~$ cd bucket/  
nifi@a848917905fd:~/bucket$ ls  
titanic.csv  
nifi@a848917905fd:~/bucket$
```

5 y 6



```
hadoop@f769ed737105:/$ hdfs dfs -ls /nifi  
Found 1 items  
-rw-r--r--  1 nifi supergroup      60353 2024-09-15 23:10 /nifi/titanic.csv  
hadoop@f769ed737105:/$
```

```

hadoop@f769ed737105:~/scripts$ cat trans_titanic.py
from pyspark.context import SparkContext
from pyspark.sql.session import SparkSession
from pyspark.sql.functions import col, avg, when
sc = SparkContext('local')
spark = SparkSession(sc)
from pyspark.sql import HiveContext
hc = HiveContext(sc)

titanic = spark.read.csv('hdfs://172.17.0.2:9000/nifi/titanic.csv', header=True, inferSchema=True)

titanic = titanic.drop('SibSp', 'Parch')

avg_ages = titanic.groupBy('Sex').agg(avg('Age').alias('avg_age')).collect()
avg_male_age = [row['avg_age'] for row in avg_ages if row['Sex'] == 'male'][0]
avg_female_age = [row['avg_age'] for row in avg_ages if row['Sex'] == 'female'][0]

titanic = titanic.withColumn('Age', when((col('Sex') == 'male') & col('Age').isNull(), avg_male_age)
                                   .when((col('Sex') == 'female') & col('Age').isNull(), avg_female_age)
                                   .otherwise(col('Age'))))

titanic = titanic.fillna({'Cabin': 0})

titanic.write.insertInto('titanic_data.titanic')

spark.stop()

```

```

from datetime import timedelta
from airflow import DAG
from airflow.operators.bash_operator import BashOperator
from airflow.operators.dummy import DummyOperator
from airflow.utils.dates import days_ago

args = {
    'owner': 'airflow',
}

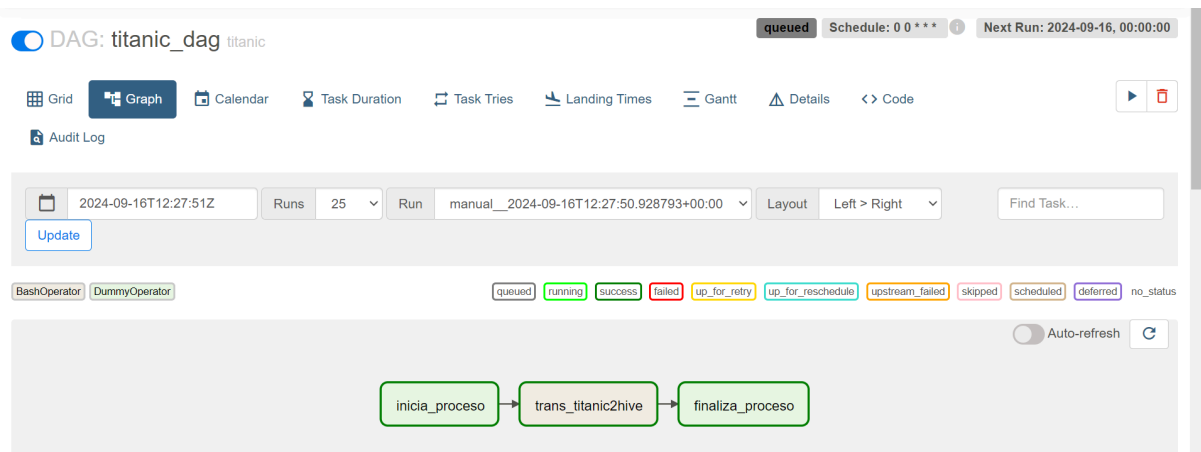
with DAG(
    dag_id='titanic_dag',
    default_args=args,
    description='titanic',
    schedule_interval='0 0 * * *',
    start_date=days_ago(2),
    dagrun_timeout=timedelta(minutes=60),
    tags=['ingest', 'transform'],
    params={"example_key": "example_value"},
) as dag:

    inicia_proceso = DummyOperator(
        task_id='inicia_proceso',
    )

    trans_titanic2hive = BashOperator(
        task_id='trans_titanic2hive',
        bash_command='ssh hadoop@172.17.0.2 /home/hadoop/spark/bin/spark-submit --files /home/hadoop/hive/conf/hive-site.xml /home/hadoop/scripts/trans_titanic.py ',
    )

    finaliza_proceso = DummyOperator(
        task_id='finaliza_proceso',
    )

```



```
hive> select * from titanic limit 5;
OK
1      0      3      Braund  Mr. Owen Harris      male  22      NULL  7.25      0
2      1      1      Cumings  Mrs. John Bradley (Florence Briggs Thayer) female 3
8      NULL  71.2833 C85
3      1      3      Heikinen  Miss. Laina      female 26      NULL  7.925     0
4      1      1      Futrelle  Mrs. Jacques Heath (Lily May Peel) female 3
5      113803.0      53.1      C123
5      0      3      Allen  Mr. William Henry      male  35      373450.0  8
.05      0
Time taken: 0.19 seconds, Fetched: 5 row(s)
```

8

Enter a part of object

- northwind - localhos
- titanic\_data - localhos
  - default
  - f1
  - northwind\_analyti
  - titanic\_data
    - Tables
      - titanic
        - Columns
          - 123 passe
          - 123 surviv
          - 123 pclass
          - A-Z name
          - A-Z sex (S
          - A-Z age (S
          - A-Z ticket
          - 123 fare (I

```
--d
select sex, COUNT(passengerid) as sobrevivientes
from titanic t
WHERE survived =1
group by sex;

--b
select pclass , COUNT(passengerid) as sobrevivientes
from titanic t
WHERE survived =1
group by pclass ;

--c
select name, cast(age as int) as age
from titanic t
WHERE survived =1
```

Results 1 x

select sex, COUNT(passengerid) as sobrevivientes

	A-Z sex	123 sobrevivientes
1	female	233
2	male	109

Project - General

Enter a part of object

- northwind - localhos
- titanic\_data - localhos
  - default
  - f1
  - northwind\_analyti
  - titanic\_data
    - Tables
      - titanic
        - Columns
          - 123 passe
          - 123 surviv
          - 123 pclass
          - A-Z name
          - A-Z sex (S
          - A-Z age (S
          - A-Z ticket
          - 123 fare (I

Project - General

Name DataSou...

Results 1 x

select pclass, COUNT(passengerid) as sobrevivientes

from titanic t

WHERE survived =1

group by sex;

--b

select pclass , COUNT(passengerid) as sobrevivientes

from titanic t

WHERE survived =1

group by pclass ;

--c

select name, cast(age as int) as age

from titanic t

WHERE survived =1

Results 1 x

select pclass, COU

Enter a SQL expression to filter re

	123 pclass	123 sobrevivientes
1	1	136
2	2	87
3	3	119

Grid

Text

from titanic t

WHERE survived =1

group by pclass ;

--c

select name, cast(age as int) as age

from titanic t

WHERE survived =1

group by name, age

order by age desc

limit 1;

-- d

select name, cast(age as int) as age

from titanic t

Results 1 x

select name, cast(a

Enter a SQL expression to filter re

	A-Z name	123 age
1	Barkworth, Mr. Algernon Henry Wilson	80

Grid

Text

Project - General

- northwind - localhos
- titanic\_data - localhos
  - default
  - f1
  - northwind\_analyti
  - titanic\_data
    - Tables
      - titanic
        - Columns
          - 123 passe
          - 123 surviv
          - 123 pclass
          - A-Z name
          - A-Z sex (S
          - A-Z age (S
          - A-Z ticket
          - 123 fare (I

```
WHERE survived =1
group by name, age
order by age desc
limit 1;

-- d

select name, cast(age as int) as age
from titanic t
WHERE survived =1
group by name, age
order by age asc
limit 1;
```

Results 1 x

select name, cast(a

Enter a SQL expression to filter re

Grid	A-Z name	123 age
1	Allison, Master. Hudson Trevor	0