

Data pipelines for PySpark

Preprocessing data for data modeling

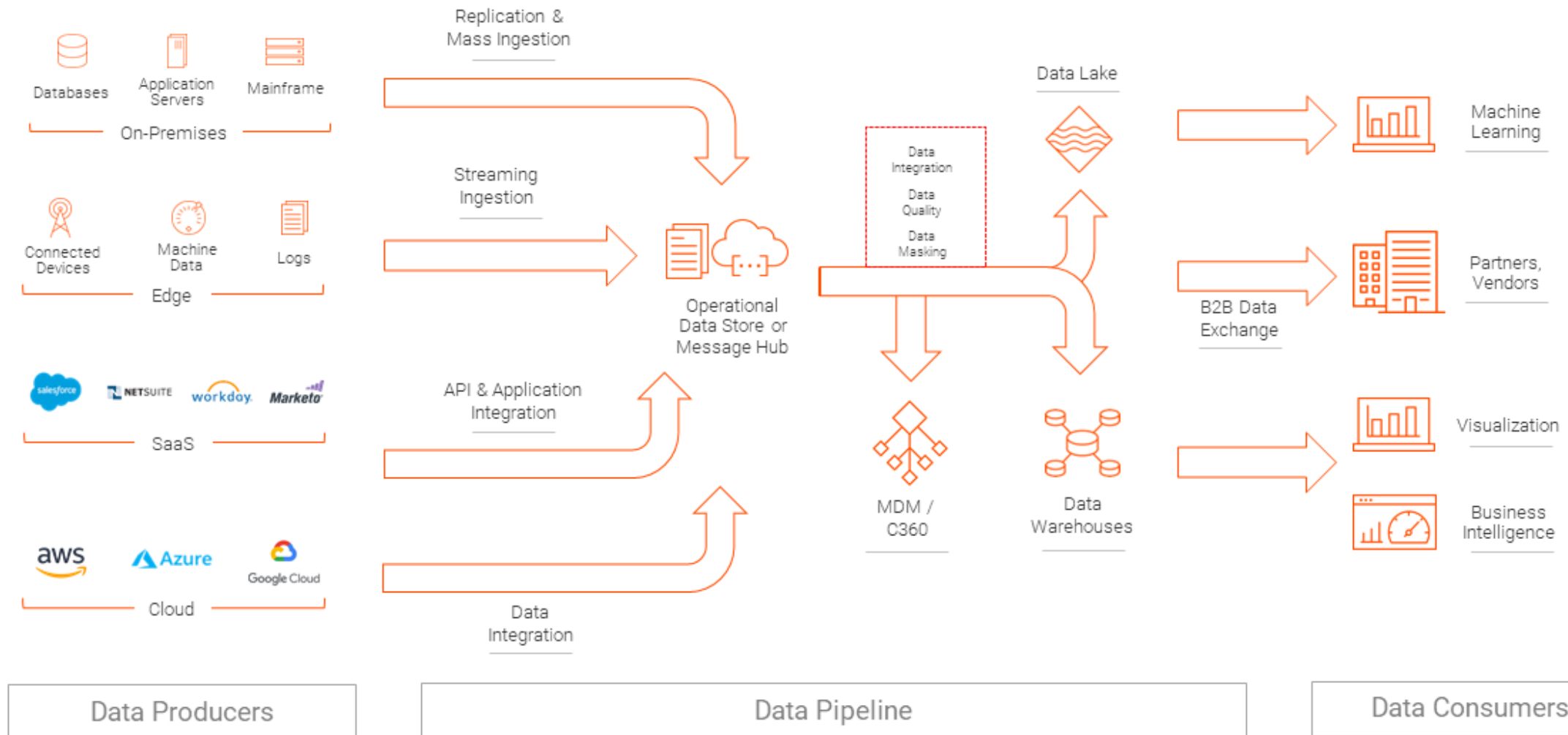
Data pipeline summary

- A data pipeline processes elements connected in series, where the output of one element is the input of the next one
- Common to
 - Ingest data using distributed data logging programs (Kafka)
 - Process data using PySpark
 - Data cleaning and eventually ML pipeline
 - Manage the ingest to modeling via a workflow framework
 - E.g., Apache Airflow
- Notebook workflow is a common for PySpark pipelines

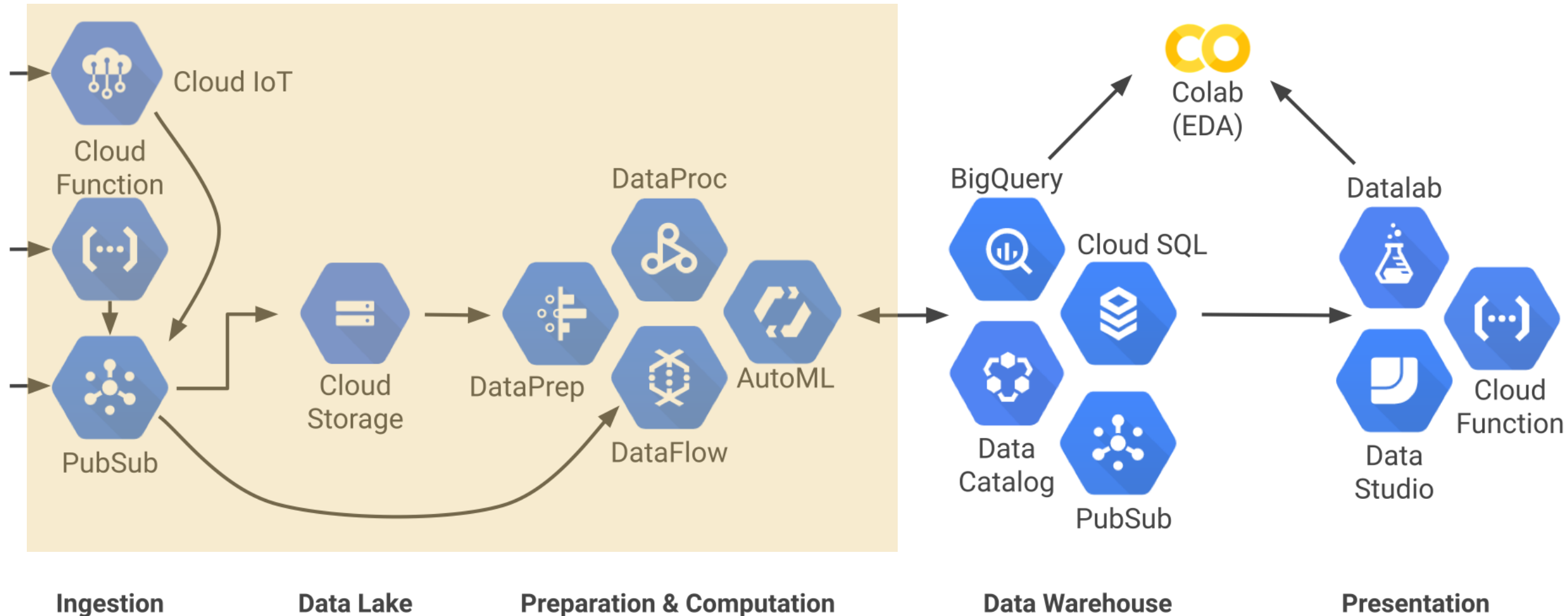
Data pipeline

- is a set of data processing elements connected in series, where the output of one element is the input of the next one.
 - The elements of a pipeline are often executed in parallel or in time-sliced fashion
- Examples
 - Machine learning pipeline
 - Data is prepared for input into ML library functions
 - Data pipeline
 - Data is acquired and processed in preparation for input into another process, such as an ML pipeline

Data pipeline patterns

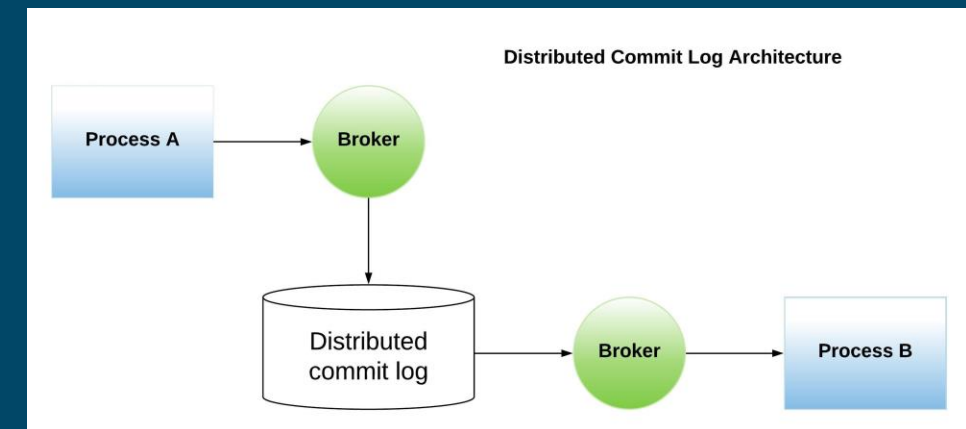
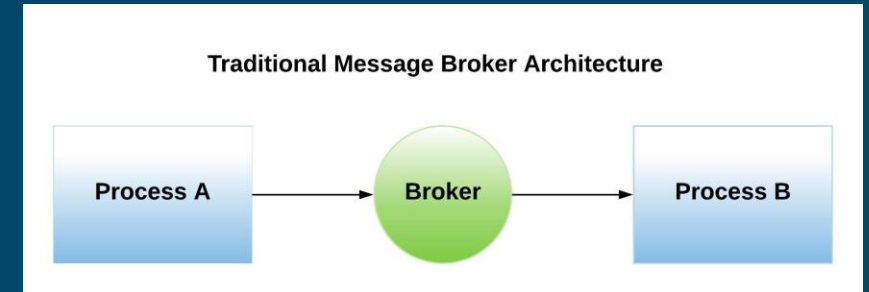


Data pipeline in Google Cloud Platform



Step 1: Ingest the data

- Traditional message broker systems (JMS) connect direct to brokers, and brokers which connect direct to processes.
 - These solutions tend to be optimized towards flexibility and configurable delivery guarantees.
- Distributed commit log technologies are similar in role to the traditional broker message systems.
 - They are optimized towards concentrated on scalability and throughput, which lesser guarantees on event ordering or arrival
 - Apache Kafka
 - Amazon Kinesis
 - Microsoft Azure Event Hubs
 - Google Pub/Sub

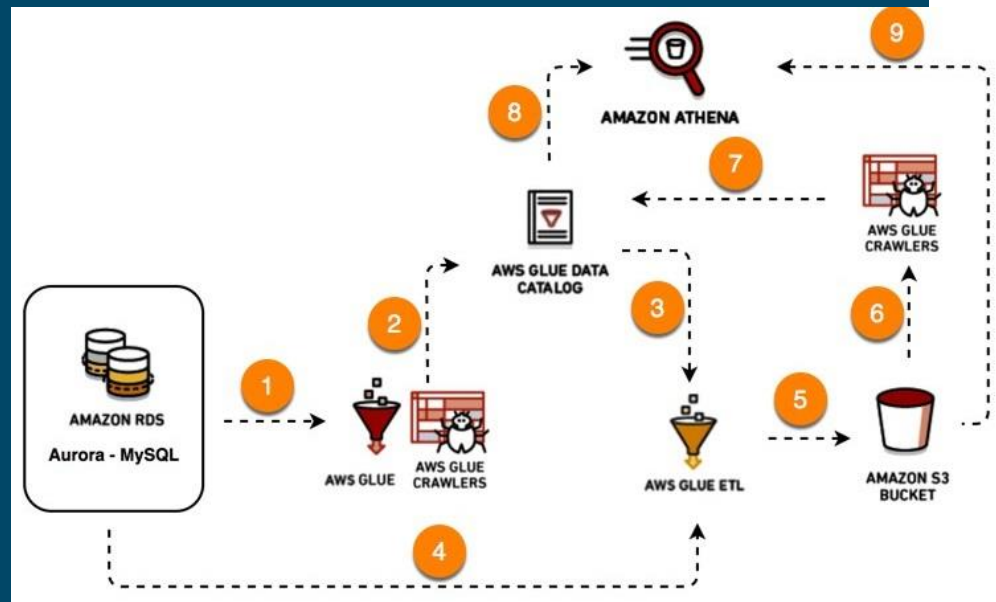
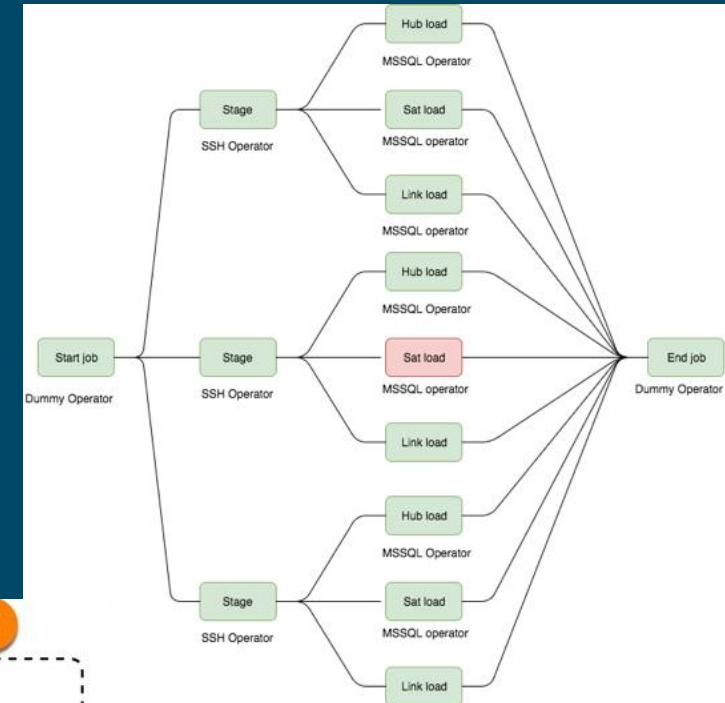


Step 2: process the data

- Transformation using
 - Python
 - Not for big data (parallelism)
 - PySpark
 - Spark notebook workflow
 - Apache Beam
 - Google Dataflow
 - AWS Data Pipeline

Step 3: schedule the ingest, data process, etc.

- Apache airflow
 - Open source, run anywhere, task scheduler
- AWS Glue
 - AWS serverless task scheduler



NEXT@:

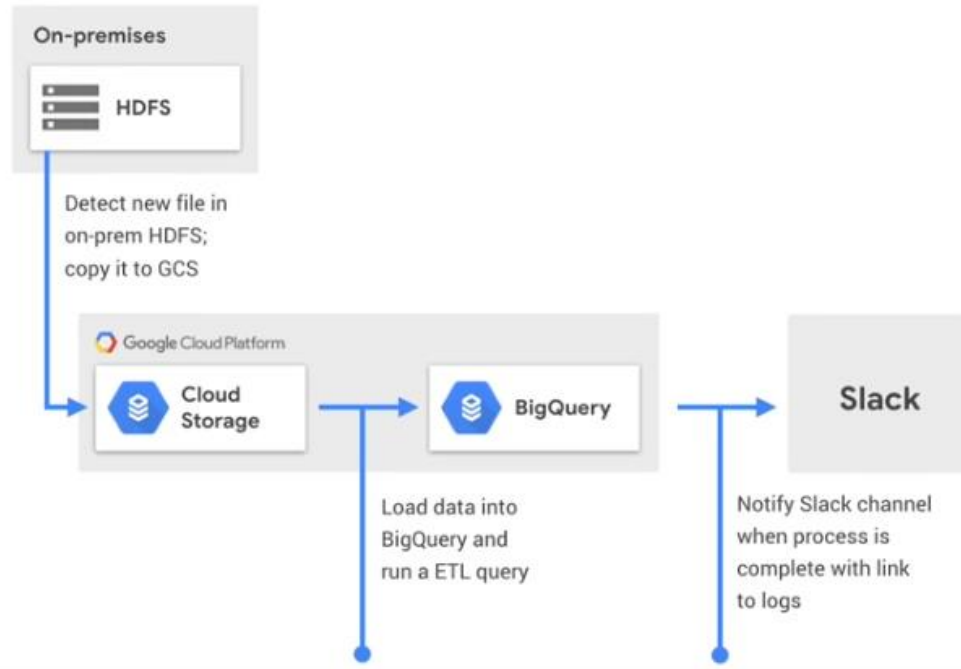
REVIEW

Flexible, Easy Data
Pipelines on Google Cloud
with Cloud Composer



Illustrative summary with GCP Cloud Composer

Example Workflow



0:41 / 2:58

CC Settings Full Screen

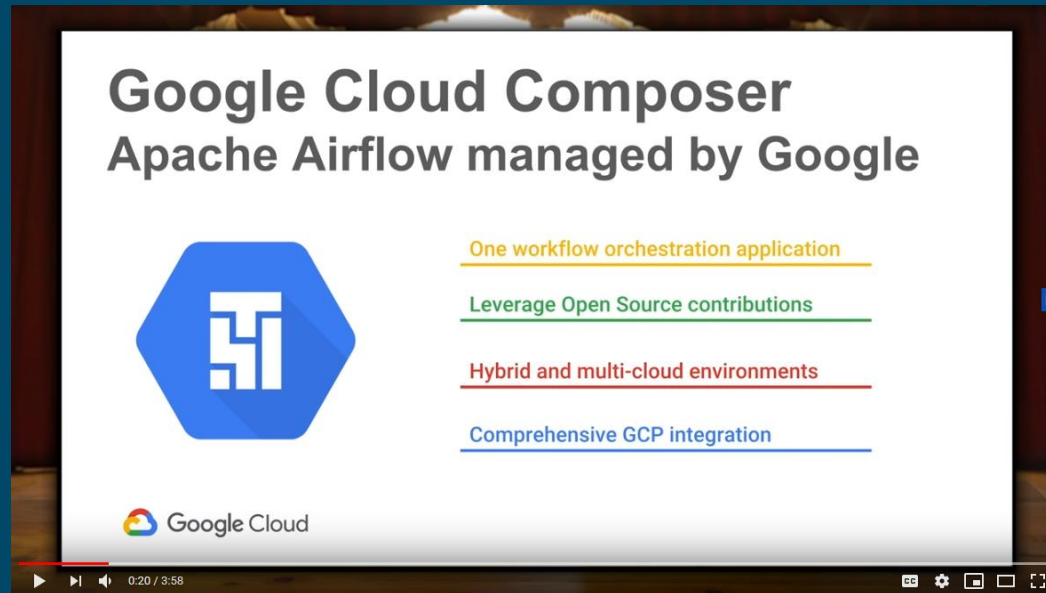
TAKE5

Cloud Composer



Example of running Cloud Composer

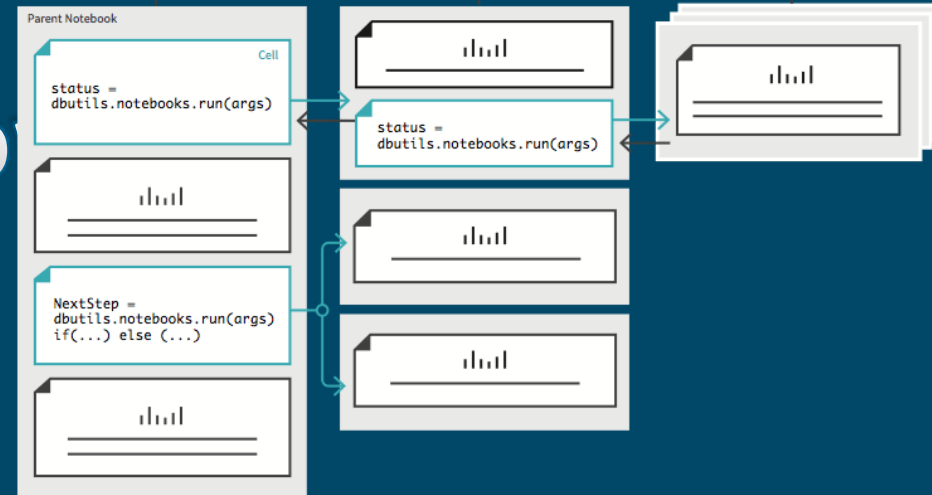
- In the example, tasks are defined as Python statement
- Tasks can also be calls to run Jupyter notebooks



Notebook workflow

Databricks Notebook Workflow

- Each notebook becomes a job to be scheduled
 - Notebooks typically do data ingest, processing, modeling, etc.
- One notebook starts the process
 - Write the workflow in PySpark
 - Can also a workflow product (Airflow, Glue)



```
dbutils.notebook.run(
    "../path/to/my/notebook",
    timeout_seconds = 60,
    arguments = {"x": "value1", "y": "value2", ...})
```

Job ID	Stage ID	Environment	Executors	Status	Duration
15	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	collect at <console>-40	2016/08/03 00:07:44	0.3 s
14	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	collect at <console>-40	2016/08/03 00:07:43	0.4 s
13	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	collect at <console>-39	2016/08/03 00:07:42	0.8 s
12	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	collect at <console>-39	2016/08/03 00:07:41	0.3 s
11	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	collect at <console>-39	2016/08/03 00:07:41	0.4 s
10	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	saveAsTable at <console>-38	2016/08/03 00:07:39	1 s
9	6122353838083489490	case class Data(id: Long) val df = sc.parallelize...	saveAsTable at <console>-37	2016/08/03 00:07:36	0.9 s
8	1886757717289855227	show tables	take at OutputAggregator.scala:80	2016/08/03 00:07:31	88 ms
7	4715541668260152232	import java.lang.management.ManagementFactory v...	first at <console>-35	2016/08/03 00:07:19	38 ms
6	4715541668260152232	import java.lang.management.ManagementFactory v...	first at <console>-35	2016/08/03 00:07:19	43 ms

Google Notebook Workflows

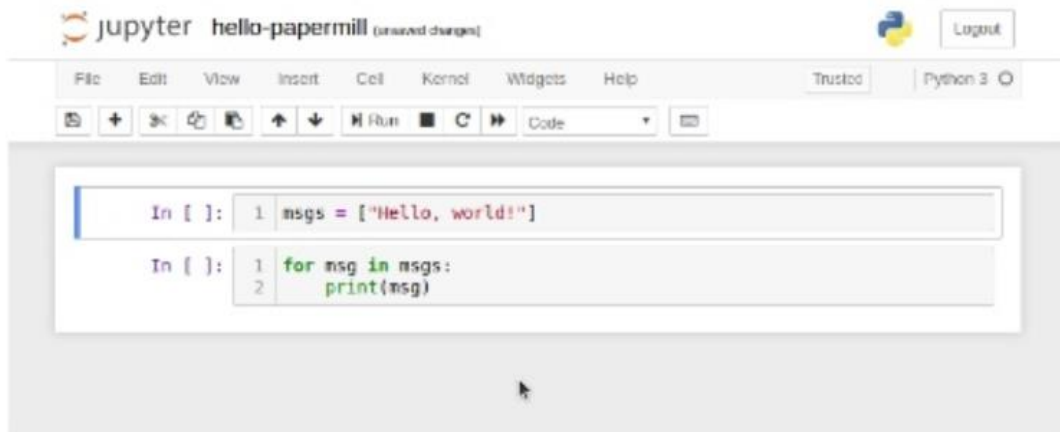
Automation

Papermill: parametrized notebooks



<https://github.com/nteract/papermill>

1. Declare some cells as a parameters cells.

A screenshot of a Jupyter Notebook interface. The title bar says 'jupyter hello-papermill (unsaved changes)'. The menu bar includes File, Edit, View, Insert, Cell, Kernel, Widgets, and Help. Below the menu is a toolbar with icons for saving, running, and other actions. The main area shows two code cells. The first cell contains 'In []: 1 msgs = ["Hello, world!"]'. The second cell contains 'In []: 1 for msg in msgs: 2 print(msg)'.

```
In [ ]: 1 msgs = ["Hello, world!"]

In [ ]: 1 for msg in msgs:
2       print(msg)
```

2. Declare parameters values in YAML file.

```
x:
  - 0.0
  - 1.0
  - 2.0
  - 3.0
linear_function:
  slope: 3.0
  intercept: 1.0
```

3. Run papermill with Python API.

```
import papermill as pm

pm.execute_notebook(
    'path/to/input.ipynb',
    'path/to/output.ipynb',
    parameters = dict(alpha=0.6, ratio=0.1)
)
```

3. Run papermill with CLI.

```
$ papermill local/input.ipynb s3://bkt/output.ipynb -f parameters.yaml
```

Notebook workflow pipeline illustrated

Call notebook

```
dbutils.notebook.run(  
    "../path/to/my/notebook",  
    timeout_seconds = 60,  
    arguments = {"x": "value1", "y": "value2", ...})
```

Notebook parameters

1. WorldData.Pipeline (Python)

File View: Code Permissions Run All Clear Publish

Cmd 4

Notebook parameters, which can be overridden by notebook workflow

```
1 params = {  
2     "file_prefix" : "WPP2019",  
3     "files" : ["_TotalPopulationBySex.csv", "_Period_Indicators_Medium.csv"],  
4     "CSV_directory": "https://population.un.org/wpp/Download/Files/1_Indicators%20(Standard)/CSV_FILES"  
5 }
```

Cmd 6

Download files

```
1 import requests  
2 for fname in params["files"]:  
3     name = params["file_prefix"] + fname  
4     url = "{}/{}/".format(params["CSV_directory"], name)  
5     print("Downloading file: ", url)  
6     r = requests.get(url)  
7     open(name, 'wb').write(r.content)  
8     print("Saved file", name)
```

Downloading file: https://population.un.org/wpp/Download/File

Data pipeline summary

- A data pipeline processes elements connected in series, where the output of one element is the input of the next one
- Common to
 - Ingest data using distributed data logging programs (Kafka)
 - Process data using PySpark
 - Data cleaning and eventually ML pipeline
 - Manage the ingest to modeling via a workflow framework
 - E.g., Apache Airflow
- Notebook workflow is a common for PySpark pipelines

Important to remember

Simplifying PySpark data transformations

Consider these kinds of DataFrame transformations

- Dropping columns
- Changing column types
- Creating a new column
 - By combining columns
 - Running a UDF

Drop duplicate columns and rename Time

```
1 dfInd = indMed.drop('Time').drop('LocID').drop('VarID').drop('Variant').withColumnRenamed('MidPeriod', 'Time')
```

Convert types of columns

```
1 # Just an illustration here because types are OK
2 dfInd = dfInd.withColumn('Births', dfInd['Births'].cast('int'))
```

Instead of in-line code, apply transformation functions

- The result looks like this

```
Cmd 10  
  
1 dfInd = indMed.transform(ind_drop_cols())\  
2         .transform(cast_to_type('Births','int'))  
3 bySex = bySex.transform(cast_to_type('Time','int'))
```

- Define the transform function, on DataFrame
 - It applies your given function to the DataFrame and returns the update data frame
 - `df = function(df)`

Monkey patch

- Allows the definition of new functions on existing classes (in Python)

```
1  from pyspark.sql.dataframe import DataFrame
2
3  #####
4  # DataFrame transformations monkey patches
5  # See https://medium.com/@mrpowers/chaining-custom-pyspark-transformations-4f38a8c7ae55
6  #####
7
8  def transform(self, f):
9      return f(self)
10 DataFrame.transform = transform
```


Two example functions to apply to a DataFrame

- Function specific to this problem
 - Dropping specific columns
- Generic function
 - Can be reused in any data pipeline

```
1 # Used to simplify the code
2 def ind_drop_cols():
3     def inner(df):
4         df = df.drop('Time').drop('LocID').drop('VarID').drop('Variant')\
5             .withColumnRenamed('MidPeriod','Time').where("Time < 2019 and Variant = 'Medium'")
6         return df
7     return inner
8
9
10 def cast_to_type(col,type):
11     def inner(df):
12         return df.withColumn(col,df[col].cast(type))
13     return inner
```

Before and after using transform

```
dfInd = indMed.drop('Time').drop('LocID').drop('VarID').drop('Variant').withColumnRenamed('MidPeriod','Time').where("Time < 2019 and Variant = 'Medium'")
```

```
# Just an illustration here because types are OK
dfInd = dfInd.withColumn('Births', dfInd['Births'].cast('int'))
bySex = bySex.withColumn('Time', bySex['Time'].cast('int'))
```



```
1 # Used to simplify the code
2 def ind_drop_cols():
3     def inner(df):
4         df = df.drop('Time').drop('LocID').drop('VarID').drop('Variant')\
5             .withColumnRenamed('MidPeriod','Time').where("Time < 2019 and Variant = 'Medium'")
6         return df
7     return inner
8
9
10 def cast_to_type(col,type):
11     def inner(df):
12         return df.withColumn(col,df[col].cast(type))
13     return inner
```

```
1 dfInd = indMed.transform(ind_drop_cols())\
2     .transform(cast_to_type('Births','int'))
3 bySex = bySex.transform(cast_to_type('Time','int'))
```

A function that returns a function

- Input

- Two parameters

- Process

- Defines a new function (inner)
- Casts the given column to the given type
- Return that function

- How to call it

```
10 def cast_to_type(col,type):  
11     def inner(df):  
12         return df.withColumn(col,df[col].cast(type))  
13     return inner
```

```
dfInd = indMed.transform(ind_drop_cols())\  
            .transform(cast_to_type('Births','int'))
```

Monkey patch for
DataFrame transform
can simplify (read, write, maintain) PySpark
data pipelines