

# **COMPLETE PySpark with ME**

**- Nishant Kolapkar**

## 1. What is Big Data and How it Started

- Cobol is the first Business Data Processing Language it allowed to store data in files, create index files, and process data efficiently. (1959)
- Then mySQL Oracle(1977) MSSQL comes in picture.
- RDBMS – What do they offer?
  - SQL- An easy Data Query Language
  - Scripting Languages such as PL/SQL and Transact SQL
  - Interface for other programming language such as JDBC and ODBC(interact with data with languages)
- Data Categories
  - Structured (rows and commas)
  - Semi-Structured

XML

VS.

JSON

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <endereco>
3   <cep>31270901</cep>
4   <city>Belo Horizonte</city>
5   <neighborhood>Pampulha</neighborhood>
6   <service>correios</service>
7   <state>MG</state>
8   <street>Av. Presidente Antônio Carlos, 6627</street>
9 </endereco>
```

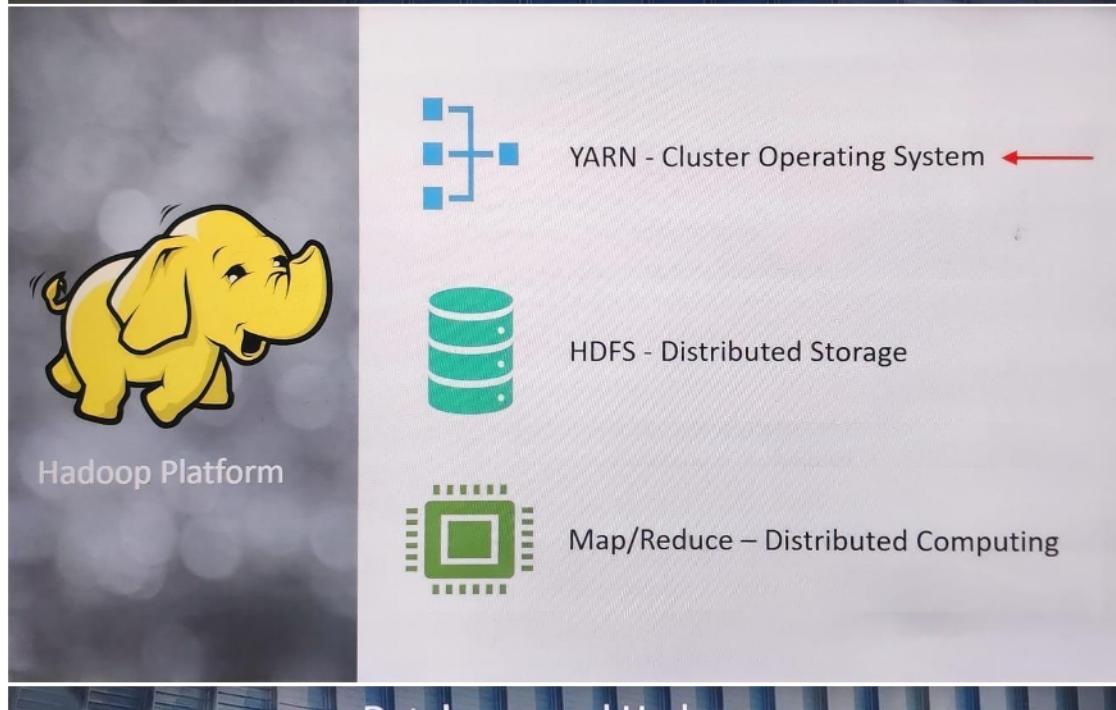
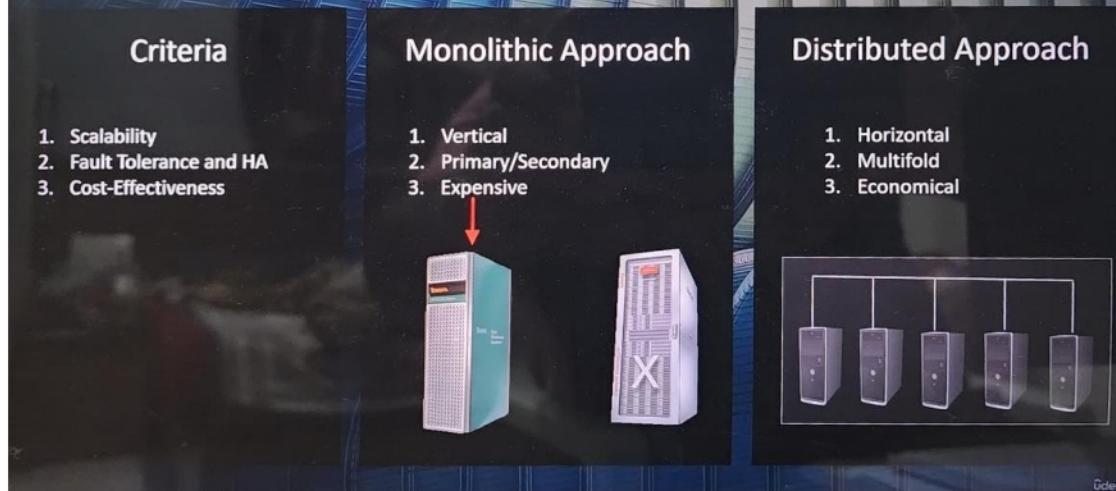
```
1 { "endereco": {
2   "cep": "31270901",
3   "city": "Belo Horizonte",
4   "neighborhood": "Pampulha",
5   "service": "correios",
6   "state": "MG",
7   "street": "Av. Presidente Antônio Carlos, 6627"
8 }
9 }
```

- Un-Structured
- Big Data Problem
  - Variety of unstructured Data
  - Huge Amount of Data Volume
  - Velocity (Huge data developed in few sec.)
- Big Data Platform Requirements
  - Store High Volumes of data arriving at higher velocity.
  - Accommodate structured, semi-structured, and unstructured data variety.
  - Process High volume of a variety of data at a higher velocity.

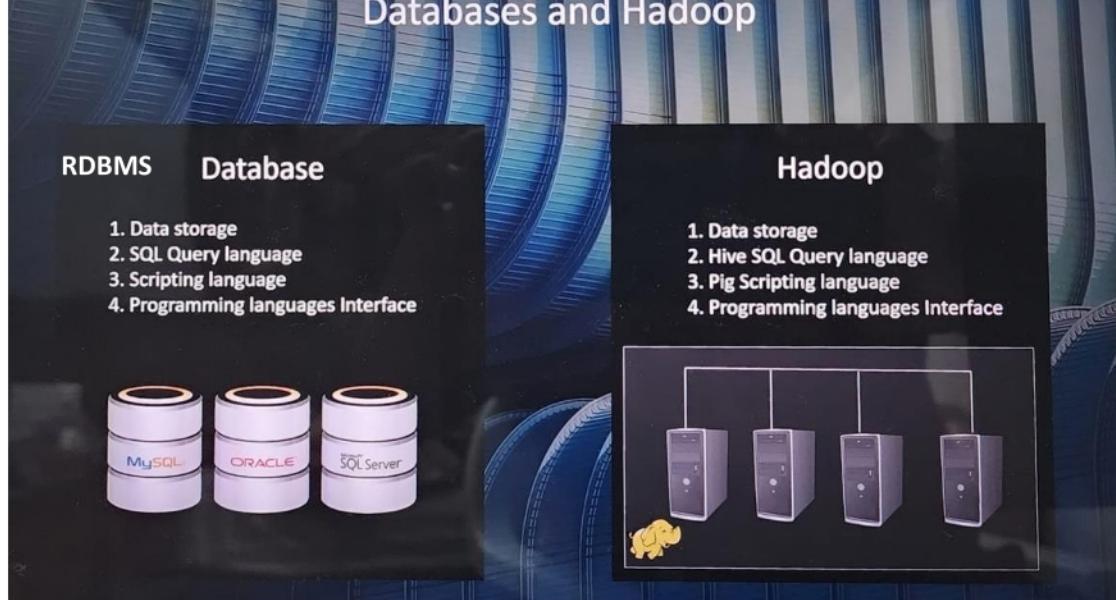
## Summary

- Data processing is one of the critical business requirements
- We used RDBMS for decades to develop data processing applications
- The advent of the internet started putting following data processing challenges
  - We started collecting a Variety of Data
    - Structured data
    - Semi-Structured data
    - Un-Structured data
  - Businesses started collecting high Data Volumes
  - Need to collect and process data at a high velocity
- The new data challenge is popularly known as Big Data Problem
- The Big Data problem was defined using the 3Vs of Big data - Variety, Volume, and Velocity
- • RDBMS failed to handle the Big Data problem
  - Industry needed a new approach or platform to handle the Big Data Problem

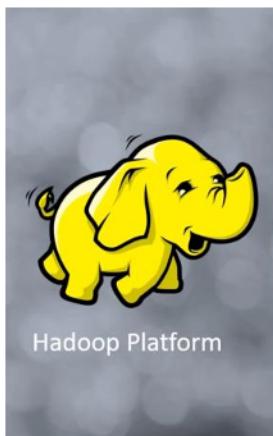
## Approaches of Big Data Solution



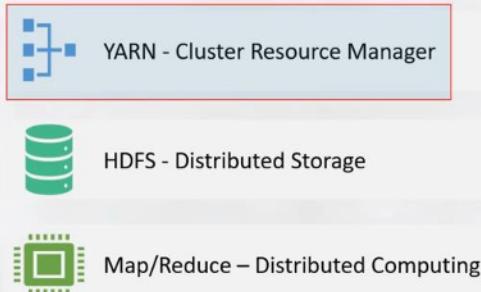
## Databases and Hadoop



- What is Hadoop?



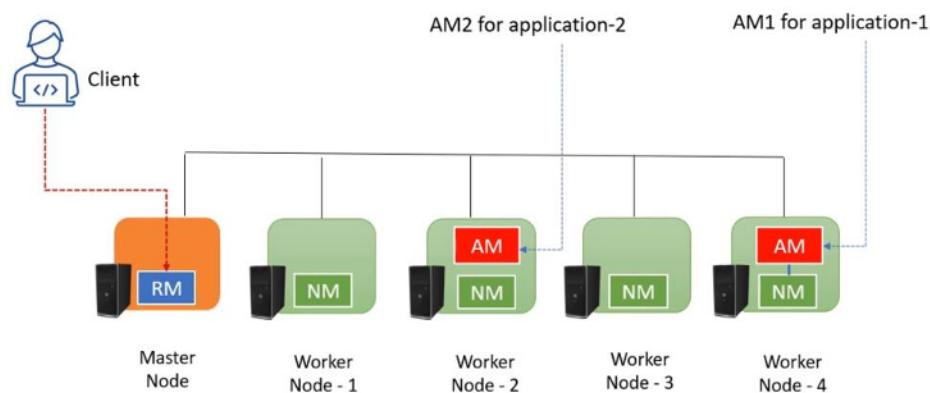
Hadoop is a distributed data processing platform that offers the following core capabilities.



- YARN is Hadoop Cluster OS but it called as Cluster Resource Manager.

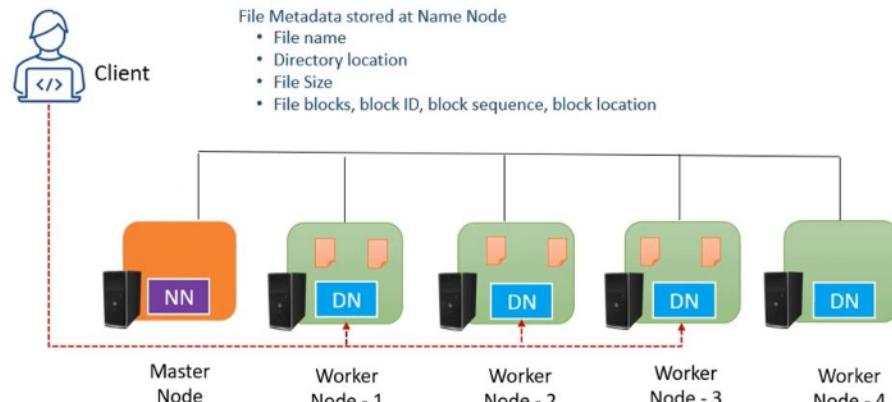
## YARN – Yet Another Resource Manager

- Hadoop cluster Operating System
- Popularly known as Hadoop Cluster Resource Manager
- Has three main components
  - RM - Resource Manager
  - NM - Node Manager
  - AM - Application Master



## HDFS – Hadoop Distributed File System

- Distributed Storage on Hadoop Cluster
- Has two main components
  - NN – Name Node
  - DN – Data Node



Data will be split into blocks and stored on nodes normally this blocks size is 128mb. And this all information of blocks are stored in Metadata and stored at Name Node. For read operation by using this metadata info all blocks are combined and retrieve to client from cluster.

## M/R – Hadoop Map Reduce Framework

- Map Reduce is
  - Programming Model
  - • Programming Framework

### Problem Statement

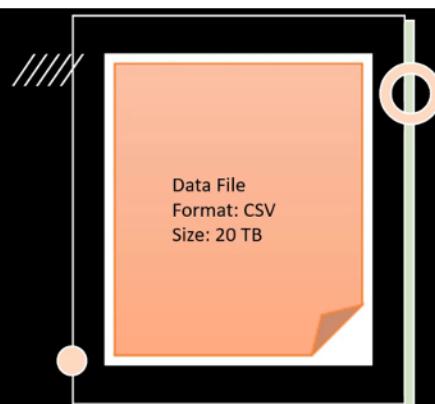
Count the lines in the given file

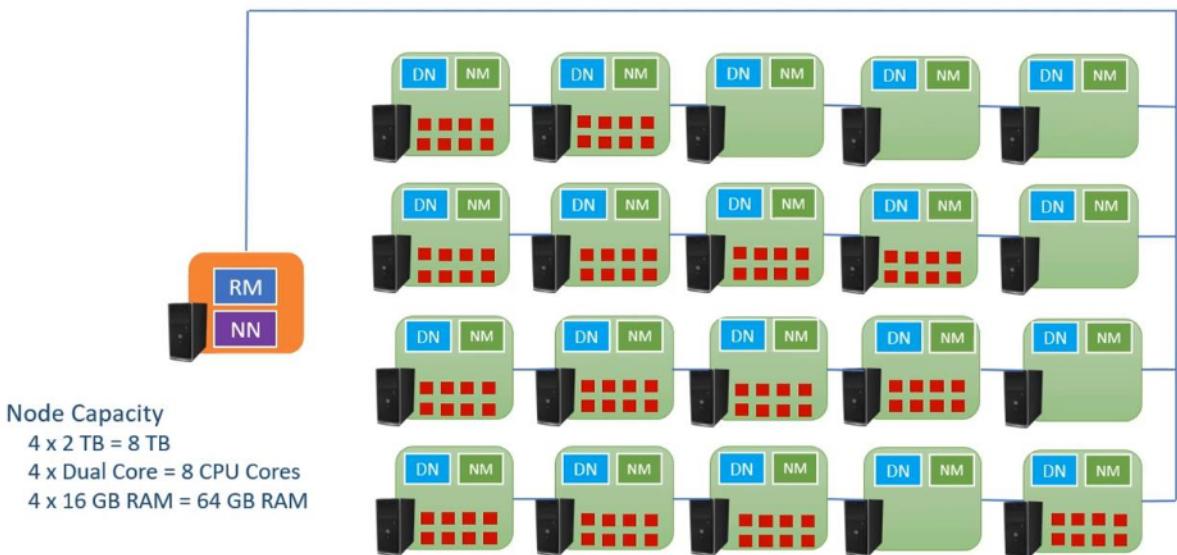
### Solution Pseudocode

```
open file as f_hd
  for each t_line in f_hd.get_line()
    n_count = n_count + 1
  close f_hd
  print n_count
```

### Challenges

1. Storage capacity
2. Processing time





## Problem Statement

Count the lines in the given file

## Distributed Solution Pseudocode

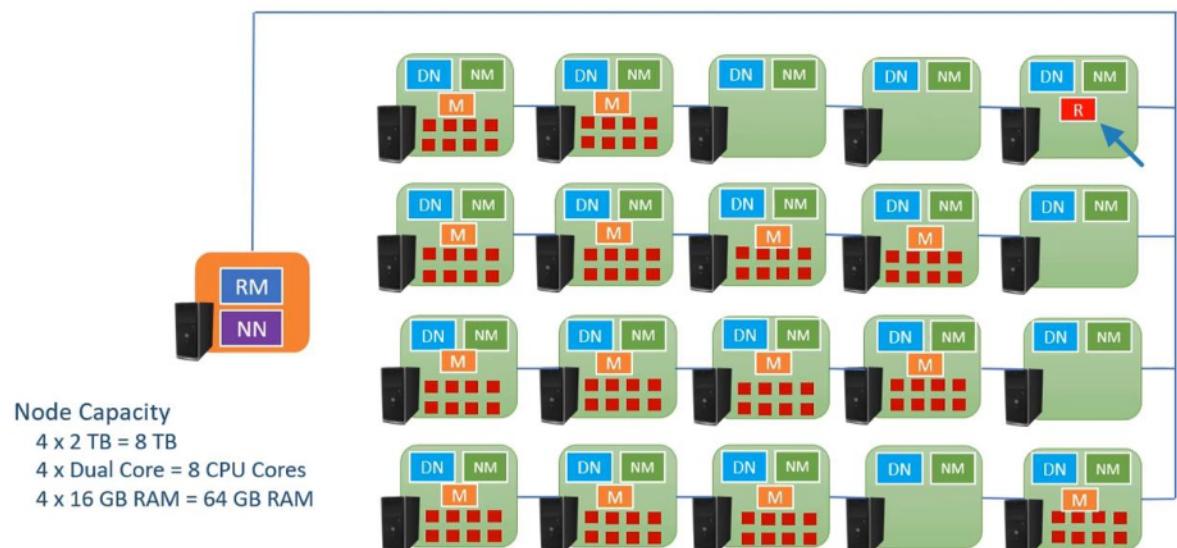
```
def map(file_block):
    open file_block as fb_hd
    for each t_line in fb_hd.get_line()
        n_count = n_count + 1

    close fb_hd
    return n_count

def reduce(list_counts):
    for each cnt in list_counts
        total_count = total_count + cnt

print total_count
```

Data File  
Format: CSV  
Size: 20 TB



## Summary

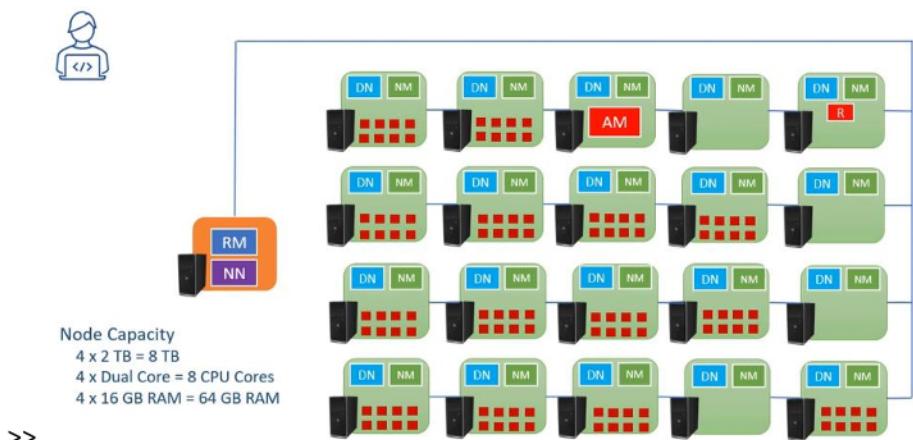
### Map Reduce Model

Implement logic in two functions

1. Map Function
  - Reads data block
  - Applies logic at block level
  - Map output is sent to Reduce
2. Reduce Function
  - Receives Map output
  - Consolidates the results

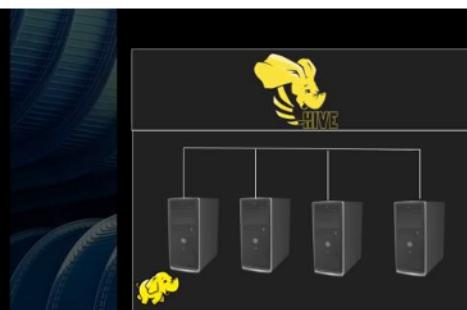
Hadoop M/R framework implement the map-reduce model.

- YARN manages resource allocation
- HDFS manages data blocks



### Hive capabilities

1. Create
  1. Databases
  2. Tables
  3. Views
2. Run SQL Queries



## Entering Apache Spark

### Advantages over Hadoop

1. Performance
  - 10 to 100 times faster than Hadoop M/R
2. Ease of development
  - Spark SQL
  - High performance SQL Engine
  - Composable Function API
3. Language support
  - Java, Scala, Python and R
4. Storage
  - HDFS Storage
  - Cloud Storage
5. Resource Management
  - YARN, Mesos, Kubernetes



Runs in two setups

1. With Hadoop (Data Lake)
2. Without Hadoop (Lakehouse)

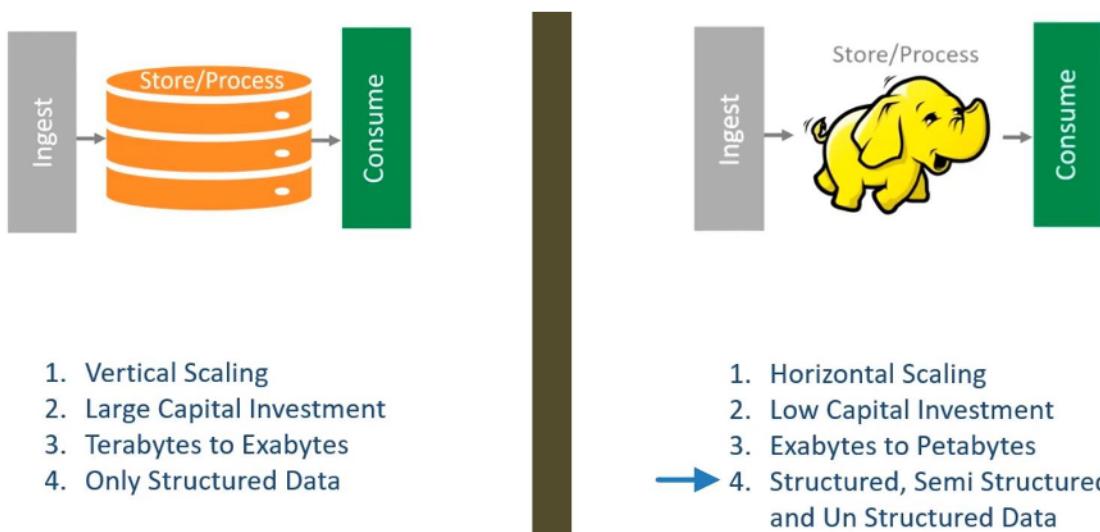
### Spark Solutions

- ▶ 1. Data Lake – On Hadoop
- 2. Lakehouse – On Cloud

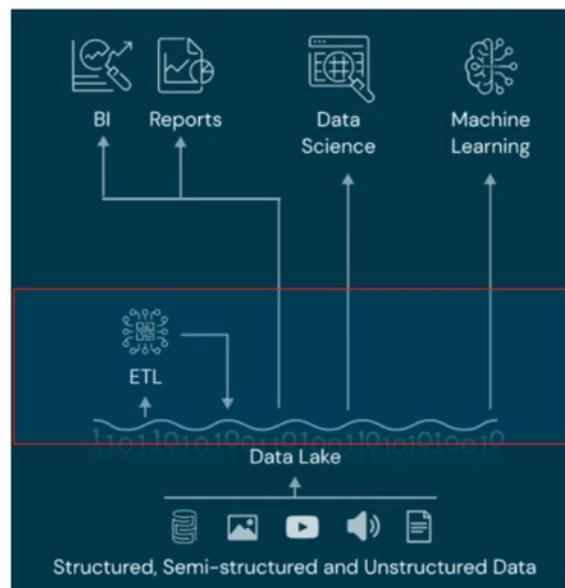


- What is Data Lake and How it works

- Google File System(GFS) >> Hadoop Distributed File System(HDFS with Map Reduce)



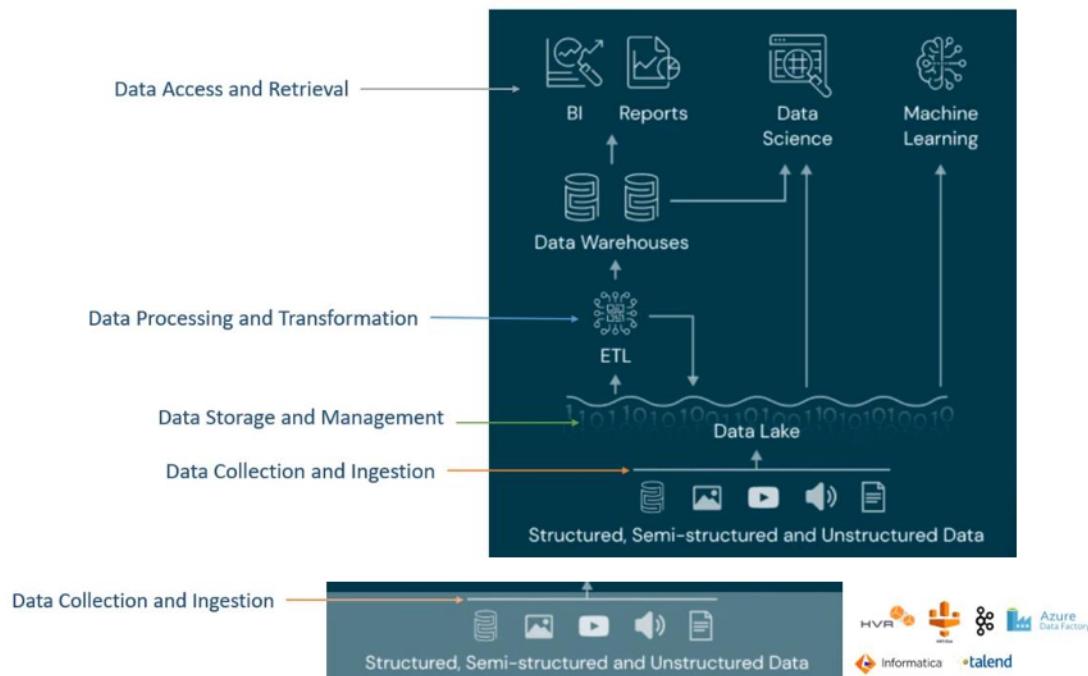
## Data Lake

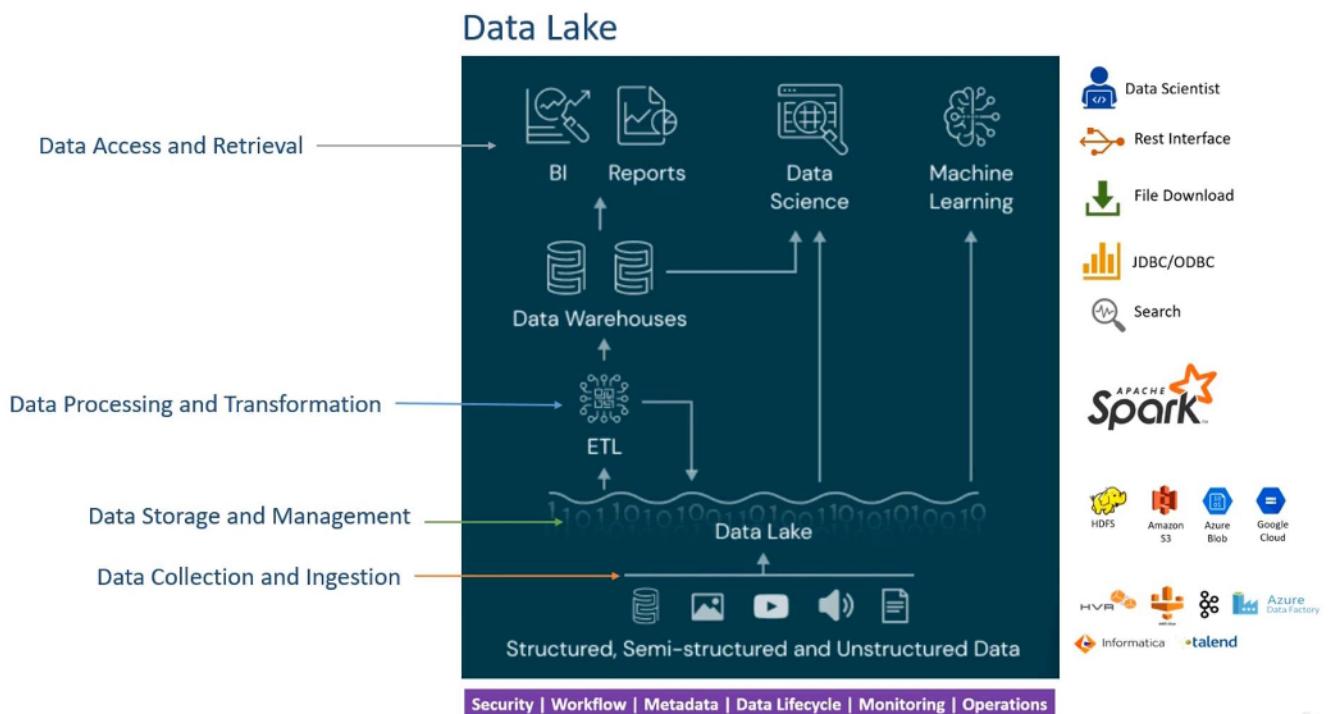
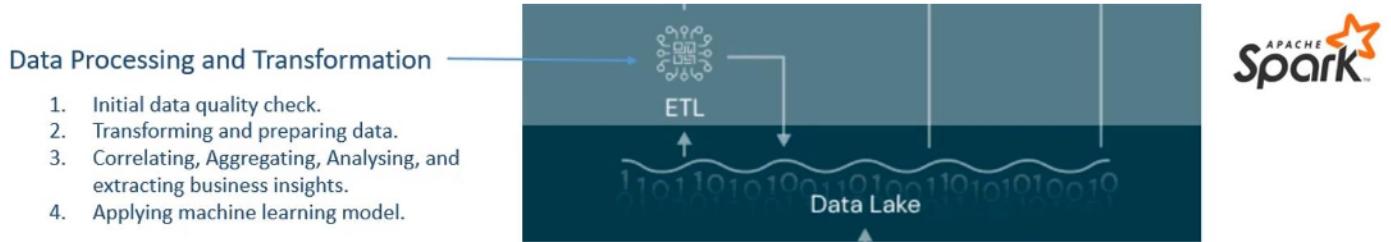


### Key Missing Features

1. Transaction and Consistency
2. Reporting Performance

## Data Lake

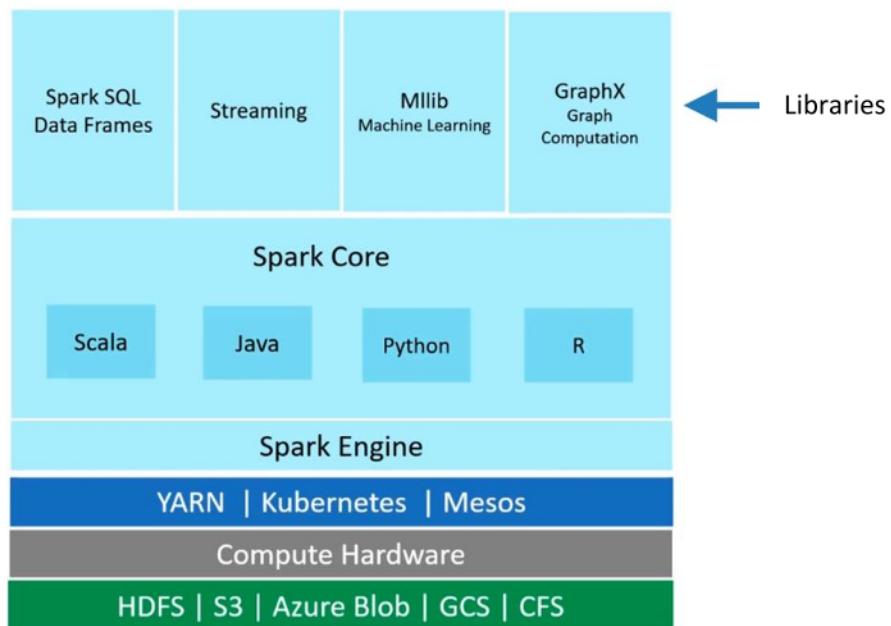




- Introducing Apache Spark and Databricks Cloud



- [What is Apache Spark?](#)



1. Spark on the Cloud Platform
2. Spark Cluster Management
3. Notebooks and Workspace
4. Administration Controls
5. Optimized Spark
6. Databases/Tables and Catalog
7. Databricks SQL Analytics
8. Delta Lake Integration
9. ML Flow
10. Industry vertical accelerators

- Installation Using Apache Spark

1. Cloud Platforms

- Notebook
- Databricks Cloud

2. On-premise Platforms

- Python IDE
- Cloudera Platform

ref. 8

```
Microsoft Windows [Version 10.0.19045.2486]
(c) Microsoft Corporation. All rights reserved.

C:\Users\innk>setx JAVA_HOME "C:\Program Files\Java\jdk-11"
SUCCESS: Specified value was saved.

C:\Users\innk>setx PATH "%PATH%;%JAVA_HOME%\bin"
WARNING: The data being saved is truncated to 1024 characters.
SUCCESS: Specified value was saved.

C:\Users\innk>_
```

```
C:\Users\innk>java -version
openjdk version "11" 2018-09-25
OpenJDK Runtime Environment 18.9 (build 11+28)
OpenJDK 64-Bit Server VM 18.9 (build 11+28, mixed mode)
```

```
C:\Users\innk>setx HADOOP_HOME "E:\pySpark_soft\hadoop-3.2.2"
SUCCESS: Specified value was saved.
```

```
| SPARK_HOME           E:\pySpark_soft\spark_3
| E:\pySpark_soft\hadoop-3.2.2\bin
| E:\pySpark_soft\spark_3\bin
```

```
C:\Users\innk>setx PYTHONPATH "E:\pySpark_soft\spark_3\python;E:\pySpark_soft\spark_3\python\libpy4j-0.10.9.5-src.zip"
SUCCESS: Specified value was saved.

C:\Users\innk>where python
C:\Program Files\Python311\python.exe

C:\Users\innk>setx PYSPARK_PYTHON "C:\Program Files\Python311\python.exe"
SUCCESS: Specified value was saved.
```

Need to install python v3.10.9

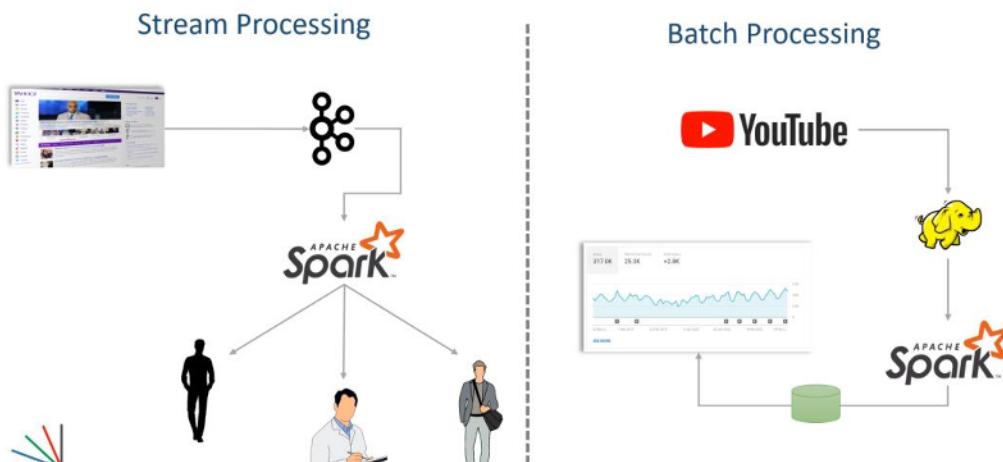
- **Spark Execution Model and Architecture**

- What is Apache Spark? A Distributed Computing Platform.
- What do we do with Spark? We create program and execute them on Cluster.
  - o How to create Spark Programs?
  - o How to execute Spark Programs?

## How to execute Spark Programs?

### 1. Interactive Clients

spark-shell, Notebook



For both scenarios output appears on dashboard.

### 2. Submit Job

spark-submit (allow to submit job to the cluster)

## How to execute Spark Programs?

### 1. Interactive Clients

spark-shell, Notebook

### 2. Submit Job

spark-submit, Databricks Notebook, Rest API

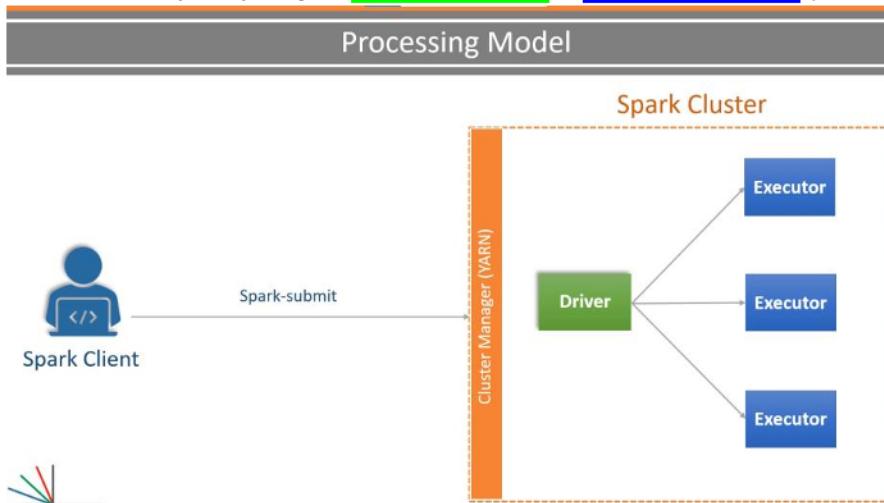
- Interactive clients are used for learning, Dev and Explorations.

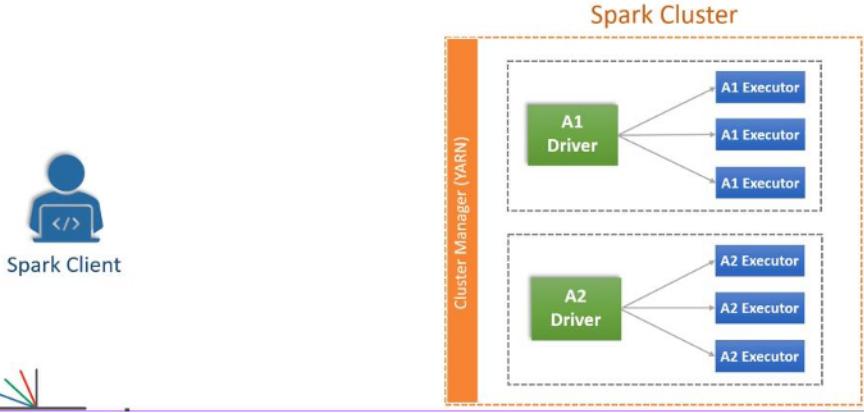
- Packaging of the Spark applications and submitting on the cluster for execution.

- Spark Distributed Processing Model – How your program runs?

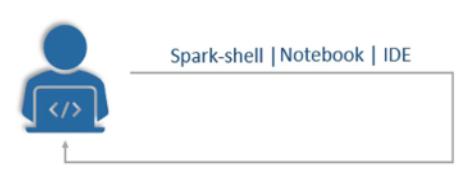
### Spark Distributed Processing Model?

- Spark applies to a master slave architecture to every application. When we submit application its create master process for the application and this master process is create bunch of slaves to distribute the work and compute your job. **Master = Drivers** & **Slaves = Executors** (with respect to clusters its different)



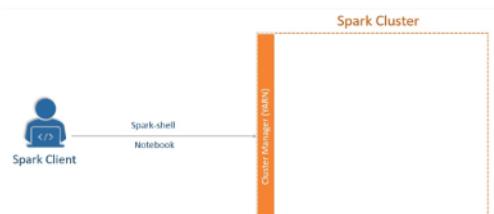


- Every Spark application applies a master slave architecture and runs independently on the cluster.
- Spark Execution Modes and Cluster Managers (What happened after you submit Spark Application?)



How Spark works on Local Machine?

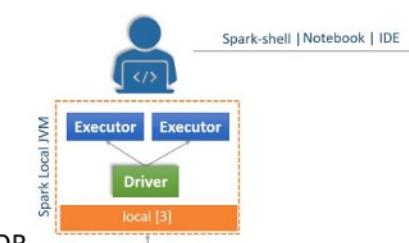
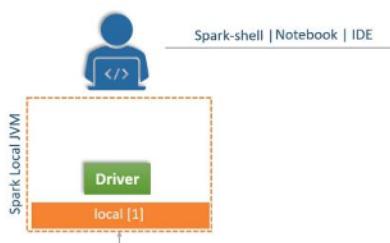
How Spark works with Interactive Clients?



- We can run Spark application on local without having cluster
- Configuration application to run on **verity of Clusters** Spark is compatible with following Clusters.

```
1. local[n]
2. YARN
3. Kubernetes
4. Mesos
5. Standalone >>
sparkconf x
[SPARK_APP_CONFIGS]
1 spark.app.name = HelloSpark
2 spark.master = local[3]
3
4 spark.sql.shuffle.partitions = 2
```

Here is configuration which tells about the target cluster manager

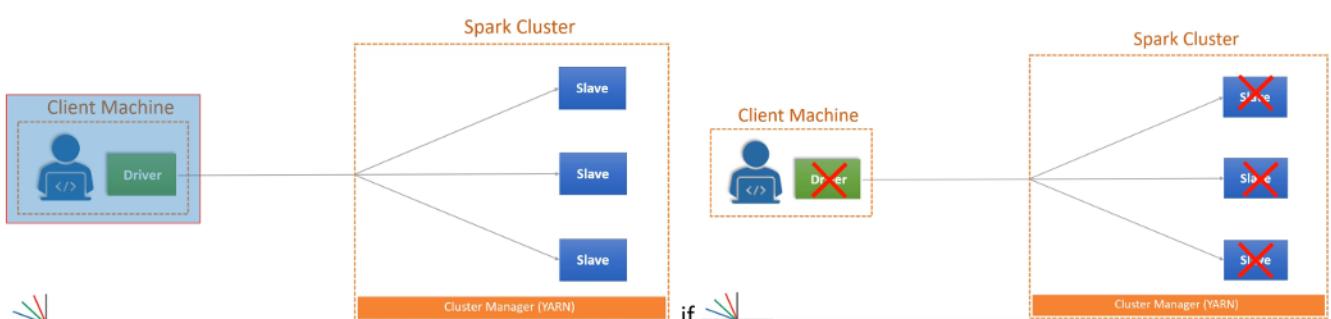


OR

(running testing debugging application locally)

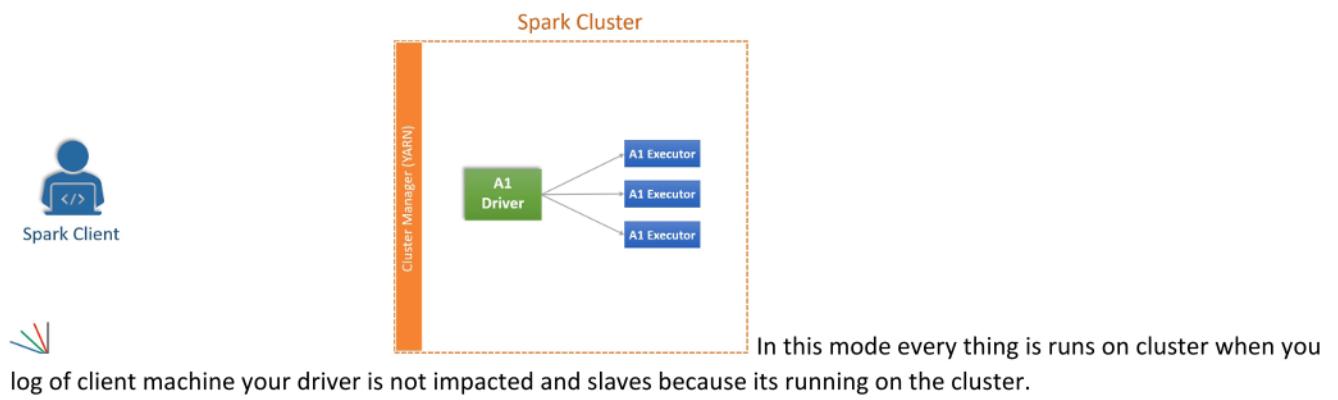
### How to Spark with interactive clients?

#### 1. Client mode



Slaves also dies of the absence of the driver. Its suitable for interactive work and not for long running jobs.

## 2. Cluster mode



- Three widely used execution models

Cluster Managers	Execution Modes	Execution Tools
1. local[n]	1. Client	1. IDE, Notebook
2. YARN	2. Client	2. Notebook, Shell
3. YARN	3. Cluster	3. Spark Submit

- Working with PySpark Shell

Cluster	Mode	Tool
Local	Client Mode	spark-shell, Notebook

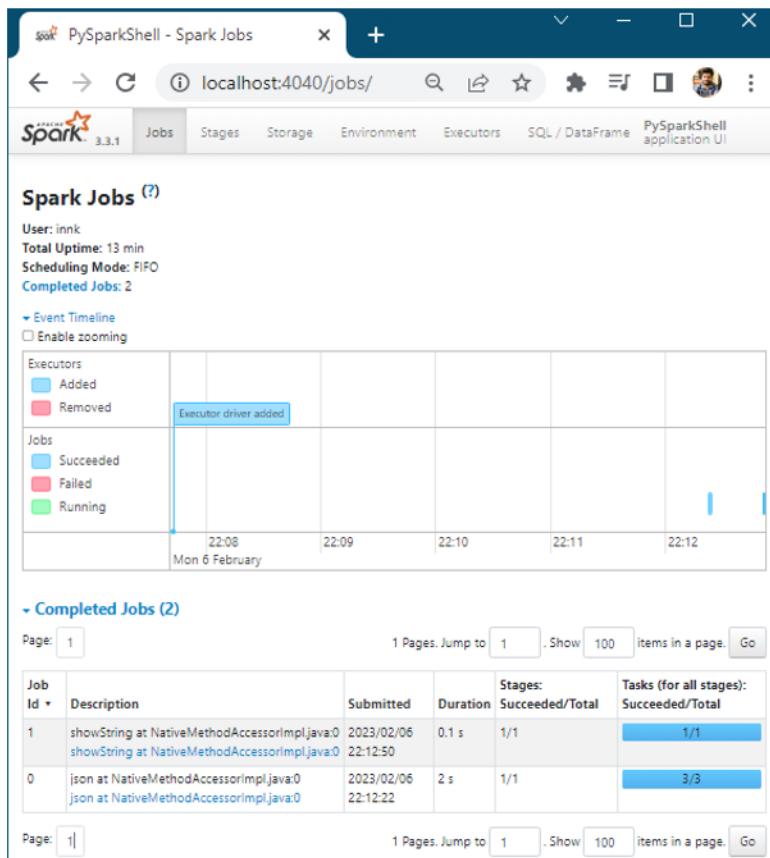
- Spark Shell on local Client Mode:

pyspark --help

```
>pyspark --master local[3] --driver-memory 2G To increase memory to 2GB (only for 1st time)
>>> df = spark.read.json("C:/demo/notebook/data/people.json")
>>> df.show()
+---+-----+
| age| name|
+---+-----+
| null| Prashant|
| 30 | Abdul|
| 19 | Justin|
| 43 | Andy|
+---+-----+
```

- Spark context Web UI:

<http://localhost:4040/jobs/> (This UI is available only when Spark application is in running state)



- Installing Multi-Node Spark Cluster

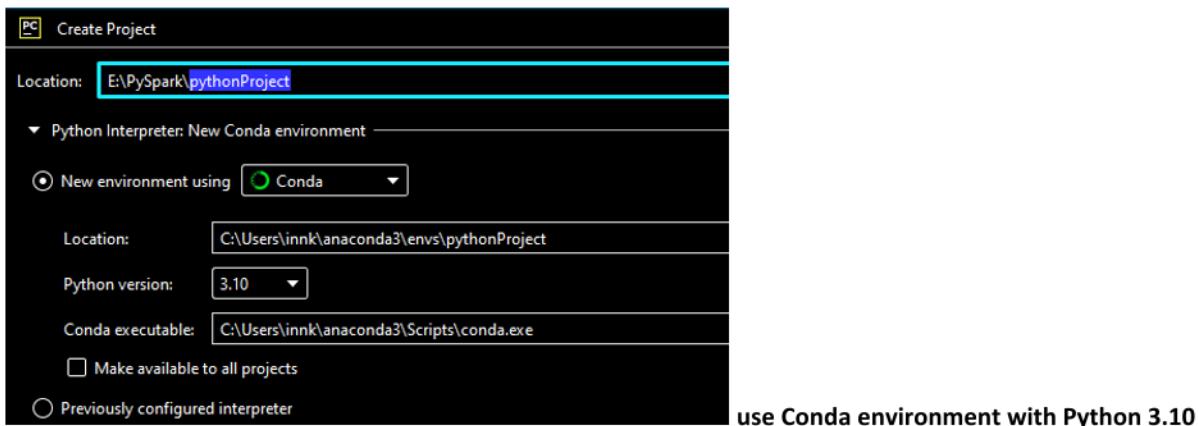
[How to setup Spark Cluster in GCP?](#)

- YARN Client Mode Using Spark-Shell (Directly on production cluster, here GCP) [ref. 21, 22](#)

Cluster	Mode	Tool
YARN	Client Mode	spark-shell, Notebook

[YARN Dynamic policy | Zeppelin](#)

- PySpark Project One



[pip install pyspark](#)

[pip install pytest](#)

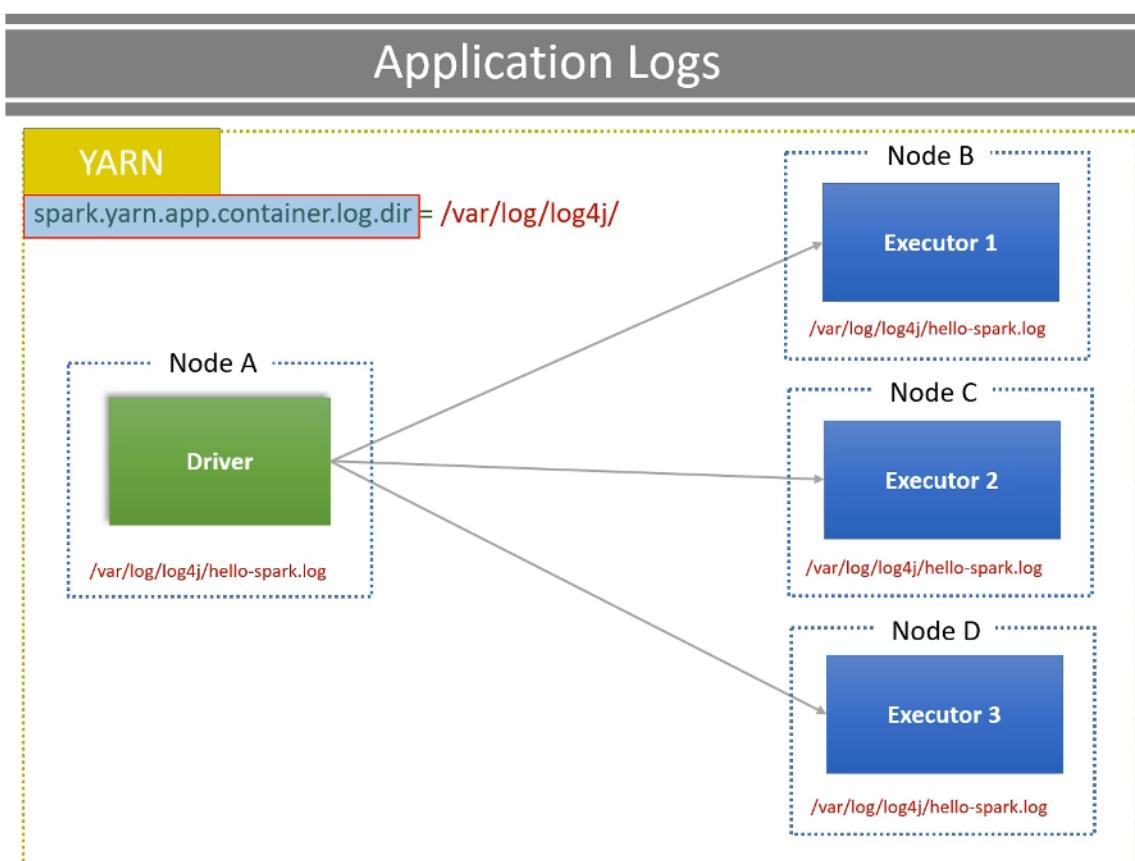
- Spark Application Logs

### Using Log4J

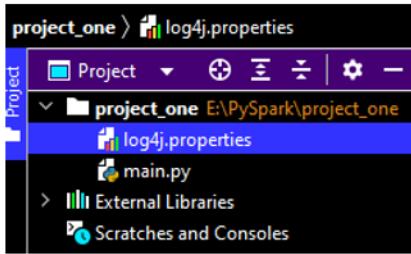
1. Create a Log4J configuration file
2. Configure Spark JVM to pickup the Log4J configuration file
3. Create a Python Class to get Spark's Log4J instance and use it

Instead of python built in logger why Log4J is used?

Because the Python Logger is not integrated with Spark.



YARN only collects logs only from [/var/log/log4j/location](#)



copy log4j.properties to main project folder

```
#define rolling file appender
log4j.appender.file=org.apache.log4j.RollingFileAppender
log4j.appender.file.File=${spark.yarn.app.container.log.dir}/${logfile.name}.log
```

## How to configure JVM variables?

Java Virtual Machine

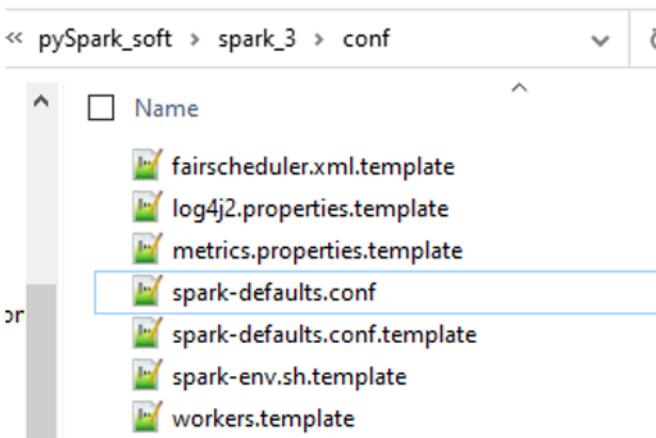
## Setting JVM Parameters

SPARK\_HOME/conf/spark-defaults.conf

Name	Date modified	Type	Size
bin	5/31/2020 12:53 PM	File folder	
conf	5/31/2020 1:07 PM	File folder	
data	5/31/2020 12:53 PM	File folder	
examples	5/31/2020 12:53 PM	File folder	
jars	5/31/2020 12:53 PM	File folder	
kubernetes	5/31/2020 12:53 PM	File folder	
licenses	5/31/2020 12:53 PM	File folder	
python	5/31/2020 12:53 PM	File folder	
R	5/31/2020 12:53 PM	File folder	
sbin	5/31/2020 12:53 PM	File folder	
yarn	5/31/2020 12:53 PM	File folder	
LICENSE	2/3/2020 1:17 AM	File	21 KB
NOTICE	2/3/2020 1:17 AM	File	42 KB
README.md	2/3/2020 1:17 AM	MD File	4 KB
RELEASE	2/3/2020 1:17 AM	File	1 KB

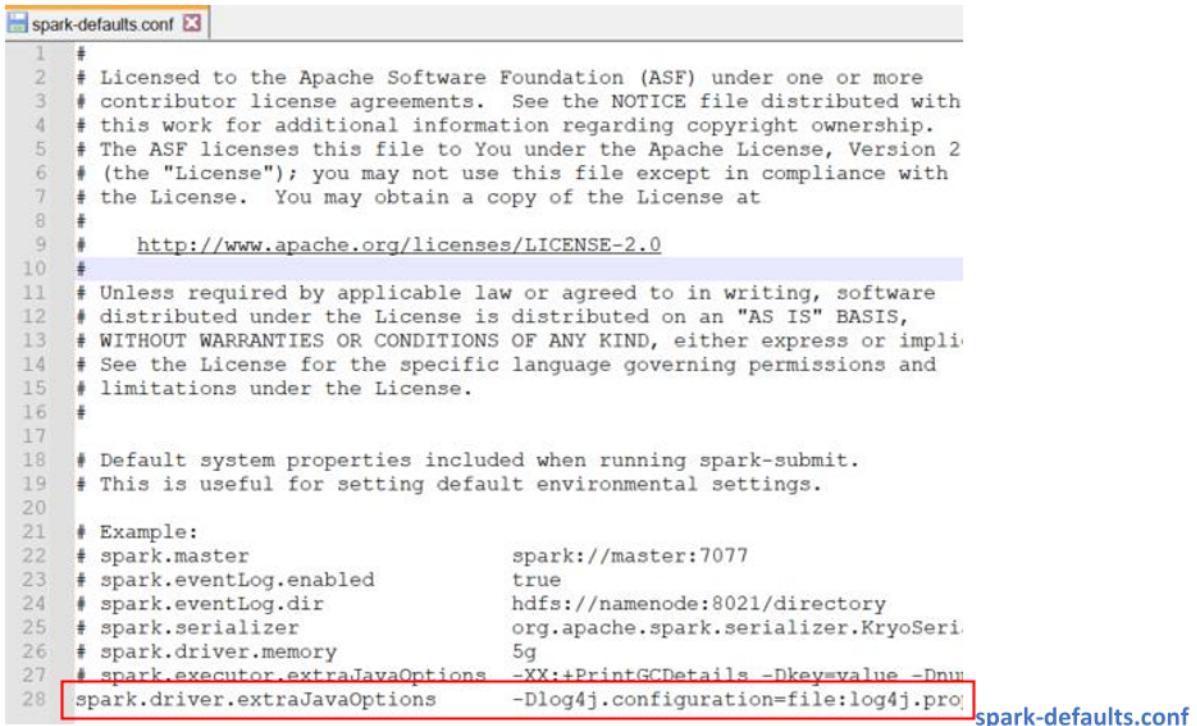
This is the place where Spark is looking

for the configuration dir. SPARK\_HOME



spark-defaults.conf

```
spark.driver.extraJavaOptions -Dlog4j.configuration=file:log4j.properties -DSpark.yarn.app.container.log.dir=app-logs -Dlogfile.name=First_Spark
```



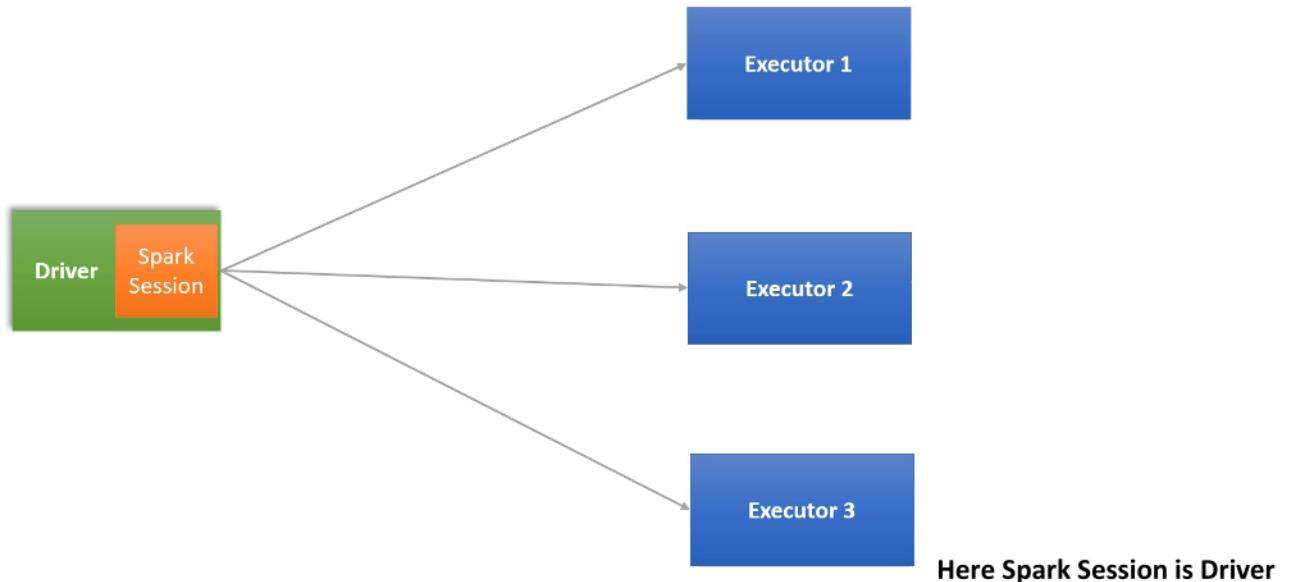
```
spark-defaults.conf
```

```
1 #
2 # Licensed to the Apache Software Foundation (ASF) under one or more
3 # contributor license agreements. See the NOTICE file distributed with
4 # this work for additional information regarding copyright ownership.
5 # The ASF licenses this file to You under the Apache License, Version 2
6 # (the "License"); you may not use this file except in compliance with
7 # the License. You may obtain a copy of the License at
8 #
9 #     http://www.apache.org/licenses/LICENSE-2.0
10 #
11 # Unless required by applicable law or agreed to in writing, software
12 # distributed under the License is distributed on an "AS IS" BASIS,
13 # WITHOUT WARRANTIES OR CONDITIONS OF ANY KIND, either express or implied.
14 # See the License for the specific language governing permissions and
15 # limitations under the License.
16 #
17 #
18 # Default system properties included when running spark-submit.
19 # This is useful for setting default environmental settings.
20 #
21 # Example:
22 # spark.master          spark://master:7077
23 # spark.eventLog.enabled true
24 # spark.eventLog.dir    hdfs://namenode:8021/directory
25 # spark.serializer      org.apache.spark.serializer.KryoSerializer
26 # spark.driver.memory   5g
27 # spark.executor.extraJavaOptions -XX:+PrintGCDetails -Dkey=value -Dnu
28 spark.driver.extraJavaOptions -Dlog4j.configuration=file:log4j.pro
```

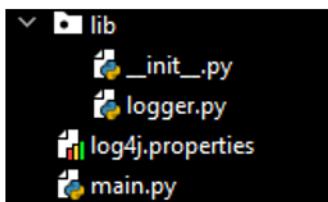
spark-defaults.conf

- Creating Spark Session

## Spark Session



Steps :



create package name lib at project level create logger.py

```
class Log4j(object):
    def __init__(self, spark):
        root_class = "guru.learningjournal.spark.examples"
        conf = spark.sparkContext.getConf()
        app_name = conf.get("spark.app.name")
        log4j = spark._jvm.org.apache.log4j
        self.logger = log4j.LogManager.getLogger(root_class + "." + app_name)

    def warn(self, message):
        self.logger.warn(message)

    def info(self, message):
        self.logger.info(message)

    def error(self, message):
        self.logger.error(message)

    def debug(self, message):
        self.logger.debug(message)

logger.py
```

## main.py

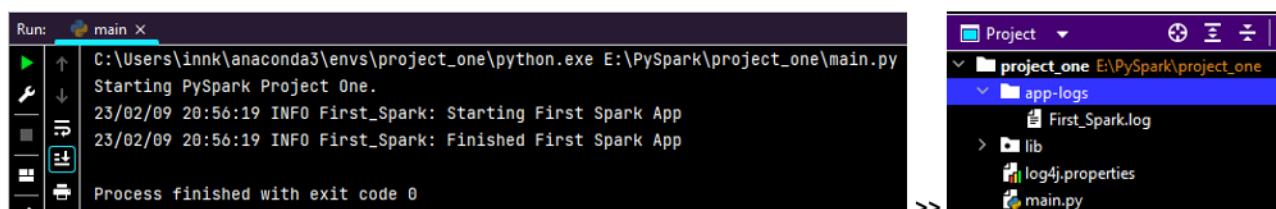
```
from pyspark.sql import *
from lib.logger import Log4j

if __name__ == "__main__":
    print("Starting PySpark Project One.")

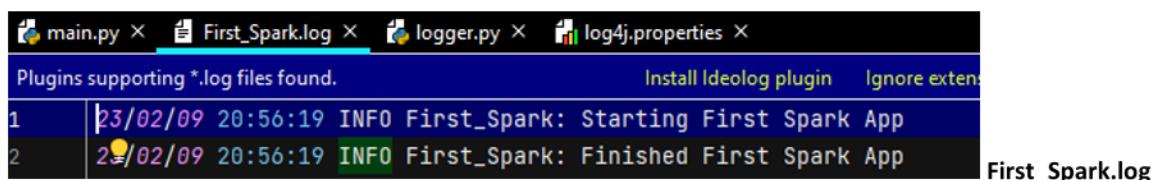
    spark = SparkSession.builder \
        .appName("First_Spark") \
        .master("local[3]") \
        .getOrCreate()

    logger = Log4j(spark)
    logger.info("Starting First Spark App")
    # Your processing code

    logger.info("Finished First Spark App")
    spark.stop()
```



New app log folder is created after executing [main.py](#)



Question 1:

SparkSession

- is an entry point to Programming Spark
- can be created using SparkSession.builder()
- both 1 and 2
- none of the above

Question 2:

A SparkContext

- represents the connection to a Spark cluster
- can be used to get SparkConf object
- was an entry point to Programming Spark in older versions
- All of the above

Question 3:

SparkConf

- is used to get/set various Spark parameters as key-value pairs
- can be retrieved from your spark session using spark.sparkContext.getConf()
- You can create a SparkConf() that loads defaults from system properties and the classpath
- setting values to SparkConf directly takes priority over system properties
- All of the above

Question 4:

You can set the configuration to your Spark application

- using spark-submit command-line options
- setting environment variable on the machine where you run spark-submit
- using SPARK\_HOME/conf/spark-defaults.conf on the machine where you run spark-submit
- setting key/value pair to SparkConf object in your application code
- all of the above

Question 5:

Following Spark application configuration precedence is correct

- environment variable -> spark-submit command line -> SparkConf -> spark-defaults.conf
- spark-submit command line -> environment variable -> spark-defaults.conf -> SparkConf
- environment variable -> spark-submit command line -> spark-defaults.conf -> SparkConf
- environment variable -> spark-defaults.conf -> spark-submit command line -> SparkConf

- Configuring Spark Session

## Configuring Spark Session

1. Environment variables
2. SPARK\_HOME\conf\spark-defaults.conf
- 3. spark-submit command-line options
4. SparkConf Object

2<sup>nd</sup> method is used for special cases. 3<sup>rd</sup> and 4<sup>th</sup> is used by Dev

3<sup>rd</sup>

```
Microsoft Windows [Version 10.0.18362.836]
(c) 2019 Microsoft Corporation. All rights reserved.
```

```
C:\demo\spark245>bin\spark-submit --master local[3] --conf "spark.app.name=Hello Spark" --conf spark.eventLog.enabled=false HelloSpark.py
```

- Config can also possible with above submit command-line which is 3<sup>rd</sup> one.
- IF you have space in config use “ ”.

#### 4<sup>th</sup> Quote the configuration in Application

<https://spark.apache.org/docs/latest/configuration.html#application-properties>

```
from pyspark import SparkConf
from pyspark.sql import *

from lib.logger import Log4J

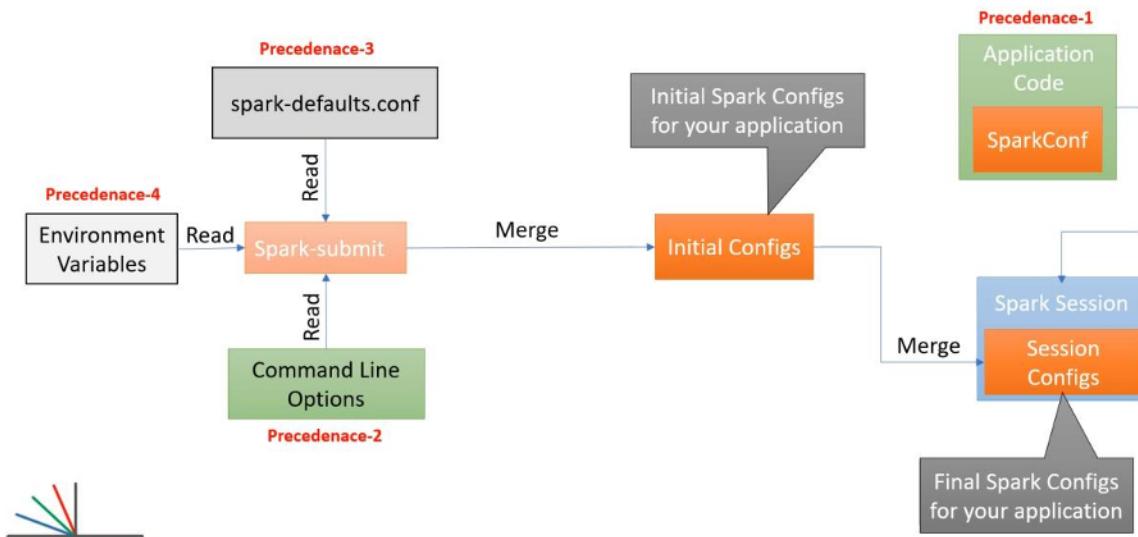
if __name__ == "__main__":
    conf = SparkConf()
    conf.set("spark.app.name", "Hello Spark")
    conf.set("spark.master", "local[3]")
    spark = SparkSession.builder \
        .config(conf=conf) \
        .getOrCreate()

    logger = Log4J(spark)

    logger.info("Starting HelloSpark")
    # Your processing code

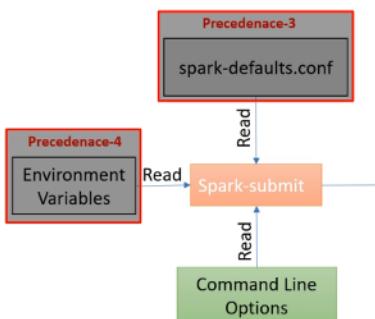
    logger.info("Finished HelloSpark")
    # spark.stop()
```

- Precedence Levels of config methods

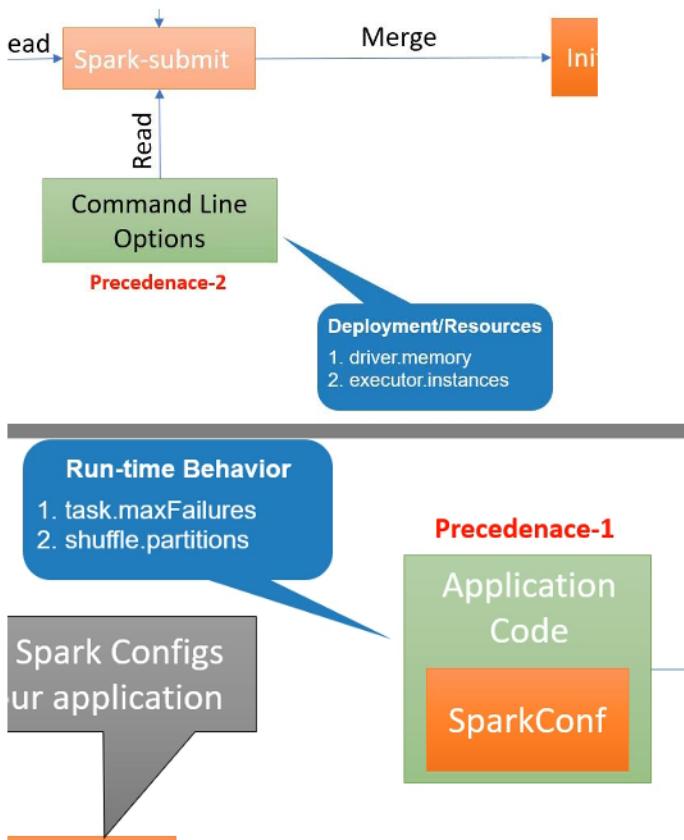


- Rule :

Leave for Cluster Admins



You should be use 3<sup>rd</sup> command line or 4<sup>th</sup> Application config



- **Steps for Application level config**

```

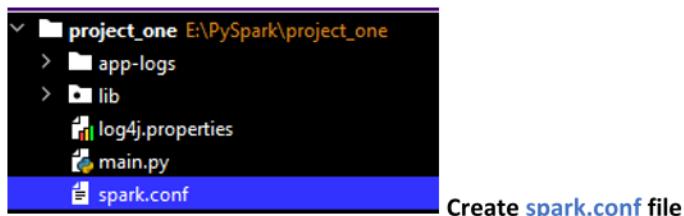
from pyspark.sql import *
from pyspark import SparkConf
from lib.logger import Log4j

if __name__ == "__main__":
    conf = SparkConf()
    conf.set("spark.app.name", "First Spark")
    conf.set("spark.master", "local[3]")
    spark = SparkSession.builder \
        .config(conf=conf) \
        .getOrCreate() \

    logger = Log4j(spark)
    logger.info("Starting First Spark App")
    # Your processing code

    logger.info("Finished First Spark App")
    spark.stop()
  
```

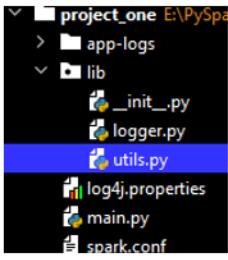
Actually in real world project to minimize complexity and errors go through following



```

[SPARK_APP_CONFIGS]
spark.app.name = Hello Spark
spark.master = local[3]
  
```

**spark.conf**



Create util.py in lib package

```
import configparser
from pyspark import SparkConf

def get_spark_app_config():
    # This function will load the config from spark.conf file and return spark_conf_object
    spark_conf = SparkConf()
    config = configparser.ConfigParser()
    config.read("spark.conf")

    for (key, val) in config.items("SPARK_APP_CONFIGS"):
        spark_conf.set(key, val)
    return spark_conf
```

util.py

```
from pyspark.sql import *
from lib.utils import get_spark_app_config
from lib.logger import Log4j

if __name__ == "__main__":
    conf = get_spark_app_config()
    spark = SparkSession.builder \
        .config(conf=conf) \
        .getOrCreate() \

    logger = Log4j(spark)
    logger.info("Starting First Spark App")
    # Your processing code

    logger.info("Finished First Spark App")
    spark.stop()
```

main.py (here hard coding is fixed)

```
[SPARK_APP_CONFIGS]
spark.app.name = Hello Spark
spark.master = local[3]
app.author = Nishant
```

spark.conf

```
1  from pyspark.sql import *
2  from lib.utils import get_spark_app_config
3  from lib.logger import Log4j
4
5  if __name__ == "__main__":
6      conf = get_spark_app_config()
7      spark = SparkSession.builder \
8          .config(conf=conf) \
9          .getOrCreate() \
10
11     logger = Log4j(spark)
12     logger.info("Starting First Spark App")
13     # Your processing code
14
15     conf_out = spark.sparkContext.getConf()
16     logger.info(conf_out.toDebugString())
17     # To see out come
18
19     logger.info("Finished First Spark App")
20     spark.stop()
```

main.py

```

C:\Users\innk\anaconda3\envs\project_one\python.exe E:\PySpark\project_one\main.py
Warning: Ignoring non-Spark config property: app.author
23/02/09 22:06:35 INFO Hello Spark: Starting First Spark App
23/02/09 22:06:35 INFO Hello Spark: app.author=Nishant
spark.app.id=local-1675960594904
spark.app.name=Hello Spark
spark.app.startTime=1675960592803
spark.app.submitTime=1675960592585
spark.driver.extraJavaOptions=-XX:+IgnoreUnrecognizedVMOptions --add-opens=java.base/sun.nio.ch=BIG_ENDIAN
spark.driver.host=DESKTOP-G5MIVKC
spark.driver.port=52075
spark.executor.extraJavaOptions=-XX:+IgnoreUnrecognizedVMOptions --add-opens=java.base/sun.nio.ch=BIG_ENDIAN
spark.executor.id=driver
spark.master=local[3]

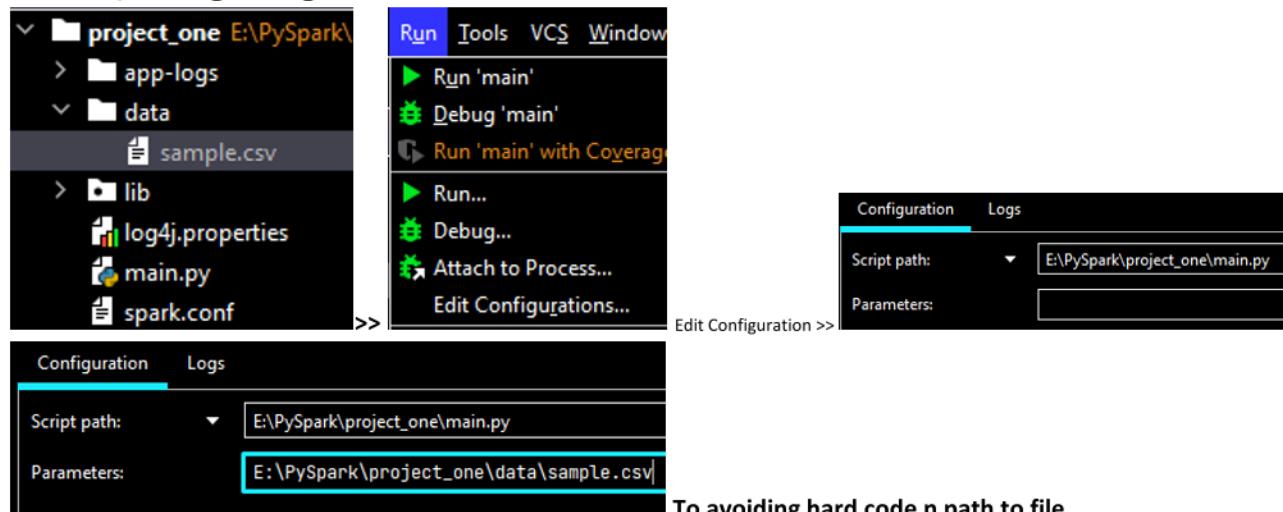
```

O/P

## • Data Frame Introduction

1. Read
2. Processes
3. Write

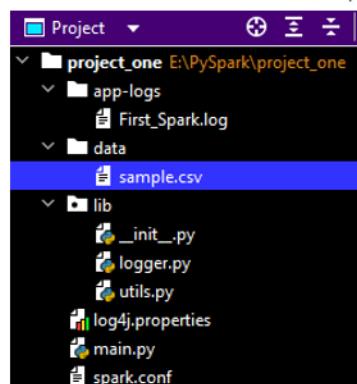
### - Run / Debug Configuration



Timestamp	Age	Gender	Country	state
2014-08-27 11:29:31	37	Female	United States	IL
2014-08-27 11:29:37	44	M	United States	IN
2014-08-27 11:29:44	32	Male	Canada	
2014-08-27 11:29:46	31	Male	United Kingdom	
2014-08-27 11:30:22	31	Male	United States	TX

### Data Frame Schema

1. Column Names
2. Data Types



complete project structure for ref.

### utils.py

```
import configparser
from pyspark import SparkConf

def get_spark_app_config():
    #This fun will load the config from spark.conf file and return spark_conf_object
    spark_conf = SparkConf()
    config = configparser.ConfigParser()
    config.read("spark.conf")

    for (key, val) in config.items("SPARK_APP_CONFIGS"):
        spark_conf.set(key, val)
    return spark_conf

def load_survey_df(spark, data_file):
    return spark.read \
        .option("header", "true") \
        .option("inferSchema", "true") \
        .csv(data_file) \
    # pandas dataframe reader : header, Schema : infraSchema
```

### main.py

```
import sys
from pyspark.sql import *
from lib.utils import get_spark_app_config, load_survey_df
from lib.logger import Log4j

if __name__ == "__main__":
    conf = get_spark_app_config()
    spark = SparkSession.builder \
        .config(conf=conf)\ \
        .getOrCreate() \

    logger = Log4j(spark)

    if len(sys.argv) != 2: # Command line argument if not provided we get error
        logger.error("Usage: HelloSpark <filename>")
        sys.exit(-1)

    logger.info("Starting First Spark App")
    # Your processing code

    survey_df = load_survey_df(spark, sys.argv[1]) # sys.argv[1] file path to csv
    survey_df.show()

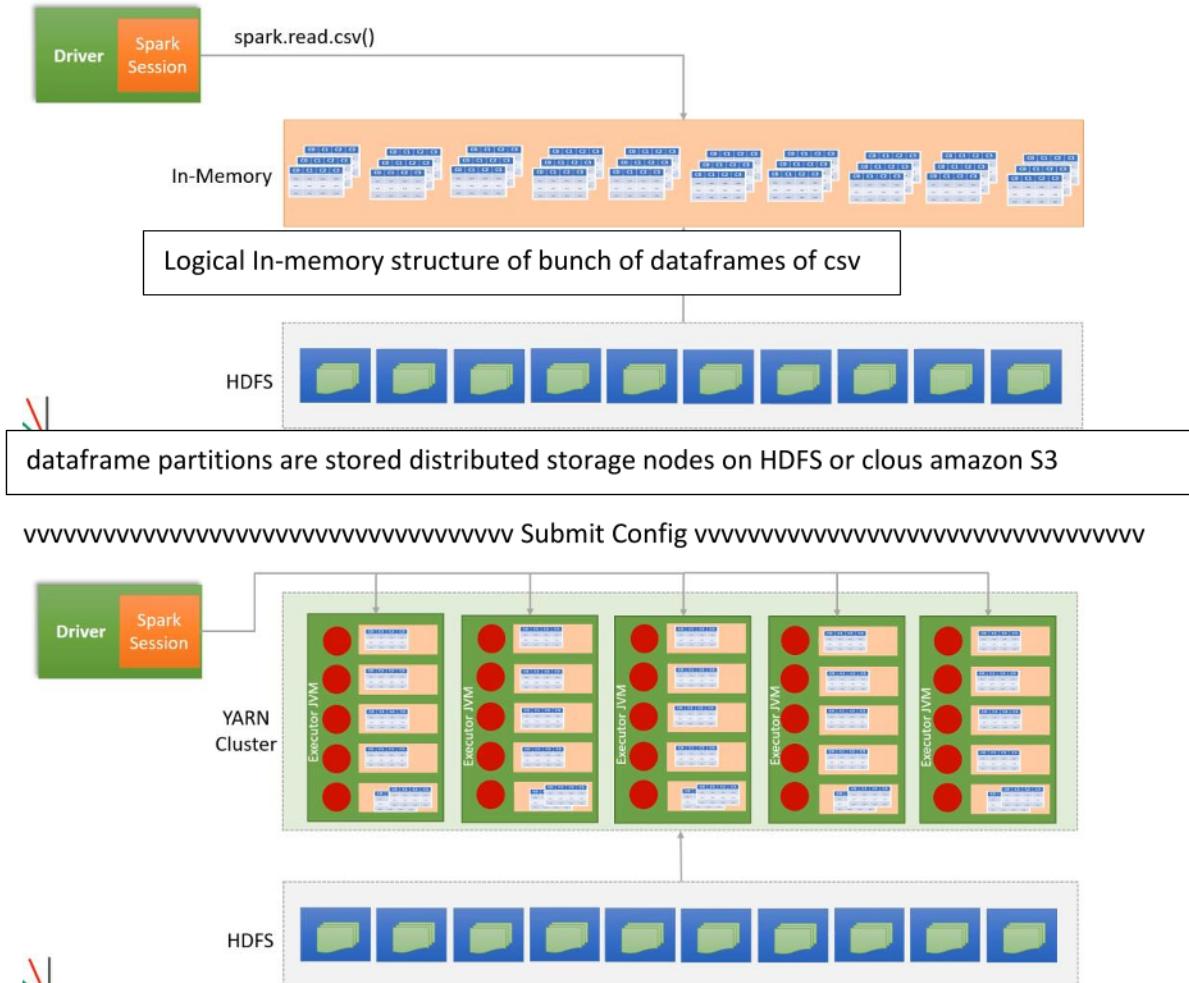
    logger.info("Finished First Spark App")
    spark.stop()
```

### O/P

```
23/02/09 23:59:06 INFO Hello Spark: Starting First Spark App
23/02/09 23:59:16 WARN package: Truncated the string representation of a Row at index 5
+-----+-----+-----+-----+
|      Timestamp|Age|Gender|      Country|state|self_employe
+-----+-----+-----+-----+
|2014-08-27 11:29:31| 37|Female| United States| IL|      N
|2014-08-27 11:29:37| 44|      M| United States| IN|      N
|2014-08-27 11:29:44| 32| Male|      Canada| NA|      N
|2014-08-27 11:29:46| 31| Male|United Kingdom| NA|      N
|2014-08-27 11:30:22| 31| Male| United States| TX|      N
|2014-08-27 11:31:22| 33| Male| United States| TN|      N
|2014-08-27 11:31:50| 35|Female| United States| MI|      N
|2014-08-27 11:32:51| 30|Male| United States| NC|      N
```

- Data Frame Partitions and Executors

- Spark data frame is distributed data structure and how it help to Spark to implement in precessing.

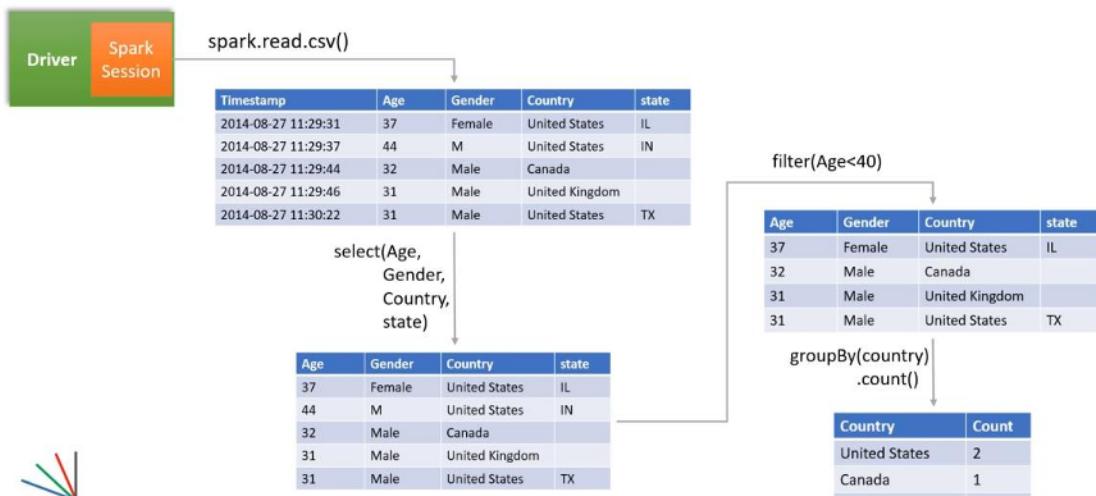


- **Spark Transformations and Action**

- Spark Data Frame is immutable data structure.

# Spark Transformations?

- You can give instruction to driver what do you want to do.
  - Driver decides it how to achieve the instructions by executors these instructions are called as Transformations. These Transformations are same like SQL operations.



## 1. Transformations

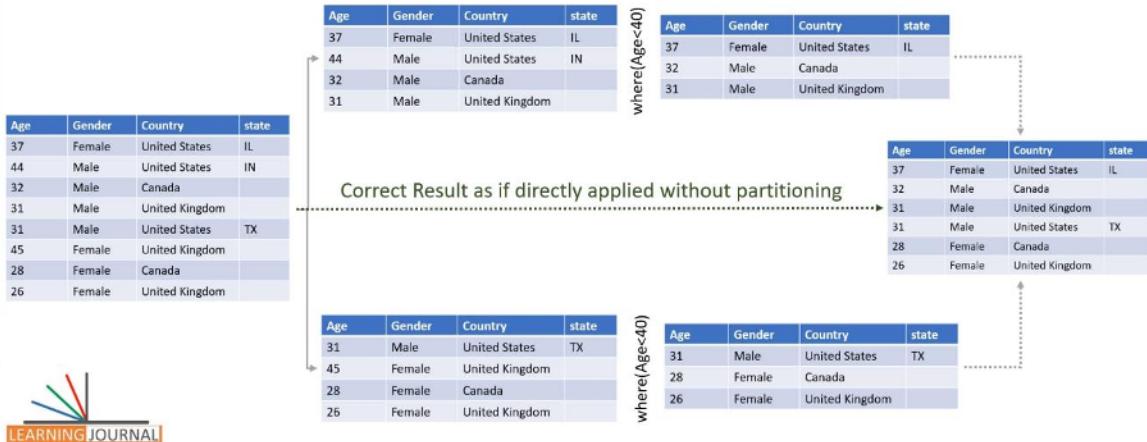
1. Narrow Dependency
  2. Wide Dependency

## Narrow Vs. Wide Transformations?

“Narrow”

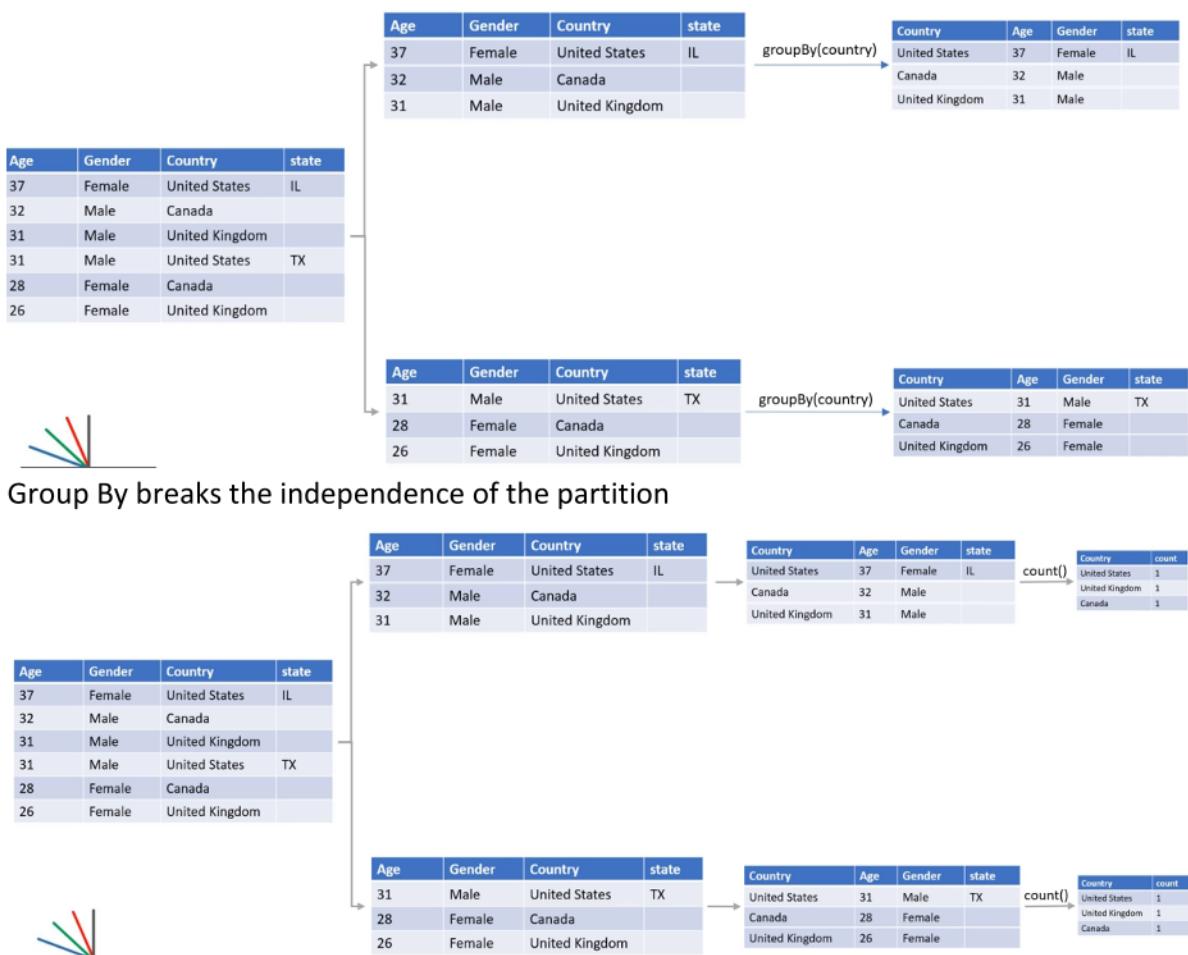
### Narrow Dependency Transformation

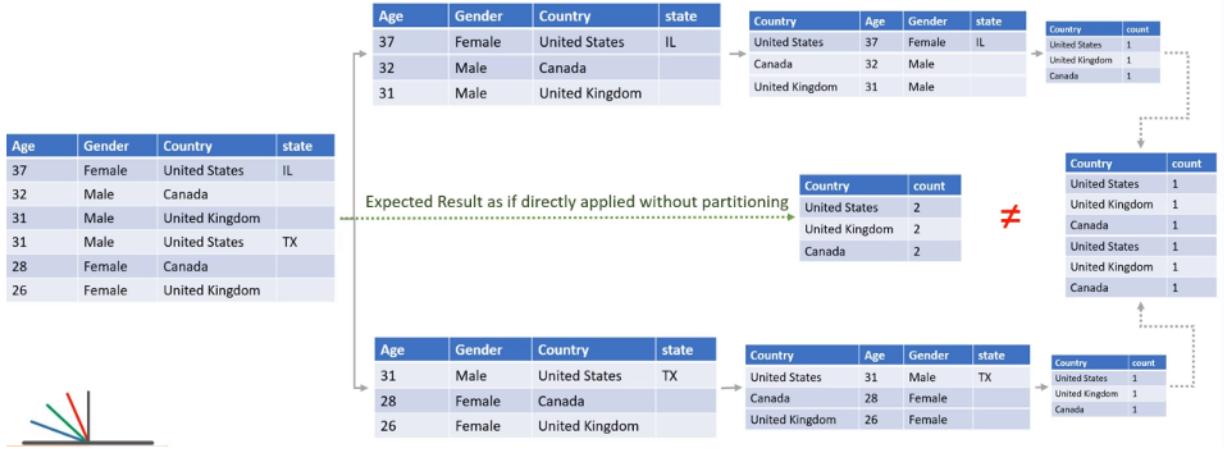
A transformation performed independently on a single partition to produce valid results.



“Wide”

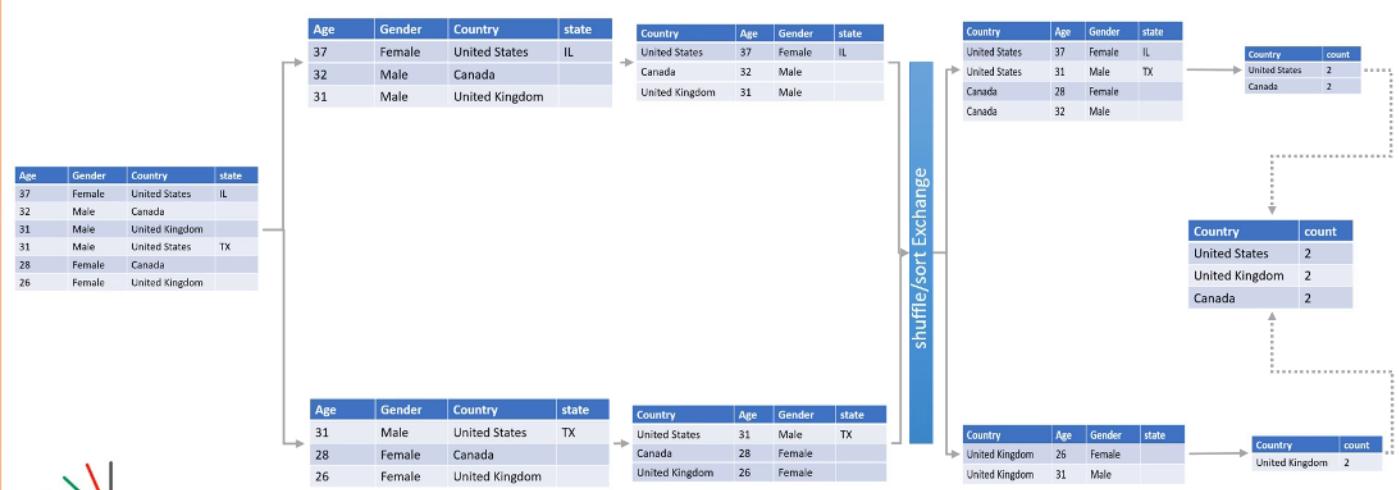
A transformation that requires data from other partitions to produce valid results.





This kind of transformation will create incorrect results how we can correct it? Combine & Re-partition data are known as Shuffle / sort Exchange operations.

Group by, Join, Order By, Distinct, etc. are wide dependency transformations.



## Lazy Evaluations?

Its functional programming technique.

```
spark = SparkSession \
    .builder \
    .config(conf=conf) \
    .getOrCreate()
```

```
survey_df = load_survey_df(spark, sys.argv[1])
filtered_df = survey_df.where("Age < 40")
selected_df = filtered_df.select("Age", "Gender", "Country", "state")
grouped_df = selected_df.groupBy("Country")
count_df = grouped_df.count()
```

```
count_df.show()
```

Typical program runs line by line. But Spark Programs are not behaves same way because in Spark builder pattern is used i.e. DAG of transformations.

```
spark = SparkSession \
    .builder \
    .config(conf=conf) \
    .getOrCreate()
```

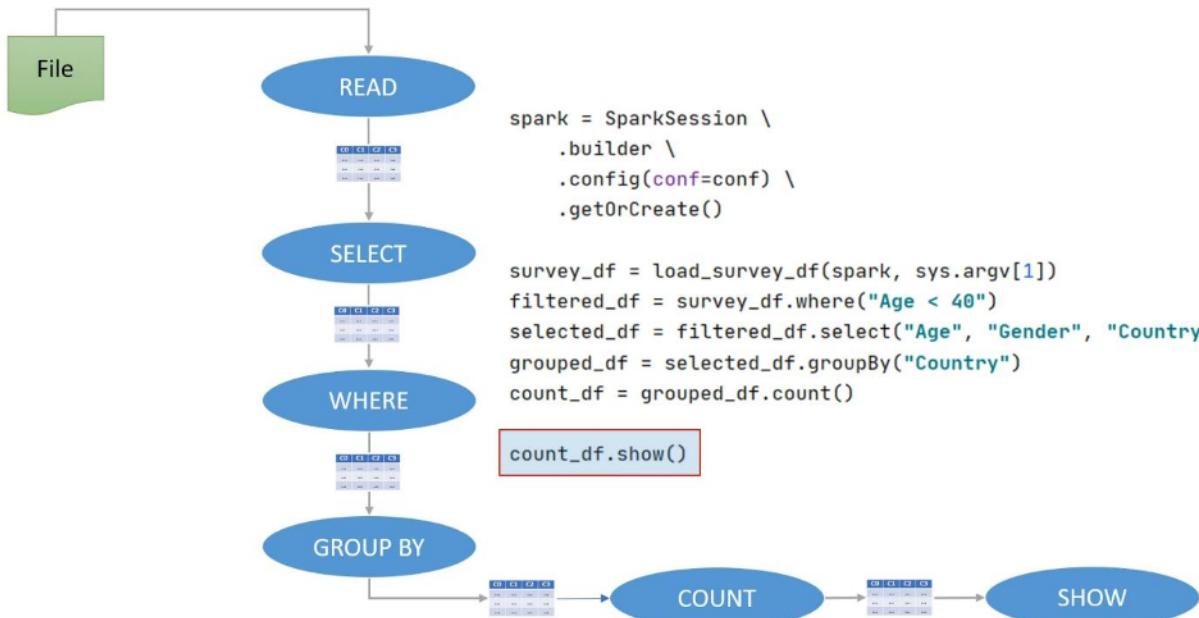
```
survey_df = load_survey_df(spark, sys.argv[1])
filtered_df = survey_df.where("Age < 40")
selected_df = filtered_df.select("Age", "Gender", "Country", "state")
grouped_df = selected_df.groupBy("Country")
count_df = grouped_df.count()

count_df.show()
```

All of these go to the spark driver and rearrange them for optimized certain activities and finally create an execution plan which will be executed by the executors which is terminated by and triggered by Action.

- What's an action? READ, WRITE, COLLECT and SHOW all these are actions.

## Spark Actions



Show( ) is an action (actions terminates the transformation DAG and trigger the Execution)

- Transformations are lazy and actions are evaluated immediately.

```
import sys
from pyspark.sql import *
from lib.utils import get_spark_app_config, load_survey_df
from lib.logger import Log4j

if __name__ == "__main__":
    conf = get_spark_app_config()
    spark = SparkSession.builder \
        .config(conf=conf) \
        .getOrCreate() \

    logger = Log4j(spark)

    if len(sys.argv) != 2: # Command line argument if not provided we get error
        logger.error("Usage: HelloSpark <filename>")
        sys.exit(-1)

    logger.info("Starting First Spark App")
    # Your processing code
```

```

survey_df = load_survey_df(spark, sys.argv[1]) # Schema and DF

count_df = survey_df\
    .where("Age < 40")\
    .select("Age", "Gender", "Country", "state") \
    .groupby("Country") \
    .count() # count_df Change of transformation

count_df.show() # Action

logger.info("Finished First Spark App")
spark.stop()
```

```

```

23/02/10 11:21:53 INFO Hello Spark: Starting First Spark App
+-----+----+
|     Country|count|
+-----+----+
United States	4
Canada	2
United Kingdom	1
+-----+----+
23/02/10 11:22:01 INFO Hello Spark: Finished First Spark App

```

- by creating function

### utils.py

```

- def count_by_country(survey_df):
    return survey_df \
        .where("Age < 40") \
        .select("Age", "Gender", "Country", "state") \
        .groupby("Country") \
        .count()

```

### main.py

```

survey_df = load_survey_df(spark, sys.argv[1]) # Schema and DF

count_df = count_by_country(survey_df)

count_df.show() # Action

```

- **Spark Jobs, Stages and Task ref. 23**

- Controlling Internal operations
- UI - Jobs

localhost:4040/jobs/

APACHE Spark 3.3.1 Jobs Stages Storage Environment Executors SQL / DataFrame Hello Spark application UI

## Spark Jobs (?)

User: innk  
Total Uptime: 3.5 min  
Scheduling Mode: FIFO  
Completed Jobs: 2

Event Timeline

Enable zooming

Executors: Added (blue), Removed (red). Jobs: Succeeded (blue), Failed (red), Running (green). Timeline from 41 to 51. A blue dot marks the event "Executor driver added" at step 43.

csv at 10 February 11:42

Completed Jobs (2)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

| Job Id | Description                                                                      | Submitted           | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|--------|----------------------------------------------------------------------------------|---------------------|----------|-------------------------|-----------------------------------------|
| 1      | csv at NativeMethodAccessorImpl.java:0<br>csv at NativeMethodAccessorImpl.java:0 | 2023/02/10 11:42:50 | 0.3 s    | 1/1                     | 1/1                                     |
| 0      | csv at NativeMethodAccessorImpl.java:0<br>csv at NativeMethodAccessorImpl.java:0 | 2023/02/10 11:42:50 | 0.5 s    | 1/1                     | 1/1                                     |

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

- UI – Stages

localhost:4040/stages/

APACHE Spark 3.3.1 Jobs Stages Storage Environment Executors SQL / DataFrame Hello Spark application UI

## Stages for All Jobs

Completed Stages: 2

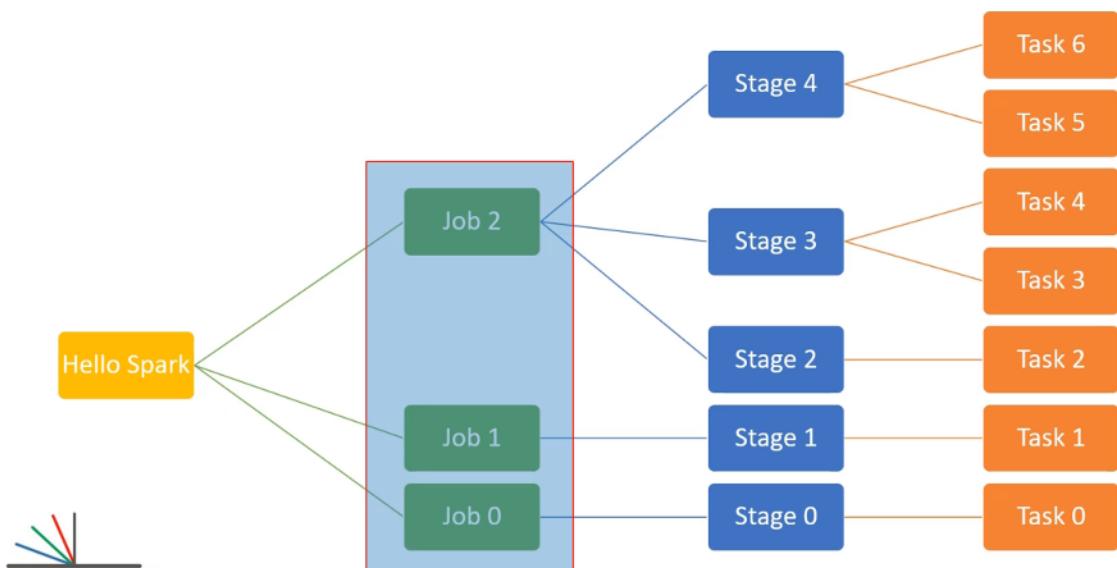
Completed Stages (2)

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

| Stage Id | Description                                        | Submitted           | Duration | Tasks: Succeeded/Total | Input   | Output | Shuffle Read | Shuffle Write |
|----------|----------------------------------------------------|---------------------|----------|------------------------|---------|--------|--------------|---------------|
| 1        | csv at NativeMethodAccessorImpl.java:0<br>+details | 2023/02/10 11:42:50 | 0.2 s    | 1/1                    | 2.2 KiB |        |              |               |
| 0        | csv at NativeMethodAccessorImpl.java:0<br>+details | 2023/02/10 11:42:50 | 0.3 s    | 1/1                    | 2.2 KiB |        |              |               |

Page: 1 1 Pages. Jump to 1 . Show 100 items in a page. Go

## Spark Execution Plan



### main.py

```
import sys
from pyspark.sql import *
from lib.utils import get_spark_app_config, load_survey_df, count_by_country
from lib.logger import Log4j

if __name__ == "__main__":
    conf = get_spark_app_config()
    spark = SparkSession.builder \
        .config(conf=conf) \
        .getOrCreate() \

    logger = Log4j(spark)

    if len(sys.argv) != 2: # Command line argument if not provided we get error
        logger.error("Usage: HelloSpark <filename>")
        sys.exit(-1)

    logger.info("Starting First Spark App")
    # Your processing code

    survey_df = load_survey_df(spark, sys.argv[1]) # Schema and DF

    partitioned_survey_df = survey_df.repartition(2) # It takes DF as I/p and produces
    another DF this new DF should have 2 partitions

    count_df = count_by_country(partitioned_survey_df) # We are using this partitions
    for rest of the transformations.

    input("Press Enter") # for local debugging in UI to hold program to check in UI

    logger.info(count_df.collect()) # Collect Action | The collect action returns the
    DF as Python List
    # show is used for printing, debugging and show method compiles down to a complex
    internal code it creates
    # unnecessary confusion
    logger.info("Finished First Spark App")
    spark.stop()

spark.conf
```

```
[SPARK_APP_CONFIGS]
spark.app.name = Hello Spark
spark.master = local[3]
spark.sql.shuffle.partitions = 2
# its used for the controlling partitions because in internal operations
# we dont know how many partitions are made so we need to control it.
```

## utils.py

```
def count_by_country(survey_df):
    return survey_df \
        .where("Age < 40") \
        .select("Age", "Gender", "Country", "state") \
        .groupby("Country") \
        .count()
```

### • Understanding Execution Plan

- Every Spark action is transmitted into a JOB as follow
- Loading only single csv file created single action that creates single action as follow

```
survey_df = load_survey_df(spark, sys.argv[1]) # now in this code we are just loading
data into DF only.
input("Press Enter")
```



## Spark Jobs (2)

User: innk

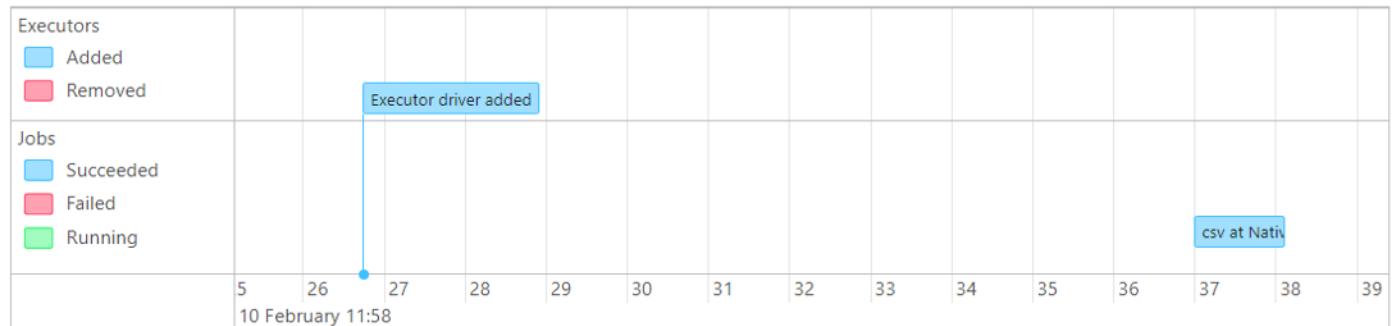
Total Uptime: 59 s

Scheduling Mode: FIFO

Completed Jobs: 1

Event Timeline

Enable zooming



### Completed Jobs (1)

Page: 1

1 Pages. Jump to  . Show  items in a page.

| Job Id | Description                                                                      | Submitted           | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|--------|----------------------------------------------------------------------------------|---------------------|----------|-------------------------|-----------------------------------------|
| 0      | csv at NativeMethodAccessorImpl.java:0<br>csv at NativeMethodAccessorImpl.java:0 | 2023/02/10 11:58:36 | 1 s      | 1/1                     | 1/1                                     |

### - Single Stage



## Stages for All Jobs

Completed Stages: 1

### Completed Stages (1)

Page: 1

1 Pages. Jump to  . Show  items in a page.

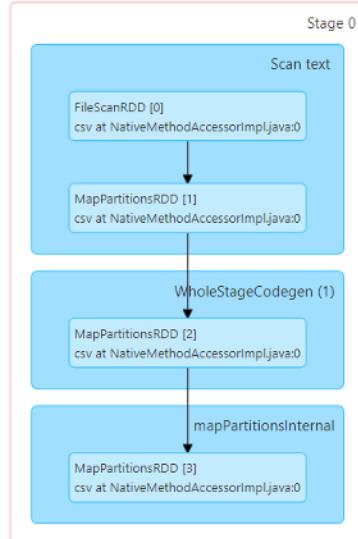
| Stage Id | Description                                        | Submitted           | Duration | Tasks: Succeeded/Total | Input   | Output | Shuffle Read | Shuffle Write |
|----------|----------------------------------------------------|---------------------|----------|------------------------|---------|--------|--------------|---------------|
| 0        | csv at NativeMethodAccessorImpl.java:0<br>+details | 2023/02/10 11:58:37 | 0.4 s    | 1/1                    | 2.2 KiB |        |              |               |

To see DAG Visualization, click on Description

## Details for Stage 0 (Attempt 0)

Resource Profile Id: 0  
 Total Time Across All Tasks: 0.2 s  
 Locality Level Summary: Process local: 1  
 Input Size / Records: 2.2 KIB / 1  
 Associated Job Ids: 0

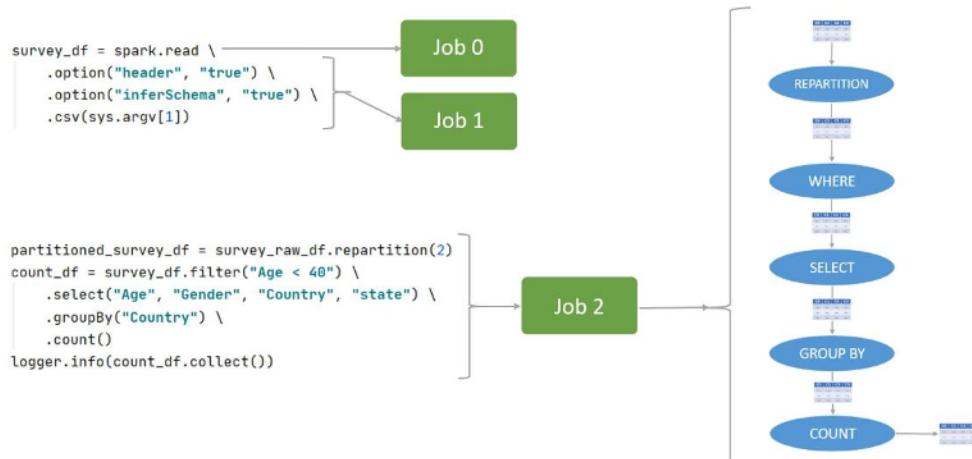
► DAG Visualization



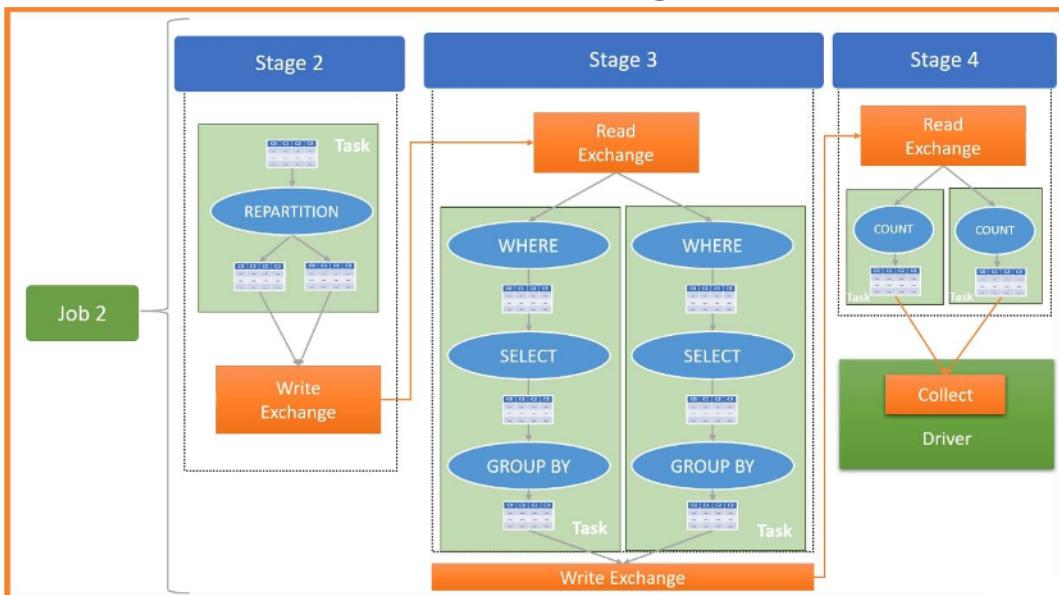
▶ Show Additional Metrics  
 ▶ Event Timeline

- DAG shows internal processes the sequence is compiled code which is generated by the spark.

## Spark Execution Plan



## Shuffle/sort Exchange

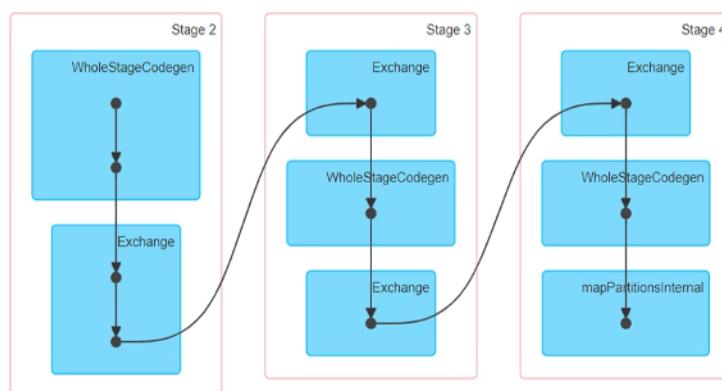


## Details for Job 2

Status: SUCCEEDED

Completed Stages: 3

- ▶ Event Timeline
- ▼ DAG Visualization



DAG 3 stages

```

23/02/10 12:30:01 INFO Hello Spark: Starting First Spark App
23/02/10 12:30:12 INFO Hello Spark: [[Canada, 2], [United States, 4], [United Kingdom, 1]]
23/02/10 12:34:16 INFO Hello Spark: Finished First Spark App

```

Log with count list

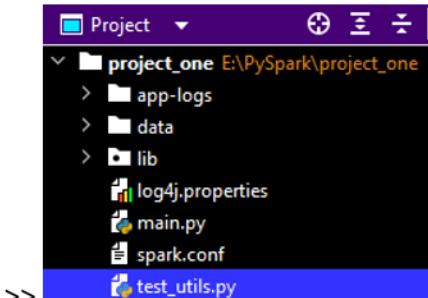
### • Unit Testing Spark Application

**Test Case:** Make sure that we are loading the data file correctly.

| S. No. | Test Step                      | Test Data        |
|--------|--------------------------------|------------------|
| 1.     | Read the data file             | sample.csv       |
| 2.     | Validate the number of records | record count = 9 |

**Test Case:** Record Count by country is computed correctly

| S. No. | Test Step                 | Test Data                                             |
|--------|---------------------------|-------------------------------------------------------|
| 1.     | Read the data file        | sample.csv                                            |
| 2.     | Validate count by country | United Kingdom = 1<br>Canada = 2<br>United States = 4 |



### test\_utils.py

```

from unittest import TestCase
from pyspark.sql import SparkSession
from lib.utils import load_survey_df, count_by_country

class UtilsTestCase(TestCase):

    @classmethod
    def setUpClass(cls) -> None:
        cls.spark = SparkSession.builder \
            .master("local[3]") \
            .appName("SparkTest") \
            .getOrCreate()

    def test_datafile_loading(self):
        sample_df = load_survey_df(self.spark, "data/sample.csv")
        result_count = sample_df.count()
        self.assertEqual(result_count, 9, "Record count should be 9")

    def test_country_count(self):
        sample_df = load_survey_df(self.spark, "data/sample.csv")
        count_list = count_by_country(sample_df).collect()
        count_dict = dict()
        for row in count_list:
            count_dict[row["Country"]] = row["count"]
        self.assertEqual(count_dict["United States"], 4, "Count for United State should be 4")
        self.assertEqual(count_dict["Canada"], 2, "Count for United State should be 2")
        self.assertEqual(count_dict["United Kingdom"], 1, "Count for United State should be 1")

    # @classmethod
    # def tearDownClass(cls) -> None:
    #     cls.spark.stop()

```

```
Run: Python tests for test_utils_UtilsTestCase X

▶ ✓ ⚡ ⚡ | ↴ ⌂ | ⌂ ⌂ | > ✓ Tests passed: 2 of 2 tests - 6 sec 817 ms

Test Results 6 sec 817 ms C:\Users\innk\anaconda3\envs\project_one\python.exe "C:/Program Files/JetBrains/PyCharm Community Edition 2020.2.3\helpers\pycharm\test_runner.py" --no-header --no-summary -q in E:\PyCharm\Project\test_utils.py

Testing started at 12:57 pm ...
Launching pytest with arguments test_utils.py::UtilsTestCase --no-header --no-summary -q in E:\PyCharm\Project\test_utils.py

=====
test session starts =====
collecting ... collected 2 items

test_utils.py::UtilsTestCase::test_country_count
test_utils.py::UtilsTestCase::test_datafile_loading

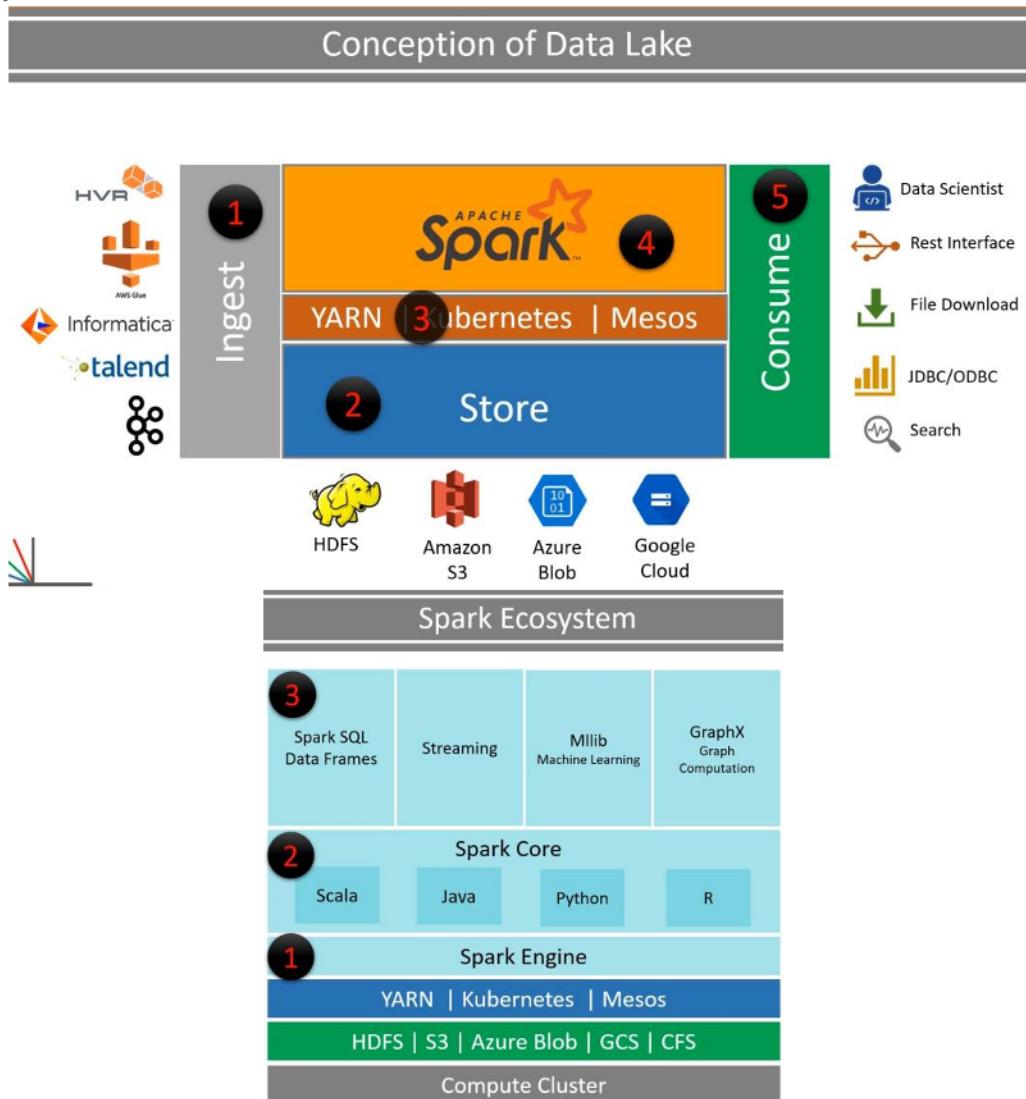
=====
2 passed in 13.14s =====

Process finished with exit code 0
PASSED [ 50%]{'United States': 4, 'Canada': 2, 'United Kingdom': 1}
PASSED [100%]
```

```
+-----+---+
|      Country|count|
+-----+---+
United States	4
Canada	2
United Kingdom	1
+-----+---+
```

### count table for ref.

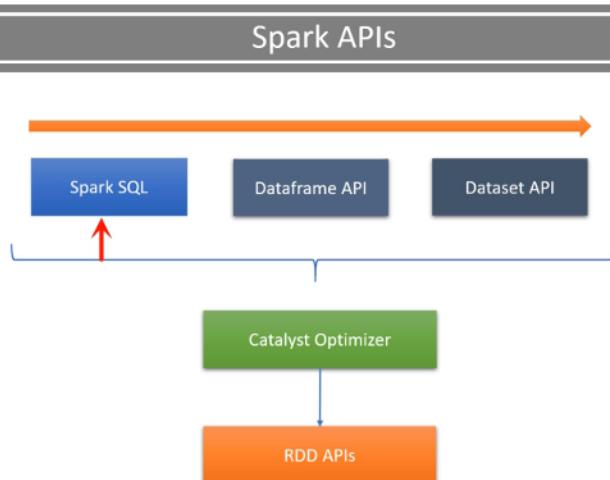
- Summary



### **Developer Experience**

1. Creating and Configuring Spark Project using your IDE
2. Configuring Log4J for your Spark Application
3. Creating and Configuring Spark Session
4. Managing your Spark Session Configurations using spark.conf
5. Creating a modular Structure for your Spark Application
6. Unit Testing Spark Application
7. Building and packaging your Spark Application
8. Deploying your Spark Application on a Cluster
9. Collecting Application Logs from Spark Cluster

- **Spark Structured API Foundation**



- Spark SQL: Most convenient to use.

- DF API\*\*: lesser than Spark SQL. Any favourite language is suitable.

- Dataset API: Lesser than above both. Language native such as Scala, java (which are JVM based languages)

- Catalyst Optimizer: SparkSQL, DF API, Dataset API. Optimizer decides how the code is executed and execution plan

- RDD not recommended by spark

- **Spark RDD APIs (Resilient Distributed Dataset) ref. 36**

### How to Create RDD?

#### Spark RDD APIs

##### Creating Data frame

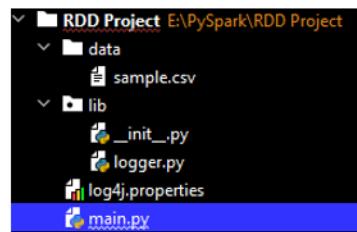
1. Create a SparkConf object
2. Create a Spark Session using your SparkConf
3. Use the spark session to read the data file

```

conf = SparkConf()
conf.setMaster("local[3]").setAppName("HelloRDD")

spark = SparkSession.builder.config(conf=conf).getOrCreate()

surveyDF = spark.read.option("header", "true").option("inferSchema", "true")
  
```



##### main.py

```

import sys
from collections import namedtuple
from pyspark import SparkConf, SparkContext
from pyspark.sql import SparkSession
from lib.logger import Log4j

SurveyRecord = namedtuple("SurveyRecord", ["Age", "Gender", "Country", "State"])

if __name__ == "__main__":
    conf = SparkConf() \
        .setMaster("local[3]") \
        .setAppName("RDDProject")

    # sc = SparkContext(conf=conf)
    # alternet method to get spark context

    spark = SparkSession \
        .builder \
        .config(conf=conf) \
        .getOrCreate()

    sc = spark.sparkContext
    logger = Log4j(spark)

    if len(sys.argv) != 2:
  
```

```

logger.error("Usage: HelloSpark <filename>")
sys.exit(-1)

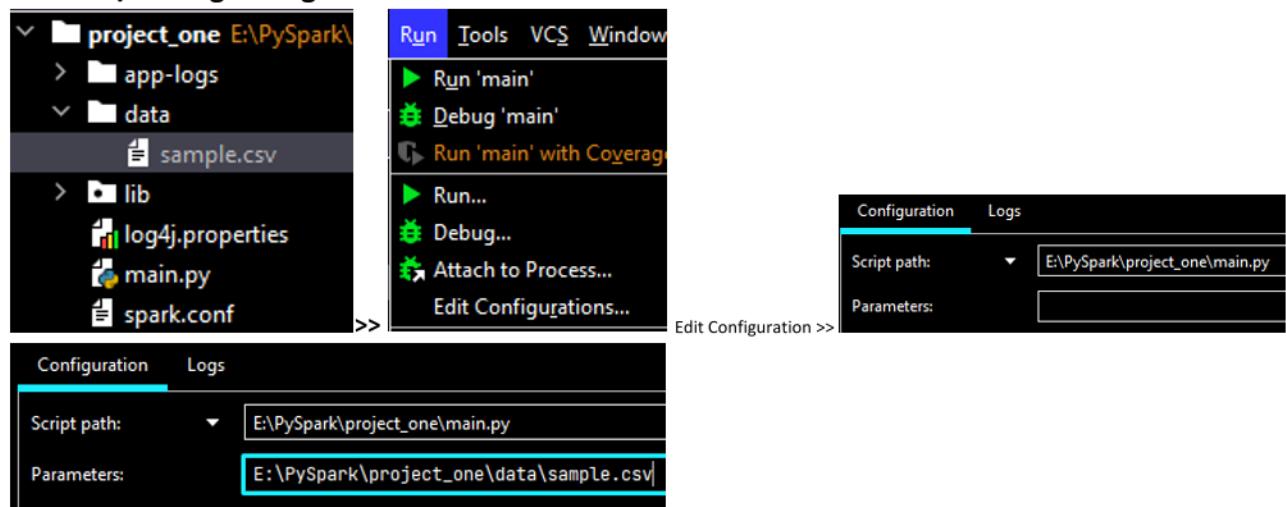
lineRDD = sc.textFile(sys.argv[1])
partitionedRDD = lineRDD.repartition(2)

colsRDD = partitionedRDD.map(lambda line: line.replace('"', '').split(","))
# removing " from csv then o/p list of text.
selectRDD = colsRDD.map(lambda cols: SurveyRecord(int(cols[1]), cols[2],
cols[3], cols[4])) # Attaching schema to RDD with help of SurveyRecord
filteredRDD = selectRDD.filter(lambda r: r.Age < 40)
kvRDD = filteredRDD.map(lambda r: (r.Country, 1)) # kv : key value where
Country becomes key and 1 hardcoded value
countRDD = kvRDD.reduceByKey(lambda v1, v2: v1 + v2)

colsList = countRDD.collect()
for x in colsList:
    logger.info(x)

```

### - Run / Debug Configuration



\*\*\*csv file path if this parameter is not config then code not able to find csv with sys.argv[1]\*\*\*

```

"C:\Users\innk\anaconda3\envs\RDD_Project\python.exe"
data/sample.csv
E:\pySpark_soft\spark_3\python\lib\pyspark.zip\pyspar
E:\pySpark_soft\spark_3\python\lib\pyspark.zip\pyspar
E:\pySpark_soft\spark_3\python\lib\pyspark.zip\pyspar
E:\pySpark_soft\spark_3\python\lib\pyspark.zip\pyspar
23/02/13 13:41:40 INFO HelloRDD: [United States, 4]
23/02/13 13:41:40 INFO HelloRDD: [Canada, 2]
23/02/13 13:41:40 INFO HelloRDD: [United Kingdom, 1]

```

- Working with Spark SQL

### main.py

```
import sys

from pyspark.sql import SparkSession

from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession \
        .builder \
        .master("local[3]") \
        .appName("SparkSQL") \
        .getOrCreate()

    logger = Log4j(spark)

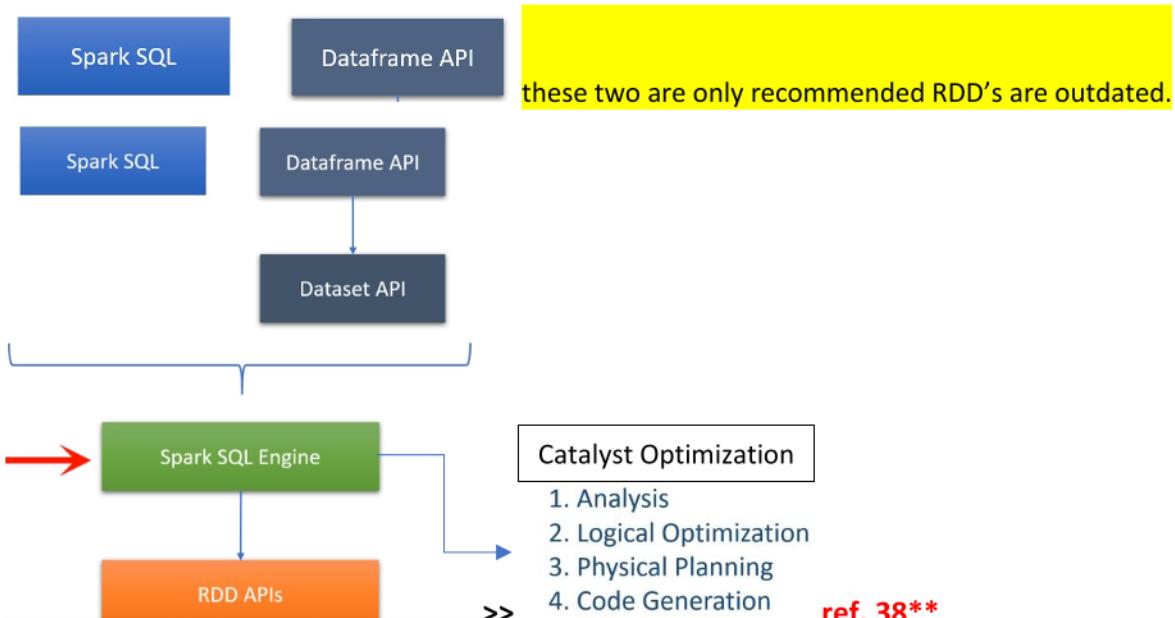
    if len(sys.argv) != 2:
        logger.error("Usage: HelloRDD <filename>")
        sys.exit(-1)

    surveyDF = spark.read \
        .option("header", "true") \
        .option("inferSchema", "true") \
        .csv(sys.argv[1])

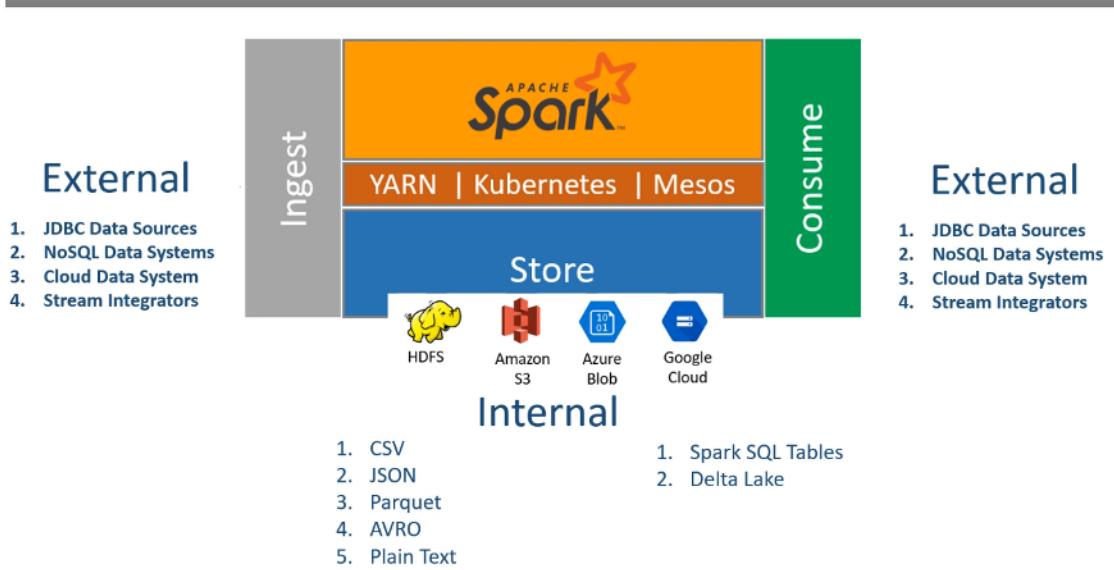
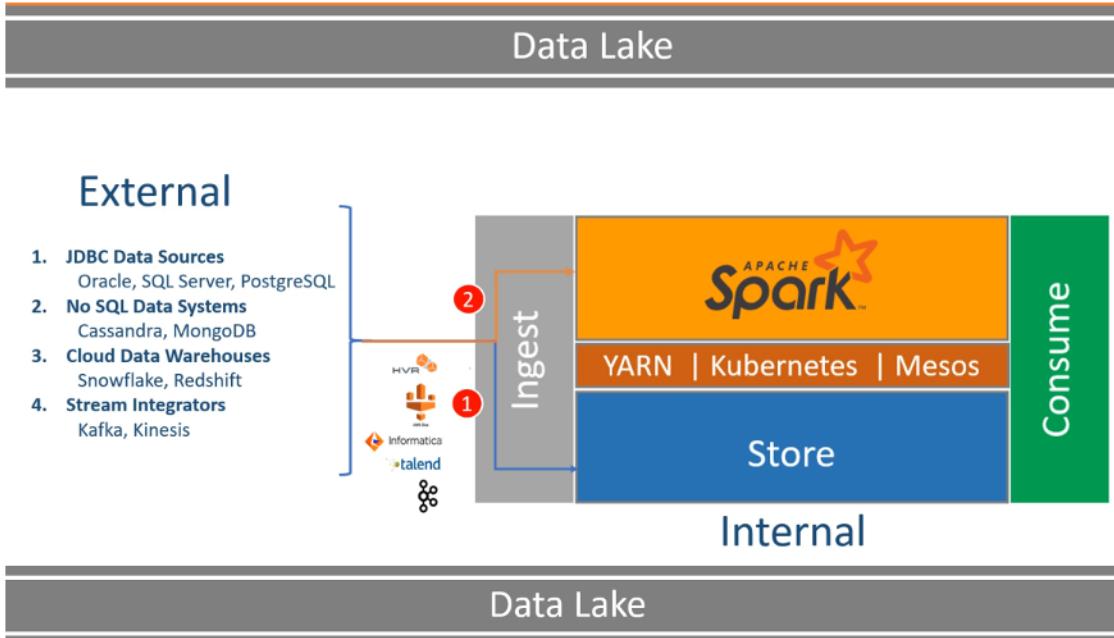
    # Spark allows you to register data frame as a view
    surveyDF.createOrReplaceTempView("survey_tbl") #View
    countDF = spark.sql("SELECT Country, COUNT(1) AS COUNT FROM survey_tbl WHERE Age<40
GROUP BY Country") # SQL code
    countDF.show()
```

| Country        | COUNT |
|----------------|-------|
| United States  | 4     |
| Canada         | 2     |
| United Kingdom | 1     |

- Spark SQL Engine and Catalyst Optimizer



- **Spark Data Sources and Sinks**



- **Spark Data Frame Reader API**



DataFrameReader API : <https://spark.apache.org/docs/latest/api/python/pyspark.sql.html?highlight=dataframereader#pyspark.sql.DataFrameReader>

#### General structure

```

DataFrameReader
  .format("...")
  .option("key", "value")
  .schema("...")
  .load()
  
```

#### Indicative Example

```

→ spark.read
  .format("csv")
  .option("header", "true")
  .option("path", "/data/mycsvfiles/")
  .option("mode", "FAILFAST")
  .schema(mySchema)
  .load()
  
```

### Indicative Example

```
spark.read  
    .format("csv")  
    .option("header", "true")  
    .option("path", "/data/mycsvfiles/")  
    .option("mode", "FAILFAST")  
    .schema(mySchema)  
    .load()
```

Community Formats  
Cassandra, MongoDB, AVRO, XML,  
HBase, Redshift

```
.schema(mySchema)  
.load()
```

Schema  
1. Explicit  
2. Infer Schema  
3. Implicit

### Indicative Example

```
spark.read  
    .format("csv")  
    .option("header", "true")  
    .option("path", "/data/mycsvfiles/")  
    .option("mode", "FAILFAST")  
    .schema(mySchema)  
    .load()
```

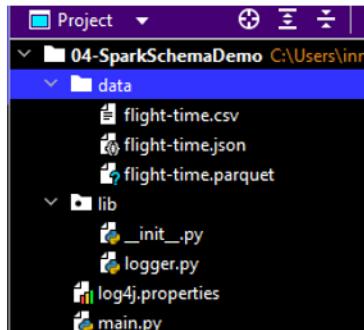
Read Mode  
1. PERMISSIVE  
2. DROPMALFORMED  
3. FAILFAST

1. **Permissive** mode is default option. This set all the field to the null when it encounters a corrupted record. Places that record is string column called `_corrupt_record`
2. **DropMalformed** mode remove the malformed record. Ignoring the malformed records only loading well formed records only.
3. **FAILFAST** raises and exceptions and terminate when detection of mal formed record.

Its standard code. Don't use shortcuts.

- **Reading CSV, JSON, and Parquet files.**

#### 1. How to use DataFrameReader for CSV, JSON, and Parquet



#### main.py (without schema)

```
from pyspark.sql import SparkSession  
  
from lib.logger import Log4j  
  
if __name__ == "__main__":  
    spark = SparkSession \  
        .builder \  
        .master("local[3]") \  
        .appName("spark_schema_demo") \  
        .getOrCreate()  
  
    logger = Log4j(spark)  
  
    flight_time_csv = spark.read \  
        .format("csv") \  
        .option("header", "true") \  
        .load("data/flight*.csv")  
  
    flight_time_csv.show(5)  
    logger.info("CSV Schema: " + flight_time_csv.schema.simpleString())
```

```

+-----+-----+-----+-----+
| FL_DATE|OP_CARRIER|OP_CARRIER_FL_NUM|ORIGIN|ORIGI
+-----+-----+-----+-----+
1/1/2000	DL	1451	BOS
1/1/2000	DL	1479	BOS
1/1/2000	DL	1857	BOS
1/1/2000	DL	1997	BOS
1/1/2000	DL	2065	BOS
+-----+-----+-----+-----+
only showing top 5 rows

```

### Result (int result as string here)

INFO spark\_schema\_demo: CSV Schema:

```
struct<FL_DATE:string,OP_CARRIER:string,OP_CARRIER_FL_NUM:string,ORIGIN:string,ORIGIN_CITY_NAME:string,DEST:string,DEST_CITY_NAME:string,CRS_DEP_TIME:string,DEP_TIME:string,WHEELS_ON:string,TAXI_IN:string,CRS_ARR_TIME:string,ARR_TIME:string,CANCELLED:string,DISTANCE:string>
```

**main.py (with schema)** in this case all numerical are result as integers.

```

flight_time_csv = spark.read \
    .format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("data/flight*.csv")

```

### Just adding Schema

INFO spark\_schema\_demo: CSV Schema:

```
struct<FL_DATE:string,OP_CARRIER:string,OP_CARRIER_FL_NUM:int,ORIGIN:string,ORIGIN_CITY_NAME:string,DEST:string,DEST_CITY_NAME:string,CRS_DEP_TIME:int,DEP_TIME:int,WHEELS_ON:int,TAXI_IN:int,CRS_ARR_TIME:int,ARR_TIME:int,CANCELLED:int,DISTANCE:int>
```

- You can't just be relied on inferSchema only. Because in result we got Date as string in above output.

## Schema for DataFrame

1. Explicit
- 2. Implicit

**main.py (for JSON)**

```

flight_time_json_df = spark.read \
    .format("json") \
    .load("data/flight*.json")

flight_time_json_df.show(5)
logger.info("JSON Schema: " + flight_time_json_df.schema.simpleString())

```

```

+-----+-----+-----+-----+-----+-----+-----+-----+
|ARR_TIME|CANCELLED|CRS_ARR_TIME|CRS_DEP_TIME|DEP_TIME|DEST|DEST_CITY_NAME|DISTANCE| FL_D
+-----+-----+-----+-----+-----+-----+-----+-----+
|  1348|      0|       1400|      1115|    1113|  ATL| Atlanta, GA|     946|1/1/2
|  1543|      0|       1559|      1315|    1311|  ATL| Atlanta, GA|     946|1/1/2
|  1651|      0|       1721|      1415|    1414|  ATL| Atlanta, GA|     946|1/1/2
|  2005|      0|       2013|      1715|    1720|  ATL| Atlanta, GA|     946|1/1/2
|  2240|      0|       2300|      2015|    2010|  ATL| Atlanta, GA|     946|1/1/2
+-----+-----+-----+-----+-----+-----+-----+-----+
only showing top 5 rows

```

23/02/26 13:48:23 INFO spark\_schema\_demo: JSON Schema: struct<ARR\_TIME:bigint,CANCELLED:bigint,CRS\_ARR\_TIME:bigint,CRS\_DEP\_TIME:bigint,DEP\_TIME:bigint,DEST:string,DEST\_CITY\_NAME:string,DISTANCE:bigint,FL\_DATE:timestamp> DF results in alphabetical format.

```
INFO spark_schema_demo: JSON Schema:  
struct<ARR_TIME:bigint,CANCELLED:bigint,CRS_ARR_TIME:bigint,CRS_DEP_TIME:bigint,DEP_TIME:bigint,DEST:string,DEST_CITY_NAME:string,DISTANCE:bigint,FL_DATE:string,OP_CARRIER:string,OP_CARRIER_FL_NUM:bigint,ORIGIN:string,ORIGIN_CITY_NAME:string,TAXI_IN:bigint,WHEELS_ON:bigint>
```

Date still in string.

## Schema for DataFrame

1. Explicit
2. Implicit

[main.py](#) (for parquet format)

```
flight_time_parquet_df = spark.read \  
.format("parquet") \  
.load("data/flight*.parquet")  
  
flight_time_parquet_df.show(5)  
logger.info("PARQUET Schema: " + flight_time_parquet_df.schema.simpleString())
```

```
+-----+-----+-----+-----+-----+-----+  
| FL_DATE|OP_CARRIER|OP_CARRIER_FL_NUM|ORIGIN|ORIGIN_CITY_NAME|DEST|DEST_CITY_NAME|CRS_DEP_TIME|  
+-----+-----+-----+-----+-----+-----+-----+  
2000-01-01	DL	1451	BOS	Boston, MA	ATL	Atlanta, GA	1115
2000-01-01	DL	1479	BOS	Boston, MA	ATL	Atlanta, GA	1315
2000-01-01	DL	1857	BOS	Boston, MA	ATL	Atlanta, GA	1415
2000-01-01	DL	1997	BOS	Boston, MA	ATL	Atlanta, GA	1715
2000-01-01	DL	2065	BOS	Boston, MA	ATL	Atlanta, GA	2015
+-----+-----+-----+-----+-----+-----+-----+  
only showing top 5 rows  
  
23/02/26 13:54:08 INFO spark_schema_demo: PARQUET Schema: struct<FL_DATE:date,OP_CARRIER:string,OP_
```

```
INFO spark_schema_demo: PARQUET Schema:
```

```
struct<FL_DATE:date,OP_CARRIER:string,OP_CARRIER_FL_NUM:int,ORIGIN:string,ORIGIN_CITY_NAME:string,DEST:string,DEST_CITY_NAME:string,CRS_DEP_TIME:int,DEP_TIME:int,WHEELS_ON:int,TAXI_IN:int,CRS_ARR_TIME:int,ARR_TIME:int,CANCELLED:int,DISTANCE:int>
```

- Parquet file format is recommended for spark.

- **Creating Spark DataFrame and Dataset Transformations**

- As we see in previous example the schema is not working with csv and json perfectly.
- Now we are **explicitly** setting schema for DF with spark data types.

## Spark Data Types : <https://spark.apache.org/docs/latest/api/python/pyspark.sql.html#module-pyspark.sql.types>

| S.No. | Spark Types   | Scala Types          | Python Types          |
|-------|---------------|----------------------|-----------------------|
| 1.    | IntegerType   | Int                  | Int                   |
| 2.    | LongType      | Long                 | long                  |
| 3.    | FloatType     | Float                | Float                 |
| 4.    | DoubleType    | Double               | Float                 |
| 5.    | StringType    | String               | String                |
| 6.    | DateType      | java.sql.Date        | datetime.date         |
| 8.    | TimestampType | java.sql.Timestamp   | datetime.datetime     |
| 9.    | ArrayType     | scala.collection.Seq | List, tuple, or array |
| 10.   | MapType       | scala.collection.Map | dict                  |

- Spark is work as compiler it complies with the high-level API code into low level RDD operations. During this compilation process it generate different execution plans and also perform bunch of optimizations it's not possible to spark engine to without maintaining its type information this kind of thing we seen in SQL every column contains a data type.

## How to define Schema?

### Spark Schema

1. Programmatically
2. Using DDL String

## 1. main.py Programmatically (used for CSV also used for other formats)

```
from pyspark.sql import SparkSession
from pyspark.sql.types import StructField, StringType, IntegerType, DateType, StructType

from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession \
        .builder \
        .master("local[3]") \
        .appName("spark_schema_demo") \
        .getOrCreate()

    logger = Log4j(spark)

    flightSchemaStruct = StructType([ # StructType:: DF row structure
        StructField("FL_DATE", DateType()), # StructField:: Column definition
        StructField("OP_CARRIER", StringType()),
        StructField("OP_CARRIER_FL_NUM", IntegerType()),
        StructField("ORIGIN", StringType()),
        StructField("ORIGIN_CITY_NAME", StringType()),
        StructField("DEST", StringType()),
        StructField("DEST_CITY_NAME", StringType()),
        StructField("CRS_DEP_TIME", IntegerType()),
        StructField("DEP_TIME", IntegerType()),
        StructField("WHEELS_ON", IntegerType()),
        StructField("TAXI_IN", IntegerType()),
        StructField("CRS_ARR_TIME", IntegerType()),
        StructField("ARR_TIME", IntegerType()),
        StructField("CANCELLED", IntegerType()),
        StructField("DISTANCE", IntegerType())
    ])

    flight_time_csv_df = spark.read \
        .format("csv") \
        .option("header", "true") \
        .schema(flightSchemaStruct) \
        .option("mode", "FAILFAST") \
        .option("dateFormat", "M/d/y") \
        .load("data/flight*.csv")
    # .option("inferSchema", "true")\ here we provided data type schema
    # .option("mode", "FAILFAST") is used to prevent error due to wrong data type we provided.
    # csv needs date format pattern

    flight_time_csv_df.show(5)
    logger.info("CSV Schema: " + flight_time_csv_df.schema.simpleString())

    flight_time_json_df = spark.read \
        .format("json") \
        .load("data/flight*.json")

    flight_time_json_df.show(5)
    logger.info("JSON Schema: " + flight_time_json_df.schema.simpleString())

    flight_time_parquet_df = spark.read \
        .format("parquet") \
        .load("data/flight*.parquet")

    flight_time_parquet_df.show(5)
    logger.info("PARQUET Schema: " + flight_time_parquet_df.schema.simpleString())
```

## Result

```
+---+---+---+---+---+  
| FL_DATE|OP_CARRIER|OP_CARRIER_FL_NUM|ORIGIN|ORIGIN_CITY_NAME|DEST|DEST  
+---+---+---+---+---+  
2000-01-01	DL	1451	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	1479	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	1857	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	1997	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	2065	BOS	Boston, MA	ATL	Atlanta, GA
+---+---+---+---+---+  
only showing top 5 rows
```

23/02/26 14:23:39 INFO spark\_schema\_demo: CSV Schema: struct<FL\_DATE:date,  
After this data type schema we  
got correct data type for csv and JSON.

## 2. DDL String for Schema

```
flightSchemaDDL = """FL_DATE DATE, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED INT, DISTANCE INT"""  
|
```

Here we are using for JSON | Also used in other formats

```
from pyspark.sql import SparkSession  
  
from lib.logger import Log4j  
  
if __name__ == "__main__":  
    spark = SparkSession \  
        .builder \  
        .master("local[3]") \  
        .appName("spark_schema_demo") \  
        .getOrCreate()  
  
    logger = Log4j(spark)  
  
    flightSchemaDDL = """FL_DATE DATE, OP_CARRIER STRING, OP_CARRIER_FL_NUM INT, ORIGIN STRING,  
    ORIGIN_CITY_NAME STRING, DEST STRING, DEST_CITY_NAME STRING, CRS_DEP_TIME INT, DEP_TIME INT,  
    WHEELS_ON INT, TAXI_IN INT, CRS_ARR_TIME INT, ARR_TIME INT, CANCELLED INT, DISTANCE INT"""  
  
    flight_time_json_df = spark.read \  
        .format("json") \  
        .schema(flightSchemaDDL) \  
        .option("dateFormat", "M/d/y") \  
        .load("data/flight*.json")  
  
    flight_time_json_df.show(5)  
    logger.info("JSON Schema: " + flight_time_json_df.schema.simpleString())
```

```
+---+---+---+---+---+  
| FL_DATE|OP_CARRIER|OP_CARRIER_FL_NUM|ORIGIN|ORIGIN_CITY_NAME|DEST|DEST_CITY_NAME|CRS_DEP_  
+---+---+---+---+---+  
2000-01-01	DL	1451	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	1479	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	1857	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	1997	BOS	Boston, MA	ATL	Atlanta, GA
2000-01-01	DL	2065	BOS	Boston, MA	ATL	Atlanta, GA
+---+---+---+---+---+  
only showing top 5 rows
```

23/02/26 14:35:51 INFO spark\_schema\_demo: JSON Schema: struct<FL\_DATE:date, OP\_CARRIER:string,

## • Spark DF Writer API

DataFrameWriter API : <https://spark.apache.org/docs/latest/api/python/pyspark.sql.html#pyspark.sql.DataFrameWriter>

### General structure

```
DataFrameWriter
    .format(...)
    .option(...)
    .partitionBy(...)
    .bucketBy(...)
    .sortBy(...)
    .save()
```

### Indicative Example

```
DataFrame.write ←
    .format("parquet")
    .mode(saveMode)
    .option("path", "/data/flights/")
    .save()
```

### Indicative Example

```
DataFrame.write
    .format("parquet")
    .mode(saveMode)
    .option("path", "/data/flights/")
    .save()
```

**Built In Formats**  
CSV, JSON, Parquet, ORC, JDBC

**parquet is default format**

### Indicative Example

```
DataFrame.write
    .format("parquet")
    .mode(saveMode)
    .option("path", "/data/flights/")
    .save()
```

**Read Mode**  
1. append  
2. overwrite  
3. errorIfExists  
4. ignore

**Saving modes**

## Spark File Layout

1. Number of files and file size
2. Partitions and Buckets
3. Sorted data

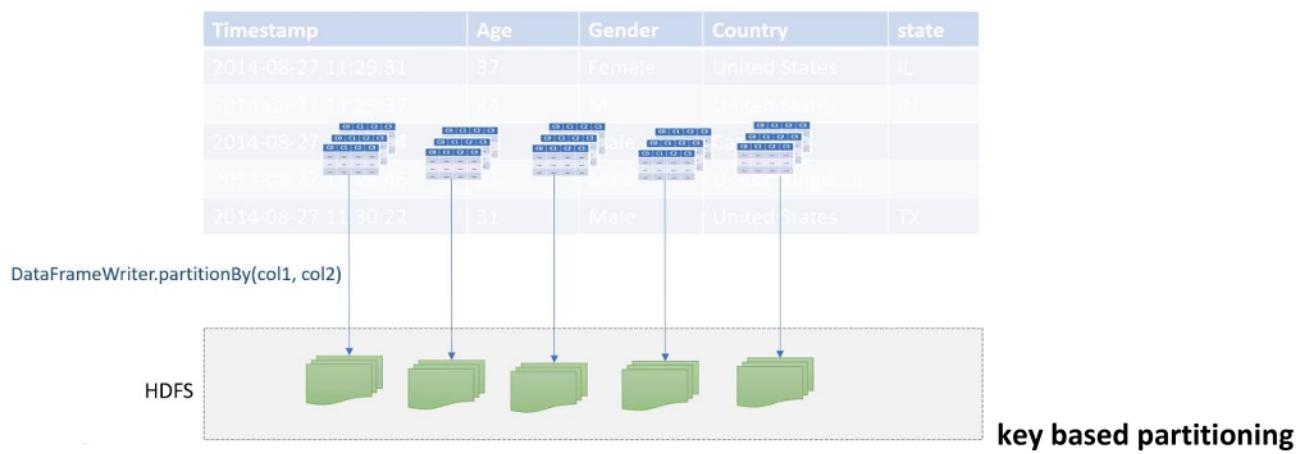
| Timestamp           | Age | Gender | Country       | state |
|---------------------|-----|--------|---------------|-------|
| 2014-08-27 11:29:31 | 37  | Female | United States | IL    |
| 2014-08-27 11:30:09 | 37  | Female | United States | TX    |
| 2014-08-27 11:30:10 | 37  | Female | United States | IL    |
| 2014-08-27 11:30:22 | 31  | Male   | United States | TX    |

DataFrame.repartition(n)

HDFS

**blind repartitioning**

## Partition by methods



DataFrameWriter.partitionBy(col1, col2) | only available on spark managed tables

- Writing your Data and Managing Layout

```
27 # spark.executor.extraJavaOptions -XX:+PrintGCDetails -Dkey=value -Dnumbers="one two"
28 spark.driver.extraJavaOptions -Dlog4j.configuration=file:log4j.properties -DSpar
29 spark.jars.packages org.apache.spark:spark-avro_2.12:3.3.2
```

### 2.12 Scala version 3.3.2 is Pyspark version

<https://spark.apache.org/docs/latest/sql-data-sources-avro.html#configuration>

Adding some Scala packages in `spark-defaults.conf` by above line. this allows to add jar in dependencies. Now you can read and write **AVRO** files.

### main.py (before default 2 partitions)

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import spark_partition_id

from lib.logger import Log4j

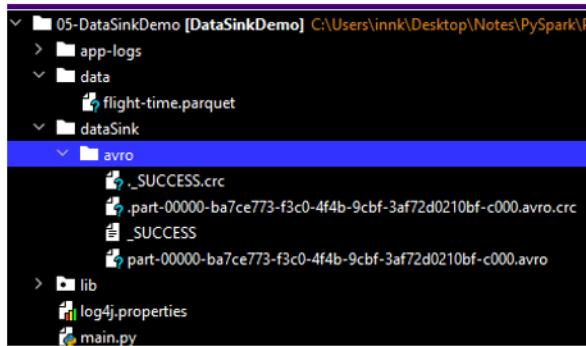
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("SparkSchemaDemo")\
        .getOrCreate()

    logger = Log4j(spark)

    flightTimeParquetDF = spark.read\
        .format("parquet")\
        .load("data/flight-time.parquet")

    flightTimeParquetDF.write\
        .format("avro")\
        .mode("overwrite")\
        .option("path", "dataSink/avro/")\
        .save()

    logger.info("Num Partition before: " + str(flightTimeParquetDF.rdd.getNumPartitions())) # to count partitions
    flightTimeParquetDF.groupBy(spark_partition_id()).count().show() # After executing this code we got only one avro file but we
expected 2 because we have 2 partitions by this count we will check that our all data is in 1st partition and 2nd don't have any data
that's why we got only one ARVO file. To resolve this problem following modifications are done.
```



```
0 artifacts copied, 3 already retrieved (0kB/29ms)
23/02/26 19:06:55 WARN ProcfsMetricsGetter: Exception when trying to read /proc/meminfo
23/02/26 19:06:57 INFO SparkSchemaDemo: Num Partition before: 2
```

```
+-----+-----+
|SPARK_PARTITION_ID()| count|
+-----+-----+
|          0|470477|
+-----+-----+
```

### main.py (After forced 5 partitions)

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import spark_partition_id

from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("SparkSchemaDemo")\
        .getOrCreate()

    logger = Log4j(spark)

    flightTimeParquetDF = spark.read\
        .format("parquet")\
        .load("data/flight-time.parquet")

    logger.info("Num Partition before: " + str(flightTimeParquetDF.rdd.getNumPartitions()))
    flightTimeParquetDF.groupBy(spark_partition_id()).count().show()

    partitionedDF = flightTimeParquetDF.repartition(5) # forced partitions
    logger.info("Num Partition after: " + str(partitionedDF.rdd.getNumPartitions()))
    partitionedDF.groupBy(spark_partition_id()).count().show()

    partitionedDF.write \
        .format("avro") \
        .mode("overwrite") \
        .option("path", "dataSink/avro/") \
        .save()
```

```
23/02/26 19:17:08 INFO SparkSchemaDemo: Num Partition after: 5
+-----+-----+
|SPARK_PARTITION_ID()|count|
+-----+-----+
0	94096
1	94095
2	94095
3	94095
4	94096
+-----+-----+
```

Repartition worked

```
+-----+
|SPARK_PARTITION_ID()|count|
+-----+
0	94096
1	94095
2	94095
3	94095
4	94096
+-----+
```

we got 5 AVRO files

Some times partitioning not suitable but partitioning improves following.

1. Parallel Processing
2. Partition Elimination

- Partition BY method

#### main.py

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import spark_partition_id

from lib.logger import Log4j

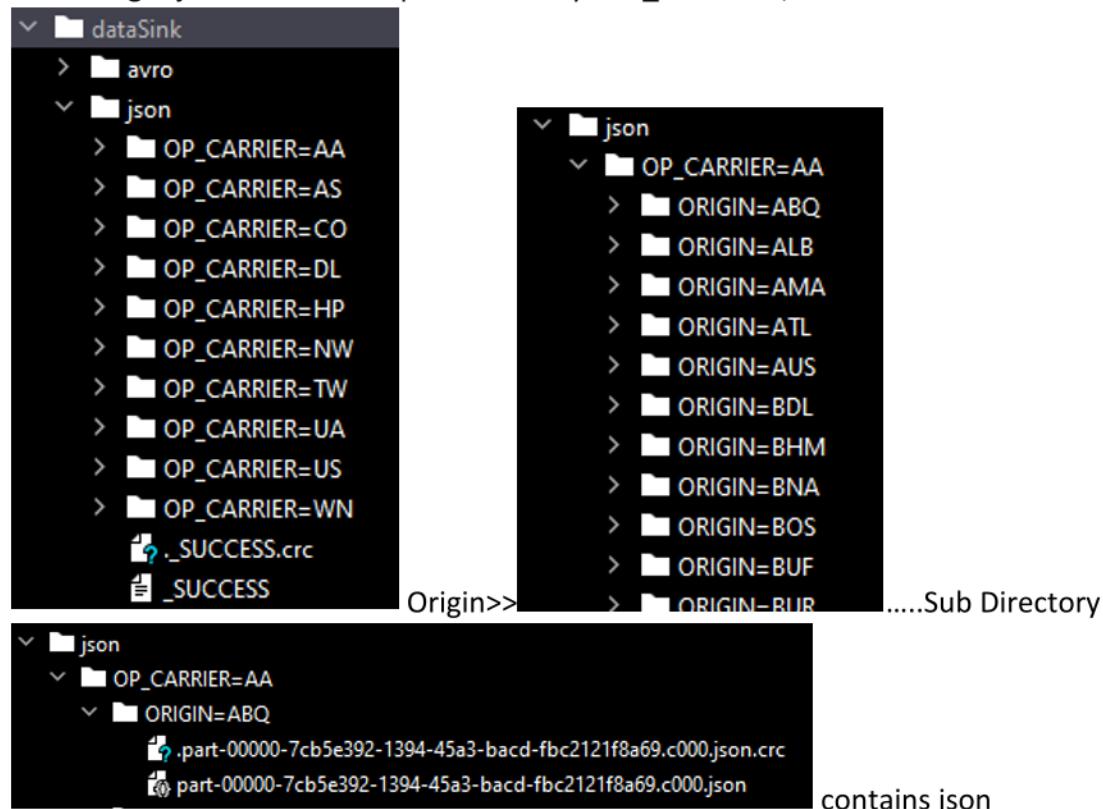
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("SparkSchemaDemo")\
        .getOrCreate()

    logger = Log4j(spark)

    flightTimeParquetDF = spark.read\
        .format("parquet")\
        .load("data/flight-time.parquet")

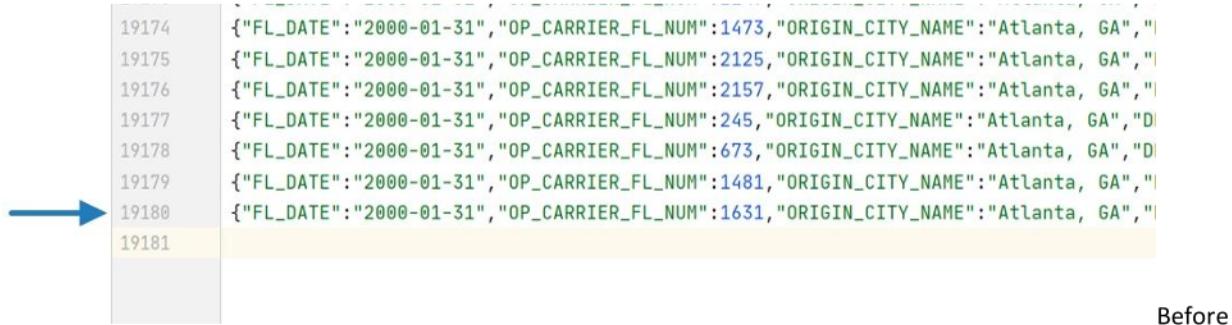
    flightTimeParquetDF.write\
        .format("json")\
        .mode("overwrite")\
        .option("path", "dataSink/json/")\
        .partitionBy("OP_CARRIER", "ORIGIN")\
        .save()
```

Here we got json format with partitioned by "OP\_CARRIER", "ORIGIN" as follow



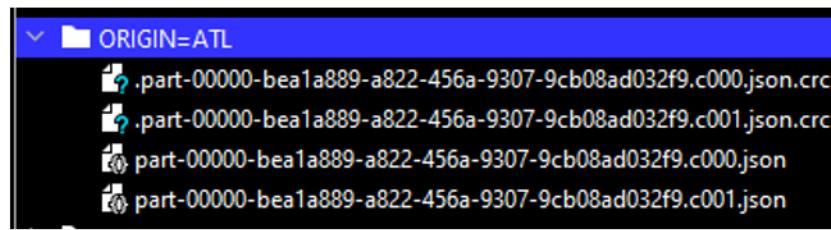
- Handling file size ( `.option("maxRecordsPerFile", 10000)` )

```
flightTimeParquetDF.write\  
.format("json")\  
.mode("overwrite")\  
.option("path","dataSink/json/")\  
.partitionBy("OP_CARRIER", "ORIGIN")\  
.option("maxRecordsPerFile", 10000)\  
.save()
```



19174 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 1473, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19175 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 2125, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19176 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 2157, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19177 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 245, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19178 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 673, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19179 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 1481, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19180 {"FL\_DATE": "2000-01-31", "OP\_CARRIER\_FL\_NUM": 1631, "ORIGIN\_CITY\_NAME": "Atlanta, GA", "D  
19181

Before



After two json files

|    |       |                                                                                                                                                                                                                                                                                                                                                                                                                            |
|----|-------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| rc | 9996  | {"FL_DATE": "2000-01-17", "OP_CARRIER_FL_NUM": 1473, "ORIGIN_CITY_NAME": "Atlanta, GA", "DISTANCE": 1000, "CARRIER": "AA", "ORIGIN": "ATL", "DESTINATION": "JFK", "DEPARTURE": "2000-01-17T00:00:00", "ARRIVAL": "2000-01-17T08:00:00", "DEPARTURE_DELAY": 0, "ARRIVAL_DELAY": 0, "CANCELLED": false, "DIVERTED": false, "CRASHED": false, "FLIGHT": "AA1473", "TAXI_IN": 0, "TAXI_OUT": 0, "FLYING": 1000, "TOTAL": 1000} |
| rc | 9997  | {"FL_DATE": "2000-01-17", "OP_CARRIER_FL_NUM": 2125, "ORIGIN_CITY_NAME": "Atlanta, GA", "DISTANCE": 1000, "CARRIER": "AA", "ORIGIN": "ATL", "DESTINATION": "JFK", "DEPARTURE": "2000-01-17T00:00:00", "ARRIVAL": "2000-01-17T08:00:00", "DEPARTURE_DELAY": 0, "ARRIVAL_DELAY": 0, "CANCELLED": false, "DIVERTED": false, "CRASHED": false, "FLIGHT": "AA2125", "TAXI_IN": 0, "TAXI_OUT": 0, "FLYING": 1000, "TOTAL": 1000} |
|    | 9998  | {"FL_DATE": "2000-01-17", "OP_CARRIER_FL_NUM": 2157, "ORIGIN_CITY_NAME": "Atlanta, GA", "DISTANCE": 1000, "CARRIER": "AA", "ORIGIN": "ATL", "DESTINATION": "JFK", "DEPARTURE": "2000-01-17T00:00:00", "ARRIVAL": "2000-01-17T08:00:00", "DEPARTURE_DELAY": 0, "ARRIVAL_DELAY": 0, "CANCELLED": false, "DIVERTED": false, "CRASHED": false, "FLIGHT": "AA2157", "TAXI_IN": 0, "TAXI_OUT": 0, "FLYING": 1000, "TOTAL": 1000} |
|    | 9999  | {"FL_DATE": "2000-01-17", "OP_CARRIER_FL_NUM": 245, "ORIGIN_CITY_NAME": "Atlanta, GA", "DISTANCE": 1000, "CARRIER": "AA", "ORIGIN": "ATL", "DESTINATION": "JFK", "DEPARTURE": "2000-01-17T00:00:00", "ARRIVAL": "2000-01-17T08:00:00", "DEPARTURE_DELAY": 0, "ARRIVAL_DELAY": 0, "CANCELLED": false, "DIVERTED": false, "CRASHED": false, "FLIGHT": "AA245", "TAXI_IN": 0, "TAXI_OUT": 0, "FLYING": 1000, "TOTAL": 1000}   |
|    | 10000 | {"FL_DATE": "2000-01-17", "OP_CARRIER_FL_NUM": 673, "ORIGIN_CITY_NAME": "Atlanta, GA", "DISTANCE": 1000, "CARRIER": "AA", "ORIGIN": "ATL", "DESTINATION": "JFK", "DEPARTURE": "2000-01-17T00:00:00", "ARRIVAL": "2000-01-17T08:00:00", "DEPARTURE_DELAY": 0, "ARRIVAL_DELAY": 0, "CANCELLED": false, "DIVERTED": false, "CRASHED": false, "FLIGHT": "AA673", "TAXI_IN": 0, "TAXI_OUT": 0, "FLYING": 1000, "TOTAL": 1000}   |
|    | 10001 |                                                                                                                                                                                                                                                                                                                                                                                                                            |

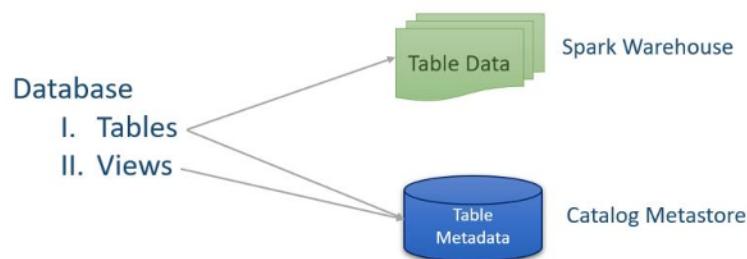
1000 limit

What we learn

- 1. How to use DataFrameWriter API  
 2. How to use PartitionBy(...)  
 3. How to control file size using maxRecordsPerFile

- **Spark Databases and Tables**

- Apache Spark is database itself you can create database in it.

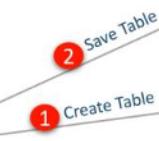


**How to Create Managed Table?**

```
dataframe.write  
.saveAsTable("your_table_name")
```

**Spark Tables**

- I. Managed Tables
- II. Unmanaged Tables (External Tables)



spark.sql.warehouse.dir



**How to Create Unmanaged Table?**

```
CREATE TABLE your_tbl_name (col1 data_type,  
                           col2 data_type,  
                           ....)  
USING PARQUET  
LOCATION "data_file_location"
```

**Spark Tables**

- I. Managed Tables
- II. Unmanaged Tables (External Tables)



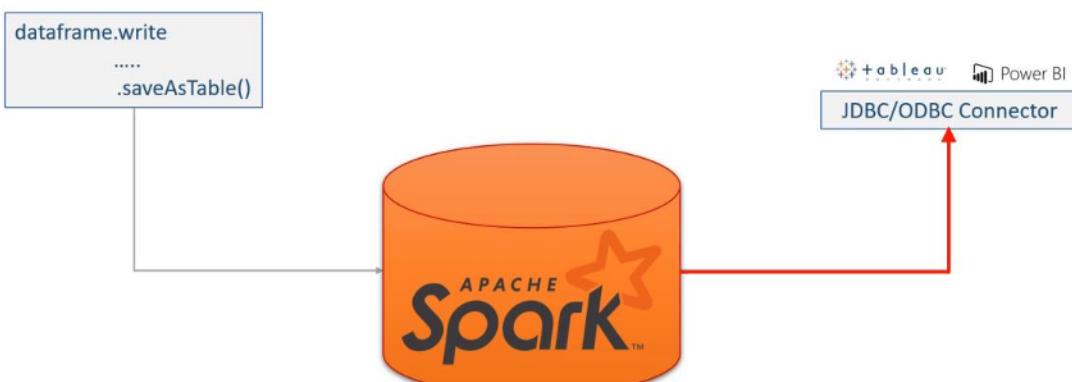
Catalog Metastore

- **Working with Spark SQL Table**

- Need Hive for this example.
- Create manage table and store DF into spark table.
- In real time scenario you will process your data and save the output data frame.



parquet or avro format



main.py

```

from pyspark.sql import *
from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("SparkSQLTableDemo")\
        .enableHiveSupport()\
        .getOrCreate()

    logger = Log4j(spark)

    flightTimeParquetDF = spark.read\
        .format("parquet")\
        .load("dataSource/")

    spark.sql("CREATE DATABASE IF NOT EXISTS AIRLINE_DB")
    spark.catalog.setCurrentDatabase("AIRLINE_DB")

    # Create manage table and store DF into spark table(without any processing)

    flightTimeParquetDF.write\
        .mode("overwrite")\
        .saveAsTable("flight_data_tbl")

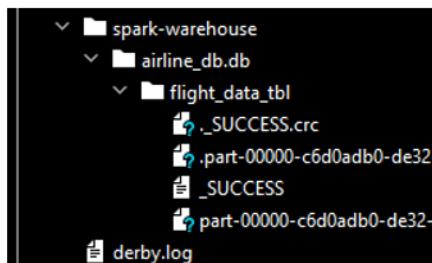
    logger.info(spark.catalog.listTables("AIRLINE_DB"))

```

```
23/02/27 14:16:59 INFO SparkSQLTableDemo: [[flight_data_tbl, airline_db, null, MANAGED, false]]
```

```
Process finished with exit code 0
```

Here is catalog information of DB



### Partition by method

```
flightTimeParquetDF.write\  
.mode("overwrite")\  
.partitionBy("ORIGIN", "OP_CARRIER")\  
.saveAsTable("flight_data_tbl")
```

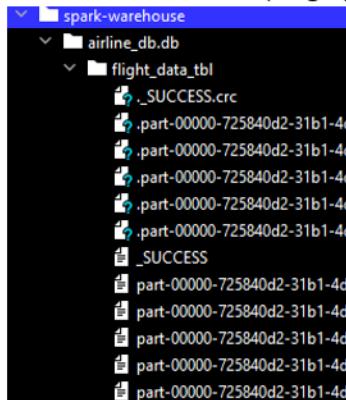


here we got 200+ dir.

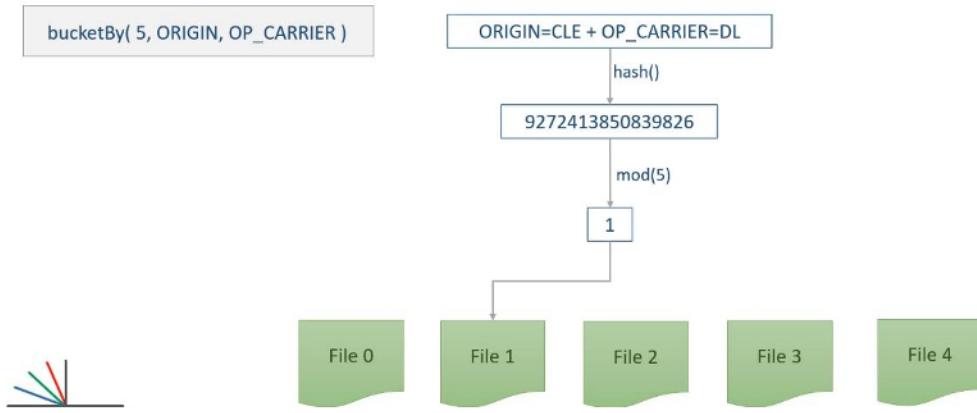
### So, When Table contains to many unique values you need to use bucket by

- The bucket by allows you to restrict the partitions.

```
flightTimeParquetDF.write\  
.format("csv") # here parquet is recommended format we are using csv only for investing data  
.mode("overwrite")  
.bucketBy(5, "ORIGIN", "OP_CARRIER")  
.saveAsTable("flight_data_tbl")
```



Here we got 5 CSV files (bucket)



Sort by

```

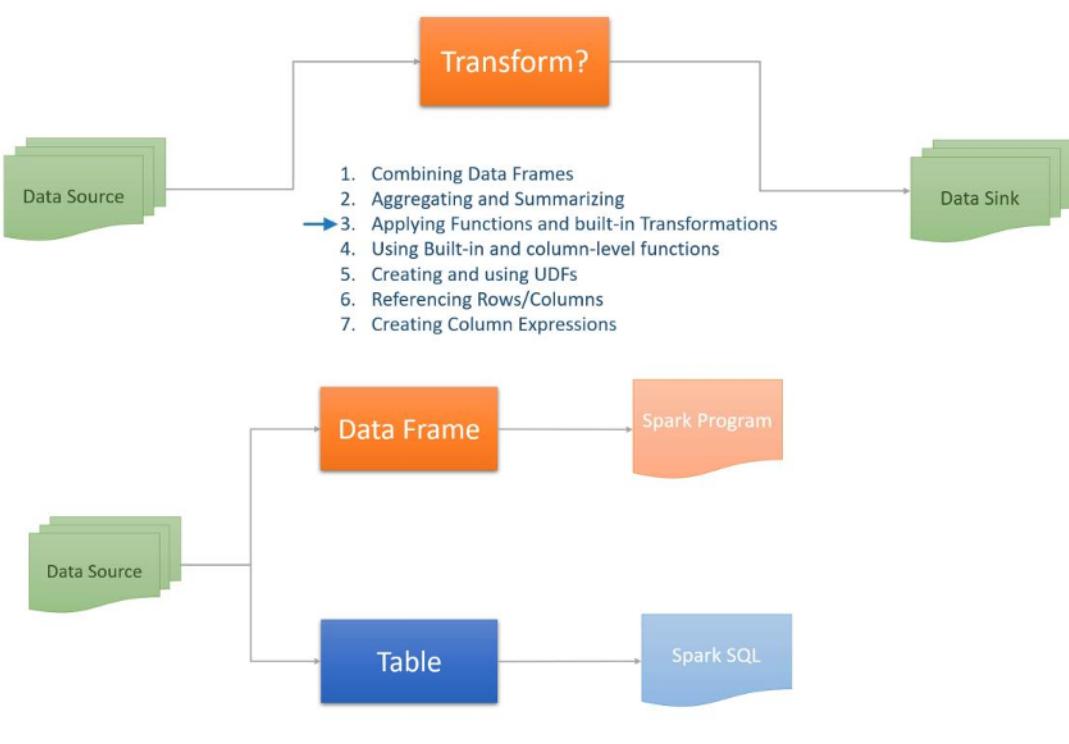
flightTimeParquetDF.write\
    .format("csv")\
    .mode("overwrite")\
    .bucketBy(5, "ORIGIN", "OP_CARRIER")\
    .sortBy("ORIGIN", "OP_CARRIER")\
    .saveAsTable("flight_data_tbl")

```

here we got sorted values columns in bucket

|    |                                          |
|----|------------------------------------------|
| 61 | 2000-01-10,AA,438,ABQ,"Albuquerque, NM"  |
| 62 | 2000-01-10,AA,1182,ABQ,"Albuquerque, NM" |
| 63 | 2000-01-10,AA,1216,ABQ,"Albuquerque, NM" |
| 64 | 2000-01-10,AA,1306,ABQ,"Albuquerque, NM" |
| 65 | 2000-01-10,AA,1400,ABQ,"Albuquerque, NM" |
| 66 | 2000-01-10,AA,1498,ABQ,"Albuquerque, NM" |
| 67 | 2000-01-11,AA,438,ABQ,"Albuquerque, NM"  |
| 68 | 2000-01-11,AA,1182,ABQ,"Albuquerque, NM" |
| 69 | 2000-01-11,AA,1216,ABQ,"Albuquerque, NM" |
| 70 | 2000-01-11,AA,1306,ABQ,"Albuquerque, NM" |

- **Spark Dataframe and Dataset Transformation** [ref.48](#)



**Dataframe = Dataset[Row]**

| Timestamp           | Age | Gender | Country        | state |
|---------------------|-----|--------|----------------|-------|
| 2014-08-27 11:29:31 | 37  | Female | United States  | IL    |
| 2014-08-27 11:29:37 | 44  | M      | United States  | IN    |
| 2014-08-27 11:29:44 | 32  | Male   | Canada         |       |
| 2014-08-27 11:29:46 | 31  | Male   | United Kingdom |       |
| 2014-08-27 11:30:22 | 31  | Male   | United States  | TX    |

1. Manually creating Rows and Dataframe.
2. Collecting Dataframe rows to the driver.
- 3. Work with an individual row in Spark Transformations.

- **Working with DataFrame Rows**

## Manually creating Rows and Dataframe

The screenshot shows the Databricks UI for managing clusters. On the left, there's a sidebar with various icons. The main area displays cluster configuration details: Databricks Runtime Version (12.1), Driver type (Community Optimized), and Instances (us-west-2b). A modal window titled 'Create Notebook' is overlaid, asking for a name ('MyPythonNotebook'), default language ('Python'), and cluster ('MyTestCluster').

```
+---+-----+
| ID| EventDate|
+---+-----+
123	04/05/2020
124	4/5/2020
125	04/5/2020
125	4/05/2020
+---+-----+
```

```
+---+-----+
| ID| EventDate|
+---+-----+
123	2020-04-05
124	2020-04-05
125	2020-04-05
125	2020-04-05
+---+-----+
```

>>

Cmd 1

```
1 def to_date_df(df, fmt, fld):
2     return df.withColumn(fld, to_date(col(fld)), fmt)
```

- 1. A large project might need hundreds of small sample data files to test your functions.
- 2. Your build pipeline might run slow due to loading hundreds of sample files and increased I/O.

### Databricks code:

```
from pyspark.sql import *
from pyspark.sql.functions import *
from pyspark.sql.types import *

def to_date_df(df, fmt, fld):
    return df.withColumn(fld, to_date(col(fld), fmt))

my_schema = StructType([
    StructField("ID", StringType()),
    StructField("EventDate", StringType())])

my_rows = [Row("123", "04/05/2020"), Row("124", "4/5/2020"), Row("125", "04/5/2020"), Row("126", "4/05/2020")]
my_rdd = spark.sparkContext.parallelize(my_rows, 2)
my_df = spark.createDataFrame(my_rdd, my_schema)

my_df.printSchema()
my_df.show()
root
|-- ID: string (nullable = true)
|-- EventDate: string (nullable = true)

+---+-----+
| ID| EventDate|
+---+-----+
123	04/05/2020
124	4/5/2020
125	04/5/2020
126	4/05/2020
+---+-----+

new_df = to_date_df(my_df, "M/d/y", "EventDate")
new_df.printSchema()
new_df.show()
```

```
root
|-- ID: string (nullable = true)
|-- EventDate: date (nullable = true)

+---+-----+
| ID| EventDate|
+---+-----+
123	2020-04-05
124	2020-04-05
125	2020-04-05
126	2020-04-05
+---+-----+
```

## Collecting Dataframe rows to the driver.

- DataFrame Rows and Unit testing

[ref.](#) **RowDemo** for **main.py** in PyCharm with above same code

## Unit Test for above code

```
from datetime import date
from unittest import TestCase

from pyspark.sql import *
from pyspark.sql.types import *

from RowDemo import to_date_df

class RowDemoTestCase(TestCase):

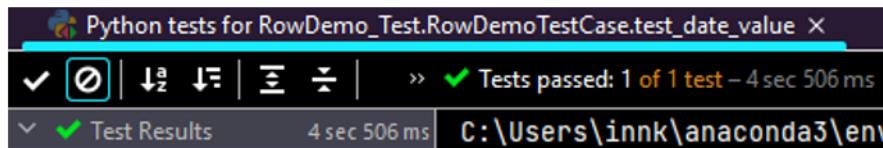
    @classmethod
    def setUpClass(cls) -> None:
        cls.spark = SparkSession.builder \
            .master("local[3]") \
            .appName("RowDemoTest") \
            .getOrCreate()

        my_schema = StructType([
            StructField("ID", StringType()),
            StructField("EventDate", StringType())])

        my_rows = [Row("123", "04/05/2020"), Row("124", "4/5/2020"), Row("125", "04/5/2020"), Row("126", "4/05/2020")]
        my_rdd = cls.spark.sparkContext.parallelize(my_rows, 2)
        cls.my_df = cls.spark.createDataFrame(my_rdd, my_schema)

    def test_data_type(self):
        rows = to_date_df(self.my_df, "M/d/y", "EventDate").collect()
        for row in rows:
            self.assertIsInstance(row["EventDate"], date)

    def test_date_value(self):
        rows = to_date_df(self.my_df, "M/d/y", "EventDate").collect()
        for row in rows:
            self.assertEqual(row["EventDate"], date(2020, 4, 5))
```



## Work with an individual row in Spark Transformations.

- DataFrame Rows and Unstructured Data

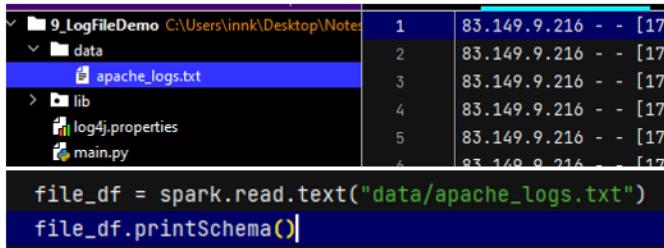
- Spark DF offers bunch of Transformation functions.

```
agg(*exprs)
cov(coll, col2)
crosstab(coll, col2)
cube(*cols)
filter(condition)
groupBy(*cols)
join(other, on=None, how=None)
orderBy(*cols, **kwargs)
replace(to_replace, value, subset)
rollup(*cols)

select(*cols)
sort(*cols, **kwargs)
where(condition)
withColumn(colName, col)
avg(*cols)
max(*cols)
mean(*cols)
min(*cols)
sum(*cols)
```

- When your DF doesn't have column structure Then you have to create columnar structure and then use Transformation functions.

Here we have Apache Log Server File(which contain some pattern and it unstructured):



```
file_df = spark.read.text("data/apache_logs.txt")
file_df.printSchema()
```

|-- value: string (nullable = true)

we have only one string field "value".

log\_reg = r'^(\S+) (\S+) (\S+) \[(\w:/]+\s[+\-]\d{4}\] "(\S+) (\S+) (\S+)" (\d{3}) (\S+) "(S+)" "([^\"]\*)"

this above regex gives following filed

1. IP
2. client.
3. user
4. datetime
5. cmd
6. request
7. protocol
8. status
9. bytes
10. referrer
11. userAgent

```
logs_df = file_df.select(regexp_extract('value', log_reg, 1).alias('ip'),
                           regexp_extract('value', log_reg, 4).alias('date'),
                           regexp_extract('value', log_reg, 6).alias('request'),
                           regexp_extract('value', log_reg, 10).alias('referrer'))
```

```
|-- ip: string (nullable = true)
|-- date: string (nullable = true)
|-- request: string (nullable = true)
|-- referrer: string (nullable = true)
```

Now we got 4 schema

```
logs_df.printSchema()

logs_df.groupBy("referrer")\
    .count()\
    .show(100, truncate=False)
```

```
+-----+
|referrer
+-----+
|http://www.semicomplete.com/blog/tags/wifi
|http://manpages.ubuntu.com/manpages/lucid/man1/
|http://s-chassis.co.nz/viewtopic.php?f=16&t=926
|https://www.google.sk/
|http://www.semicomplete.com/projects/keynav/
|http://superuser.com/questions/355151/is-it-pos
|https://www.google.ca/
|http://www.google.de/uol?sa=t&st=ts&sq=ss
```

but we got complete urls so following

## main.py

```
from pyspark.sql import *
from pyspark.sql.functions import regexp_extract, substring_index

if __name__ == "__main__":
    spark = SparkSession \
        .builder \
        .master("local[3]") \
        .appName("LogFileDemo") \
        .getOrCreate()

    file_df = spark.read.text("data/apache_logs.txt")

    log_reg = r'^(\S+) (\S+) (\S+) \[(\w:/+\s[+-]\d{4})\] "(\S+) (\S+) (\S+)" (\d{3}) (\S+) "([""]*)'

    logs_df = file_df.select(regexp_extract('value', log_reg, 1).alias('ip'),
                             regexp_extract('value', log_reg, 4).alias('date'),
                             regexp_extract('value', log_reg, 6).alias('request'),
                             regexp_extract('value', log_reg, 10).alias('referrer'))

    logs_df.printSchema()

    logs_df \
        .where("trim(referrer) != '-'") \
        .withColumn("referrer", substring_index("referrer", "/", 3)) \
        .groupBy("referrer") \
        .count() \
        .show(100, truncate=False)
```

| referrer                                                                          | count |
|-----------------------------------------------------------------------------------|-------|
| <a href="http://ijavascript.cn">http://ijavascript.cn</a>                         | 1     |
| <a href="http://www.google.co.tz">http://www.google.co.tz</a>                     | 1     |
| <a href="http://www.google.ca">http://www.google.ca</a>                           | 6     |
| <a href="https://www.qooole.hr">https://www.qooole.hr</a>                         | 2     |
| <a href="https://www.qooole.ch">https://www.qooole.ch</a>                         | 1     |
| <a href="http://www.qooole.ru">http://www.qooole.ru</a>                           | 6     |
| <a href="http://www.raspberrypi-spanish.es">http://www.raspberrypi-spanish.es</a> | 1     |
| <a href="http://semicomplete.com">http://semicomplete.com</a>                     | 2001  |
| <a href="http://manpages.ubuntu.com">http://manpages.ubuntu.com</a>               | 2     |
| <a href="http://www.google.fi">http://www.google.fi</a>                           | 4     |
| -                                                                                 | 4073  |
| <a href="https://www.google.co.za">https://www.google.co.za</a>                   | 1     |

- Working with DataFrame Columns

**Column**

| Timestamp           | Age | Gender | Country        | state |
|---------------------|-----|--------|----------------|-------|
| 2014-08-27 11:29:31 | 37  | Female | United States  | IL    |
| 2014-08-27 11:29:37 | 44  | M      | United States  | IN    |
| 2014-08-27 11:29:44 | 32  | Male   | Canada         |       |
| 2014-08-27 11:29:46 | 31  | Male   | United Kingdom |       |
| 2014-08-27 11:30:22 | 31  | Male   | United States  | TX    |

1. What is a Column and How to reference it?
2. How to create column expressions?

## DATA BRICKS

Cmd 1

```
1 %fs ls /databricks-datasets/
```

Table +

| path                                        |
|---------------------------------------------|
| 1 dbfs:/databricks-datasets/                |
| 2 dbfs:/databricks-datasets/COVID/          |
| 3 dbfs:/databricks-datasets/README.md       |
| 4 dbfs:/databricks-datasets/Rdatasets/      |
| 5 dbfs:/databricks-datasets/SPARK_README.md |
| 6 dbfs:/databricks-datasets/adult/          |
| 7 dbfs:/databricks-datasets/airlines/       |

↓ 56 rows | 51.75 seconds runtime

Command took 51.75 seconds -- by innksn@gmail.com at 2/27/2023, 9 we can get this dataset which are provided by databricks.

Cmd 2

```
1 %fs ls /databricks-datasets/airlines/
```

Table +

| path                                            | nam   |
|-------------------------------------------------|-------|
| 1 dbfs:/databricks-datasets/airlines/README.md  | REAL  |
| 2 dbfs:/databricks-datasets/airlines/_SUCCESS   | _SUCC |
| 3 dbfs:/databricks-datasets/airlines/part-00000 | part  |
| 4 dbfs:/databricks-datasets/airlines/part-00001 | part  |
| 5 dbfs:/databricks-datasets/airlines/part-00002 | part  |
| 6 dbfs:/databricks-datasets/airlines/part-00003 | part  |
| 7 dbfs:/databricks-datasets/airlines/part-00004 | part  |

↓ ▾ 1,000 rows | Truncated data ▾ | 2.86 seconds runti

Command took 2.86 seconds -- by innksn@gmail.com at 2/27/2023

Cmd 3

```
1 %fs head /databricks-datasets/airlines/part-00000
```

[Truncated to first 65536 bytes]

```
Year,Month,DayofMonth,DayOfWeek,DeptTime,CRSDepTime,ArrTime,CRSArrTime,UniqueCarrier,FlightNum,TailNum,ActualElapsedTime,CarrierDelay,WeatherDelay,NASDelay,SecurityDelay,LateAircraftDelay,IsArrDelayed,IsDepDelayed
1987,10,14,3,741,730,912,849,PS,1451,NA,91,79,NA,23,11,SAN,SFO,447,NA,NA,0,NA,NA,NA,NA,NA,YES,YES
1987,10,15,4,729,730,903,849,PS,1451,NA,94,79,NA,14,-1,SAN,SFO,447,NA,NA,0,NA,NA,NA,NA,NA,NA,NO,YES,NO
1987,10,17,6,741,730,918,849,PS,1451,NA,97,79,NA,29,11,SAN,SFO,447,NA,NA,0,NA,NA,NA,NA,NA,YES,YES
1987,10,18,7,729,730,847,849,PS,1451,NA,78,79,NA,-2,-1,SAN,SFO,447,NA,NA,0,NA,NA,NA,NA,NA,NO,NO
1987.10.19.1.749.730.922.849.PS.1451.NA.93.79.NA.33.19.SAN.SFO.447.NA.NA.0.NA.NA.NA.NA.YES.YES
```

Cmd 4

```
1 airlinesDF = spark.read \
2   .format("csv") \
3   .option("header", "true") \
4   .option("inferSchema", "true") \
5   .option("samplingRatio", "0.0001") \
6   .load("/databricks-datasets/airlines/part-00000")
```

Cmd 5

```
1 airlinesDF.select("Origin", "Dest", "Distance").show(10)
```

▶ (1) Spark Jobs

```
+-----+-----+
|Origin|Dest|Distance|
+-----+-----+
| SAN| SFO|    447|
+-----+-----+
```

only showing top 10 rows

Column String

Cmd 6

```
1 from pyspark.sql.functions import *
2 airlinesDF.select(column("Origin"), col("Dest"), "Distance").show(10)
```

▶ (1) Spark Jobs

```
+-----+-----+
|Origin|Dest|Distance|
+-----+-----+
| SAN| SFO|    447|
+-----+-----+
```

only showing top 10 rows

Column Object

Cmd 7

```
1 airlinesDF.select("Origin", "Dest", "Distance", "Year", "Month", "DayofMonth").show(10)
```

▶ (1) Spark Jobs

```
+-----+-----+-----+-----+
|Origin|Dest|Distance|Year|Month|DayofMonth|
+-----+-----+-----+-----+
SAN	SFO	447	1987	10	14
SAN	SFO	447	1987	10	15
SAN	SFO	447	1987	10	17
SAN	SFO	447	1987	10	18
SAN	SFO	447	1987	10	19
SAN	SFO	447	1987	10	21
SAN	SFO	447	1987	10	22
SAN	SFO	447	1987	10	23
SAN	SFO	447	1987	10	24
SAN	SFO	447	1987	10	25
+-----+-----+-----+-----+
```

only showing top 10 rows

by using column object

→ 1. String Expressions or SQL Expressions  
2. Column Object Expressions

Cmd 8

```
1 airlinesDF.select("Origin", "Dest", "Distance", expr("to_date(concat(Year,Month,DayofMonth),'yyyyMMdd') as FlightDate")).show(10)

▶ (1) Spark Jobs

+-----+-----+-----+
|Origin|Dest|Distance|FlightDate|
+-----+-----+-----+
SAN	SFO	447	1987-10-14
SAN	SFO	447	1987-10-15
SAN	SFO	447	1987-10-17
SAN	SFO	447	1987-10-18
SAN	SFO	447	1987-10-19
SAN	SFO	447	1987-10-21
SAN	SFO	447	1987-10-22
SAN	SFO	447	1987-10-23
SAN	SFO	447	1987-10-24
SAN	SFO	447	1987-10-25
+-----+
only showing top 10 rows
```

Cmd 8

```
1 airlinesDF.selectExpr("Origin", "Dest", "Distance", "to_date(concat(Year,Month,DayofMonth),'yyyyMMdd') as FlightDate").show(10)

▶ (1) Spark Jobs

+-----+-----+-----+
|Origin|Dest|Distance|FlightDate|
+-----+-----+-----+
SAN	SFO	447	1987-10-14
SAN	SFO	447	1987-10-15
SAN	SFO	447	1987-10-17
SAN	SFO	447	1987-10-18
SAN	SFO	447	1987-10-19
SAN	SFO	447	1987-10-21
SAN	SFO	447	1987-10-22
SAN	SFO	447	1987-10-23
SAN	SFO	447	1987-10-24
SAN	SFO	447	1987-10-25
+-----+
only showing top 10 rows
```

By using Column Object expression ref.52 / 10 folder

Cmd 9

```
1 airlinesDF.select("Origin", "Dest", "Distance", to_date(concat("Year","Month","DayofMonth"),"yyyyMMdd").alias("FlightDate")).show(10)

▶ (1) Spark Jobs

+-----+-----+-----+
|Origin|Dest|Distance|FlightDate|
+-----+-----+-----+
SAN	SFO	447	1987-10-14
SAN	SFO	447	1987-10-15
SAN	SFO	447	1987-10-17
SAN	SFO	447	1987-10-18
SAN	SFO	447	1987-10-19
SAN	SFO	447	1987-10-21
SAN	SFO	447	1987-10-22
SAN	SFO	447	1987-10-23
SAN	SFO	447	1987-10-24
SAN	SFO	447	1987-10-25
+-----+
only showing top 10 rows
```

- **Creating and Using UDF(User Defined Functions)**

### To show DF

```
from pyspark.sql import SparkSession
```

```
from lib.logger import Log4j
```

```
if __name__ == "__main__":
    spark = SparkSession \
        .builder \
        .appName("UDF Demo") \
        .master("local[2]") \
        .getOrCreate()
```

```
logger = Log4j(spark)
```

```
survey_df = spark.read \
    .option("header", "true") \
    .option("inferSchema", "true") \
    .csv("data/survey.csv")
```

```
survey_df.show(10)
```

```
+-----+-----+-----+
| 2014-08-27 11:29:31| 37|Female| United States| I
| 2014-08-27 11:29:37| 44|      M| United States| I
| 2014-08-27 11:29:44| 32|  Male|      Canada| N
| 2014-08-27 11:29:46| 31|  Male|United Kingdom| N
| 2014-08-27 11:30:22| 31|  Male| United States| T
| 2014-08-27 11:31:22| 33|  Male| United States| T
| 2014-08-27 11:31:50| 35|Female| United States| M
| 2014-08-27 11:32:05| 39|      M|      Canada| N
| 2014-08-27 11:32:39| 42|Female| United States| I
| 2014-08-27 11:32:43| 23|  Male|      Canada| N
+-----+-----+-----+
only showing top 10 rows
```

```

import re

from pyspark.sql import SparkSession
from pyspark.sql.functions import udf, expr
from pyspark.sql.types import StringType

from lib.logger import Log4j


def parse_gender(gender): #UDF
    female_pattern = r"^\f\$|f.m|w.m" # female ^f$ or f.m or w.m women
    male_pattern = r"^\m\$|ma|m.l" # male ^m$ or ma or m.l
    if re.search(female_pattern, gender.lower()):
        return "Female"
    elif re.search(male_pattern, gender.lower()):
        return "Male"
    else:
        return "Unknown"

if __name__ == "__main__":
    spark = SparkSession \
        .builder \
        .appName("UDF Demo") \
        .master("local[2]") \
        .getOrCreate()

    logger = Log4j(spark)

    survey_df = spark.read \
        .option("header", "true") \
        .option("inferSchema", "true") \
        .csv("data/survey.csv")

    survey_df.show(10)

    #Column Object Expression
    # registering to column object expression.
    parse_gender_udf = udf(parse_gender, returnType=StringType())
    logger.info("Catalog Entry:")
    [logger.info(r) for r in spark.catalog.listFunctions() if "parse_gender" in r.name]

    # withcolum is allows as to transform the column without impaction other column of DF.
    survey_df2 = survey_df.withColumn("Gender", parse_gender_udf("Gender"))
    survey_df2.show(10)

    #UDF reg. this registered as SQL function and also create one entry in catalog.
    spark.udf.register("parse_gender_udf", parse_gender, StringType())
    logger.info("Catalog Entry:")
    [logger.info(r) for r in spark.catalog.listFunctions() if "parse_gender" in r.name]

    survey_df3 = survey_df.withColumn("Gender", expr("parse_gender_udf(Gender)"))
    survey_df3.show(10)

```

```
+-----+---+-----+
|           Timestamp|Age|Gender|
+-----+---+-----+
2014-08-27 11:29:31	37	Female
2014-08-27 11:29:37	44	M
2014-08-27 11:29:44	32	Male
2014-08-27 11:29:46	31	Male
2014-08-27 11:30:22	31	Male
2014-08-27 11:31:22	33	Male
2014-08-27 11:31:50	35	Female
2014-08-27 11:32:05	39	M
2014-08-27 11:32:39	42	Female
2014-08-27 11:32:43	23	Male
+-----+---+-----+
only showing top 10 rows
```

Original DF O/p

```
+-----+---+-----+
|           Timestamp|Age|Gender|
+-----+---+-----+
2014-08-27 11:29:31	37	Female
2014-08-27 11:29:37	44	Male
2014-08-27 11:29:44	32	Male
2014-08-27 11:29:46	31	Male
2014-08-27 11:30:22	31	Male
2014-08-27 11:31:22	33	Male
2014-08-27 11:31:50	35	Female
2014-08-27 11:32:05	39	M
2014-08-27 11:32:39	42	Female
2014-08-27 11:32:43	23	Male
+-----+---+-----+
only showing top 10 rows
```

Transformed DF with Column Object

```
+-----+---+-----+
|           Timestamp|Age|Gender|
+-----+---+-----+
2014-08-27 11:29:31	37	Female
2014-08-27 11:29:37	44	Male
2014-08-27 11:29:44	32	Male
2014-08-27 11:29:46	31	Male
2014-08-27 11:30:22	31	Male
2014-08-27 11:31:22	33	Male
2014-08-27 11:31:50	35	Female
2014-08-27 11:32:05	39	M
2014-08-27 11:32:39	42	Female
2014-08-27 11:32:43	23	Male
+-----+---+-----+
only showing top 10 rows
```

Transformed DF with SQL Expression

- **Misc(Miscellaneous) Transformation**

### Miscellaneous

1. Quick method to create Dataframes
2. Adding monotonically increasing id
3. Using Case When Then transformation
4. Casting your columns
5. Adding columns to Dataframes
6. Dropping Columns
7. Dropping duplicate rows
8. Sorting Dataframes

**ref. 54**

- **Aggregation in Apache Spark**

1. Simple Aggregations
2. Grouping Aggregations
3. Windowing Aggregations

this are built works via built in function.

The screenshot shows the 'Welcome to Spark Python API Docs!' page. In the sidebar under 'Contents', the 'pyspark.sql.functions' module is highlighted with a red box.

## Aggregating Functions

avg()  
count()  
max()  
min()  
sum()

## Window Functions

lead()  
lag()  
rank()  
dense\_rank()  
cume\_dist()

### main.py

```
from pyspark.sql import SparkSession, functions as f
```

```
from lib.logger import Log4j
```

```
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[2]")\
        .appName("AggFun")\
        .getOrCreate()
```

```
logger = Log4j(spark)
```

```
invoice_df = spark.read\
    .format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("data/invoices.csv")
```

```
invoice_df.select(f.count("*").alias("Count *"),
                  f.sum("Quantity").alias("TotalQuantity"),
                  f.avg("UnitPrice").alias("AvgPrice"),
                  f.countDistinct("InvoiceNo").alias("CountDistinct")).show() By Object Function
```

| Count * | TotalQuantity | AvgPrice          | CountDistinct |
|---------|---------------|-------------------|---------------|
| 541909  | 5176450       | 4.611113626088498 | 25900         |

Simple Aggregation

## main.py

```
from pyspark.sql import SparkSession, functions as f
from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[2]")\
        .appName("AggFun")\
        .getOrCreate()

    logger = Log4j(spark)

    invoice_df = spark.read\
        .format("csv")\
        .option("header", "true")\
        .option("inferSchema", "true")\
        .load("data/invoices.csv")

    invoice_df.selectExpr(
        "count(1) as `count 1`",
        "count(StockCode) as `count field`",
        "sum(Quantity) as TotalQuantity",
        "avg(UnitPrice) as AvgPrice"
    ).show()
```

By using SQL Expression

```
+-----+-----+-----+
|count 1|count field|TotalQuantity|      AvgPrice|
+-----+-----+-----+
| 541909|      541908|      5176450|4.61111362608295|
+-----+-----+-----+
```

Count field is lesser than above because of in SQL count null are not countable.

## Grouping Aggregations?

```
invoice_df.createOrReplaceTempView("sales")
summary_sql = spark.sql("""
    SELECT Country, InvoiceNo,
           sum(Quantity) as TotalQuantity,
           round(sum(Quantity*UnitPrice),2) as InvoiceValue
    FROM sales
    GROUP BY Country, InvoiceNo""")
```

summary\_sql.show()

Group BY using Spark.SQL

```
+-----+-----+-----+
|      Country|InvoiceNo|TotalQuantity|InvoiceValue|
+-----+-----+-----+
United Kingdom	536446	329	440.89
United Kingdom	536508	216	155.52
United Kingdom	537018	-3	0.0
United Kingdom	537401	-24	0.0
United Kingdom	537811	74	268.86
+-----+-----+-----+
```

## Above grouping by DF Expressions

```
summary_df = invoice_df \
    .groupBy("Country", "InvoiceNo") \
    .agg(f.sum("Quantity").alias("TotalQuantity"), #agg for aggregation
         f.round(f.sum(f.expr("Quantity * UnitPrice")), 2).alias("InvoiceValue"),
         f.expr("round(sum(Quantity * UnitPrice),2) as InvoiceValueExpr")
    )
```

```
summary_df.show()
```

### ref. 56 Grouping Aggregation\*\*\*

from pyspark.sql import SparkSession, functions as f

```
from lib.logger import Log4j
```

```
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[2]")\
        .appName("AggFun")\
        .getOrCreate()
```

```
logger = Log4j(spark)
```

```
invoice_df = spark.read\
    .format("csv")\
    .option("header", "true")\
    .option("inferSchema", "true")\
    .load("data/invoices.csv")
```

```
NumInvoices = f.countDistinct("InvoiceNo").alias("NumInvoices")
```

```
TotalQuantity = f.sum("Quantity").alias("TotalQuantity")
```

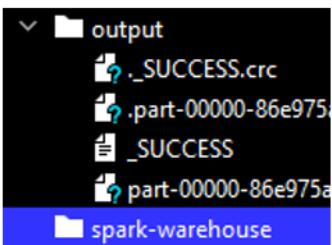
```
InvoiceValue = f.expr("round(sum(Quantity * UnitPrice),2) as InvoiceValue")
```

```
exSummary_df = invoice_df \
    .withColumn("InvoiceDate", f.to_date(f.col("InvoiceDate"), "dd-MM-yyyy H.mm")) \
    .where("year(InvoiceDate) == 2010") \
    .withColumn("WeekNumber", f.weekofyear(f.col("InvoiceDate"))) \
    .groupBy("Country", "WeekNumber") \
    .agg(NumInvoices, TotalQuantity, InvoiceValue)
```

```
exSummary_df.coalesce(1) \
    .write \
    .format("parquet") \
    .mode("overwrite") \
    .save("output")
```

```
exSummary_df.sort("Country", "WeekNumber").show()
```

| Country   | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue |
|-----------|------------|-------------|---------------|--------------|
| Australia | 48         | 1           | 107           | 358.25       |
| Australia | 49         | 1           | 214           | 258.9        |
| Australia | 50         | 2           | 133           | 387.95       |
| Austria   | 50         | 2           | 3             | 257.04       |
| Bahrain   | 51         | 1           | 54            | 205.74       |
| Belgium   | 48         | 1           | 528           | 346.1        |
| Belgium   | 50         | 2           | 285           | 625.16       |



- **Windowing Aggregations**

| Country | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue |
|---------|------------|-------------|---------------|--------------|
| EIRE    | 48         | 7           | 2822          | 3147.23      |
| EIRE    | 49         | 5           | 1280          | 3284.1       |
| EIRE    | 50         | 5           | 1184          | 2321.78      |
| EIRE    | 51         | 5           | 95            | 276.84       |
| France  | 48         | 4           | 1299          | 2808.16      |
| France  | 49         | 9           | 2303          | 4527.01      |
| France  | 50         | 6           | 529           | 537.32       |
| France  | 51         | 5           | 847           | 1702.87      |

>>

| Country | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue |
|---------|------------|-------------|---------------|--------------|
| EIRE    | 48         | 7           | 2822          | 3147.23      |
| EIRE    | 49         | 5           | 1280          | 3284.1       |
| EIRE    | 50         | 5           | 1184          | 2321.78      |
| EIRE    | 51         | 5           | 95            | 276.84       |

| Country | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue |
|---------|------------|-------------|---------------|--------------|
| France  | 48         | 4           | 1299          | 2808.16      |
| France  | 49         | 9           | 2303          | 4527.01      |
| France  | 50         | 6           | 529           | 537.32       |
| France  | 51         | 5           | 847           | 1702.87      |

| Country | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue | RunningTotal |
|---------|------------|-------------|---------------|--------------|--------------|
| EIRE    | 48         | 7           | 2822          | 3147.23      | 3147.23      |
| EIRE    | 49         | 5           | 1280          | 3284.1       | 6431.33      |
| EIRE    | 50         | 5           | 1184          | 2321.78      | 8753.11      |
| EIRE    | 51         | 5           | 95            | 276.84       | 9029.95      |

| Country | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue | RunningTotal |
|---------|------------|-------------|---------------|--------------|--------------|
| France  | 48         | 4           | 1299          | 2808.16      | 2808.16      |
| France  | 49         | 9           | 2303          | 4527.01      | 7335.17      |
| France  | 50         | 6           | 529           | 537.32       | 7872.49      |
| France  | 51         | 5           | 847           | 1702.87      | 9575.36      |

1. Identify your partitioning columns.
2. Identify your ordering requirement.
3. Define your window start and end.

## main.py ref. 56

```
from pyspark.sql import SparkSession, Window
from pyspark.sql import functions as f

from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession \
        .builder \
        .appName("Agg Demo") \
        .master("local[2]") \
        .getOrCreate()

    logger = Log4j(spark)

    summary_df = spark.read.parquet("data/summary.parquet")

    running_total_window = Window.partitionBy("Country") \
        .orderBy("WeekNumber") \
        .rowsBetween(-2, Window.currentRow)

    summary_df.withColumn("RunningTotal",
                          f.sum("InvoiceValue").over(running_total_window)) \
        .show()
```

| Country   | WeekNumber | NumInvoices | TotalQuantity | InvoiceValue | RunningTotal      |
|-----------|------------|-------------|---------------|--------------|-------------------|
| Australia | 48         | 1           | 107           | 358.25       | 358.25            |
| Australia | 49         | 1           | 214           | 258.9        | 617.15            |
| Australia | 50         | 2           | 133           | 387.95       | 1005.099999999999 |
| Austria   | 50         | 2           | 3             | 257.04       | 257.04            |
| Bahrain   | 51         | 1           | 54            | 205.74       | 205.74            |
| Belgium   | 48         | 1           | 528           | 346.1        | 346.1             |
| Belgium   | 50         | 2           | 285           | 625.16       | 971.26            |
| Belgium   | 51         | 2           | 942           | 838.65       | 1809.909999999999 |

- **Spark Dataframe JOINS**

- **Dataframe Joins and Column name ambiguity**

| orderDF  |         |            |     |
|----------|---------|------------|-----|
| order_id | prod_id | unit_price | qty |
| 01       | 02      | 350        | 1   |
| 01       | 04      | 580        | 1   |
| 01       | 07      | 320        | 2   |
| 02       | 03      | 450        | 1   |
| 02       | 06      | 220        | 1   |
| 03       | 01      | 195        | 1   |
| 04       | 09      | 270        | 3   |
| 04       | 08      | 410        | 2   |
| 05       | 02      | 350        | 1   |

| productDF |                     |            |     |
|-----------|---------------------|------------|-----|
| prod_id   | name                | list_price | qty |
| 01        | Scroll Mouse        | 250        | 20  |
| 02        | Optical Mouse       | 350        | 20  |
| 03        | Wireless Mouse      | 450        | 50  |
| 04        | Wireless Keyboard   | 580        | 50  |
| 05        | Standard Keyboard   | 360        | 10  |
| 06        | 16 GB Flash Storage | 240        | 100 |
| 07        | 32 GB Flash Storage | 320        | 50  |
| 08        | 64 GB Flash Storage | 430        | 25  |

Left and Right DF

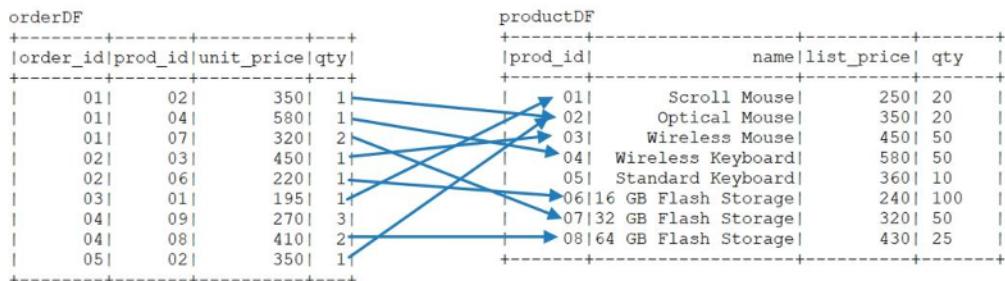
1. Join Condition/Expression : `join_expr = order_df.prod_id == product_df.prod_id`

2. Join Type :

```
order_df.join(product_renamed_df, join_expr, "inner")
```

Always start with Left DF i.e., `order_df` ^ and passing the right DF ^ `product_renamed_df` which contains `join_expression` and join mode which is "inner". (inner join is default you can skip it.)

- Inner :: find matches left table 1 in all elements in table to :: after matching data it combines both DF as like SQL inner function.



`main.py`

```
from pyspark.sql import SparkSession
```

```
from lib.logger import Log4j
```

```
if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .appName("SparkJoins")\
        .master("local[3]")\
        .getOrCreate()
```

```
logger = Log4j(spark)
```

```
orders_list = [
    ("01", "02", 350, 1),
    ("01", "04", 580, 1),
    ("01", "07", 320, 2),
    ("02", "03", 450, 1),
    ("02", "06", 220, 1),
    ("03", "01", 195, 1),
    ("04", "09", 270, 3),
    ("04", "08", 410, 2),
```

```

("05", "02", 350, 1)
]

order_df = spark.createDataFrame(orders_list).toDF("order_id", "prod_id", "unit_price", "qty")

product_list = [
    ("01", "Scroll Mouse", 250, 20),
    ("02", "Optical Mouse", 350, 20),
    ("03", "Wireless Mouse", 450, 50),
    ("04", "Wireless Keyboard", 580, 50),
    ("05", "Standard Keyboard", 360, 10),
    ("06", "16 GB Flash Storage", 240, 100),
    ("07", "32 GB Flash Storage", 320, 50),
    ("08", "64 GB Flash Storage", 430, 25)
]
]

product_df = spark.createDataFrame(product_list).toDF("prod_df", "prod_name", "list_price", "qty")

product_df.show()
order_df.show()

```

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX  
**PYSPARK\_DRIVER\_PYTHON** Environment variable for JNotebook

User variables for innk

| Variable                     | Value                                                            |
|------------------------------|------------------------------------------------------------------|
| HADOOP_HOME                  | E:\pySpark_soft\hadoop-3.2.2                                     |
| JAVA_HOME                    | C:\Program Files\Java\jdk-11                                     |
| OneDrive                     | C:\Users\innk\OneDrive                                           |
| Path                         | C:\Users\innk\AppData\Local\Programs\Python\Python39\Script...   |
| <b>PYSPARK_DRIVER_PYTHON</b> | C:\Users\innk\anaconda3\python.exe                               |
| PYSPARK_PYTHON               | C:\Users\innk\AppData\Local\Programs\Python\Python39\pytho...    |
| PYTHONPATH                   | E:\pySpark_soft\spark_3\python;E:\pySpark_soft\spark_3\python... |
| SPARK_HOME                   | E:\pySpark_soft\spark_3                                          |

Python in worker has different version 3.9 than that in driver 3.10 :: in PyCharm (NEED TO CHANGE

**PYSPARK\_PYTHON** TO 3.10 AND DELETE **DRIVER\_PYTHON**)

Python in worker has different version 3.10 than that in driver 3.9 :: in JNotebook

Note : anaconda running on 3.9 and PyCharm on 3.10

XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

Result:

| prod_df  |         | order_df   |     |
|----------|---------|------------|-----|
| order_id | prod_id | unit_price | qty |
| 01       | 01      | 250        | 20  |
| 02       | 01      | 350        | 20  |
| 03       | 02      | 450        | 50  |
| 04       | 02      | 580        | 50  |
| 05       | 03      | 360        | 10  |
| 06       | 03      | 240        | 100 |
| 07       | 04      | 320        | 50  |
| 08       | 04      | 430        | 25  |
|          | 05      | 195        | 1   |
|          | 05      | 220        | 1   |
|          | 06      | 270        | 3   |
|          | 06      | 320        | 2   |
|          | 07      | 410        | 1   |
|          | 07      | 350        | 1   |
|          | 08      | 400        | 1   |

prod\_df (\* 1<sup>st</sup> column header is prod\_id)

order\_df

Performing **Inner Join** using **innerExpression** on above prod\_df and order\_df

```
join_exp = order_df.prod_id == product_df.prod_id
```

```
order_df.join(product_df, join_exp, "inner")\
.show()
```

| order_id | prod_id | unit_price | qty | prod_id | prod_name           | list_price | qty |
|----------|---------|------------|-----|---------|---------------------|------------|-----|
| 03       | 01      | 195        | 1   | 01      | Scroll Mouse        | 250        | 20  |
| 01       | 02      | 350        | 1   | 02      | Optical Mouse       | 350        | 20  |
| 05       | 02      | 350        | 1   | 02      | Optical Mouse       | 350        | 20  |
| 02       | 03      | 450        | 1   | 03      | Wireless Mouse      | 450        | 50  |
| 01       | 04      | 580        | 1   | 04      | Wireless Keyboard   | 580        | 50  |
| 02       | 06      | 220        | 1   | 06      | 16 GB Flash Storage | 240        | 100 |
| 01       | 07      | 320        | 2   | 07      | 32 GB Flash Storage | 320        | 50  |
| 04       | 08      | 410        | 2   | 08      | 64 GB Flash Storage | 430        | 25  |

- **Column name ambiguity**

```
order_df.join(product_df, join_exp, "inner")\
.select("qty")\
.show()
```

```
pyspark.sql.utils.AnalysisException: Reference 'qty' is ambiguous, could be: qty, qty.
```

We got an error for the selected column name “qty”. because spark doesn't have ambiguity to differ same name column in two different DF. But, when we use select “\*” it shows columns with same names because the spark have its own internal ids that will differ these same column name by that ids.

```
order_df.join(product_df, join_exp, "inner")\
.select("*")\
.show()
```

| order_id | prod_id | unit_price | qty | prod_id | prod_name           | list_price | qty |
|----------|---------|------------|-----|---------|---------------------|------------|-----|
| 03       | 01      | 195        | 1   | 01      | Scroll Mouse        | 250        | 20  |
| 01       | 02      | 350        | 1   | 02      | Optical Mouse       | 350        | 20  |
| 05       | 02      | 350        | 1   | 02      | Optical Mouse       | 350        | 20  |
| 02       | 03      | 450        | 1   | 03      | Wireless Mouse      | 450        | 50  |
| 01       | 04      | 580        | 1   | 04      | Wireless Keyboard   | 580        | 50  |
| 02       | 06      | 220        | 1   | 06      | 16 GB Flash Storage | 240        | 100 |
| 01       | 07      | 320        | 2   | 07      | 32 GB Flash Storage | 320        | 50  |
| 04       | 08      | 410        | 2   | 08      | 64 GB Flash Storage | 430        | 25  |

We have to take care like **rename** that same column names before execution.

```
join_exp = order_df.prod_id == product_df.prod_id
```

```
product_renamed_df = product_df.withColumnRenamed("qty", "reorder_qty")
order_df.join(product_renamed_df, join_exp, "inner")\
    .drop(product_renamed_df.prod_id)\ # we also have same prod_id column in both DF so drop.
    .select("prod_id", "order_id", "prod_name", "unit_price", "qty")\
    .show()
```

| prod_id | order_id | prod_name           | unit_price | qty |
|---------|----------|---------------------|------------|-----|
| 01      | 03       | Scroll Mouse        | 195        | 1   |
| 02      | 01       | Optical Mouse       | 350        | 1   |
| 02      | 05       | Optical Mouse       | 350        | 1   |
| 03      | 02       | Wireless Mouse      | 450        | 1   |
| 04      | 01       | Wireless Keyboard   | 580        | 1   |
| 06      | 02       | 16 GB Flash Storage | 220        | 1   |
| 07      | 01       | 32 GB Flash Storage | 320        | 2   |
| 08      | 04       | 64 GB Flash Storage | 410        | 2   |

- **Outer Joins in DF**

- Outer joins are joins that return matched values and unmatched values from either or both tables.

```
join_exp = order_df.prod_id == product_df.prod_id
```

```
product_renamed_df = product_df.withColumnRenamed("qty", "reorder_qty")
order_df.join(product_renamed_df, join_exp, "outer")\
    .drop(product_renamed_df.prod_id)\
    .select("*")\
    .sort("order_id")\
    .show()
```

| order_id | prod_id | unit_price | qty  | prod_name           | list_price | reorder_qty |
|----------|---------|------------|------|---------------------|------------|-------------|
| null     | null    | null       | null | Standard Keyboard   | 360        | 10          |
| 01       | 04      | 580        | 1    | Wireless Keyboard   | 580        | 50          |
| 01       | 02      | 350        | 1    | Optical Mouse       | 350        | 20          |
| 01       | 07      | 320        | 2    | 32 GB Flash Storage | 320        | 50          |
| 02       | 03      | 450        | 1    | Wireless Mouse      | 450        | 50          |
| 02       | 06      | 220        | 1    | 16 GB Flash Storage | 240        | 100         |
| 03       | 01      | 195        | 1    | Scroll Mouse        | 250        | 20          |
| 04       | 08      | 410        | 2    | 64 GB Flash Storage | 430        | 25          |
| 04       | 09      | 270        | 3    | null                | null       | null        |
| 05       | 02      | 350        | 1    | Optical Mouse       | 350        | 20          |

FULL OUTER JOIN

1. Outer Join – Full Outer
2. Left Join – Left Outer
3. Right Join – Right Outer

```

product_renamed_df = product_df.withColumnRenamed("qty", "reorder_qty")
order_df.join(product_renamed_df, join_exp, "left")\
    .drop(product_renamed_df.prod_id)\n    .select("order_id", "prod_id", "prod_name", "unit_price", "list_price", "qty")\
    .sort("order_id")\
    .show()

```

| order_id | prod_id | prod_name           | unit_price | list_price | qty |
|----------|---------|---------------------|------------|------------|-----|
| 01       | 07      | 32 GB Flash Storage | 320        | 320        | 2   |
| 01       | 04      | Wireless Keyboard   | 580        | 580        | 1   |
| 01       | 02      | Optical Mouse       | 350        | 350        | 1   |
| 02       | 03      | Wireless Mouse      | 450        | 450        | 1   |
| 02       | 06      | 16 GB Flash Storage | 220        | 240        | 1   |
| 03       | 01      | Scroll Mouse        | 195        | 250        | 1   |
| 04       | 09      | null                | 270        | null       | 3   |
| 04       | 08      | 64 GB Flash Storage | 410        | 430        | 2   |
| 05       | 02      | Optical Mouse       | 350        | 350        | 1   |

LEFT OUTER JOIN

```

join_exp = order_df.prod_id == product_df.prod_id\n\nproduct_renamed_df = product_df.withColumnRenamed("qty", "reorder_qty")\norder_df.join(product_renamed_df, join_exp, "left")\
    .drop(product_renamed_df.prod_id)\n    .select("order_id", "prod_id", "prod_name", "unit_price", "list_price", "qty")\
    .withColumn("prod_name", expr("coalesce(prod_name, prod_id)"))\
    .withColumn("list_price", expr("coalesce(list_price, unit_price)"))\
    .sort("order_id")\
    .show()\n#coalesce takes 1st null value replace with 2nd value

```

| order_id | prod_id | prod_name           | unit_price | list_price | qty |
|----------|---------|---------------------|------------|------------|-----|
| 01       | 07      | 32 GB Flash Storage | 320        | 320        | 2   |
| 01       | 04      | Wireless Keyboard   | 580        | 580        | 1   |
| 01       | 02      | Optical Mouse       | 350        | 350        | 1   |
| 02       | 03      | Wireless Mouse      | 450        | 450        | 1   |
| 02       | 06      | 16 GB Flash Storage | 220        | 240        | 1   |
| 03       | 01      | Scroll Mouse        | 195        | 250        | 1   |
| 04       | 09      | 09                  | 270        | 270        | 3   |
| 04       | 08      | 64 GB Flash Storage | 410        | 430        | 2   |
| 05       | 02      | Optical Mouse       | 350        | 350        | 1   |

Using coalesce LEFT JOIN

## • Internals of Spark Join and Shuffle

- Spark JOIN one of the most commonly causes to slowing down your application.

### 1. Shuffle Join

### 2. Broadcast Join

Shuffle JOIN most commonly used.

| id | FL_DATE  | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|----|----------|------------|------------|--------|--------|------------------|------|----------------|
| 00 | 1/1/2000 | DL         |            | 1451   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 01 | 1/1/2000 | DL         |            | 1479   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 02 | 1/1/2000 | DL         |            | 1857   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 03 | 1/1/2000 | DL         |            | 1997   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 04 | 1/1/2000 | DL         |            | 2065   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 05 | 1/1/2000 | US         |            | 2619   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 06 | 1/1/2000 | US         |            | 2821   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 07 | 1/1/2000 | DL         |            | 346    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 08 | 1/1/2000 | DL         |            | 412    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 09 | 1/1/2000 | DL         |            | 299    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 10 | 1/1/2000 | DL         |            | 495    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 11 | 1/1/2000 | DL         |            | 677    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 12 | 1/1/2000 | DL         |            | 251    | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 13 | 1/1/2000 | DL         |            | 1003   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 14 | 1/1/2000 | DL         |            | 1501   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 15 | 1/1/2000 | DL         |            | 1907   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 16 | 1/1/2000 | DL         |            | 2063   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 17 | 1/1/2000 | DL         |            | 2111   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 18 | 1/1/2000 | US         |            | 2632   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 19 | 1/1/2000 | US         |            | 2967   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 20 | 1/1/2000 | DL         |            | 540    | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 21 | 1/1/2000 | DL         |            | 544    | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 22 | 1/1/2000 | DL         |            | 1289   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 23 | 1/1/2000 | DL         |            | 1518   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 24 | 1/1/2000 | DL         |            | 1519   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 25 | 1/1/2000 | DL         |            | 1564   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 26 | 1/1/2000 | DL         |            | 2039   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 27 | 1/1/2000 | DL         |            | 391    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 28 | 1/1/2000 | DL         |            | 423    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 29 | 1/1/2000 | DL         |            | 717    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 30 | 1/1/2000 | DL         |            | 771    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |

| id | CRS_DEP_TIME | DEP_TIME | WHEELS_ON | TAXI_IN | CRS_ARR_TIME | ARR_TIME | CANCELLED | DISTANCE |
|----|--------------|----------|-----------|---------|--------------|----------|-----------|----------|
| 00 | 1115         | 1113     | 1343      | 5       | 1400         | 1348     | 0         | 946      |
| 01 | 1315         | 1311     | 1536      | 7       | 1559         | 1543     | 0         | 946      |
| 02 | 1415         | 1414     | 1642      | 9       | 1721         | 1651     | 0         | 946      |
| 03 | 1715         | 1720     | 1955      | 10      | 2013         | 2005     | 0         | 946      |
| 04 | 2015         | 2010     | 2230      | 10      | 2300         | 2240     | 0         | 946      |
| 05 | 650          | 649      | 956       | 7       | 955          | 1003     | 0         | 946      |
| 06 | 1440         | 1446     | 1713      | 4       | 1738         | 1717     | 0         | 946      |
| 07 | 1740         | 1744     | 1957      | 9       | 2008         | 2006     | 0         | 449      |
| 08 | 1345         | 1345     | 1552      | 9       | 1622         | 1601     | 0         | 449      |
| 09 | 1245         | 1245     | 1443      | 5       | 1455         | 1448     | 0         | 712      |
| 10 | 2035         | 2035     | 2226      | 9       | 2241         | 2235     | 0         | 712      |
| 11 | 710          | 710      | 940       | 7       | 925          | 947      | 0         | 712      |
| 12 | 2040         | 2100     | 2235      | 7       | 2233         | 2242     | 0         | 576      |
| 13 | 1635         | 1838     | 2020      | 12      | 1832         | 2032     | 0         | 576      |
| 14 | 1400         | 1435     | 1623      | 12      | 1634         | 1635     | 0         | 576      |
| 15 | 530          | 530      | 716       | 4       | 723          | 720      | 0         | 576      |
| 16 | 1250         |          |           |         | 1449         |          | 1         | 576      |
| 17 | 1845         | 1855     | 2041      | 9       | 2046         | 2050     | 0         | 576      |
| 18 | 710          | 710      |           |         | 965          |          | 0         | 576      |
| 19 | 1700         | 1700     | 1845      | 6       | 1851         | 1851     | 0         | 576      |
| 20 | 655          | 652      | 1052      | 5       | 1104         | 1057     | 0         | 294      |
| 21 | 2135         | 2125     | 2219      | 5       | 2242         | 2224     | 0         | 294      |

^^^^^^^^^^^^^^^^^^^^^^^^^ Partition of above DF ^^^^^^^^^^^^^^^^^^^^^^

| FL_DATE | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|---------|------------|------------|--------|--------|------------------|------|----------------|
| 00      | 1/1/2000   | DL         | 677    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 01      | 1/1/2000   | DL         | 251    | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 02      | 1/1/2000   | DL         | 1003   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 03      | 1/1/2000   | DL         | 1501   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 04      | 1/1/2000   | DL         | 1907   | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 05      | 1/1/2000   | US         | 2619   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 06      | 1/1/2000   | US         | 2821   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 07      | 1/1/2000   | DL         | 346    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 08      | 1/1/2000   | DL         | 412    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 09      | 1/1/2000   | DL         | 299    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 10      | 1/1/2000   | DL         | 495    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |

| FL_DATE | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|---------|------------|------------|--------|--------|------------------|------|----------------|
| 21      | 1/1/2000   | DL         | 544    | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 22      | 1/1/2000   | DL         | 1289   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 23      | 1/1/2000   | DL         | 1515   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 24      | 1/1/2000   | DL         | 1519   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 25      | 1/1/2000   | DL         | 1564   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 26      | 1/1/2000   | DL         | 2039   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 27      | 1/1/2000   | DL         | 391    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 28      | 1/1/2000   | DL         | 423    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 29      | 1/1/2000   | DL         | 717    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 30      | 1/1/2000   | DL         | 771    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |

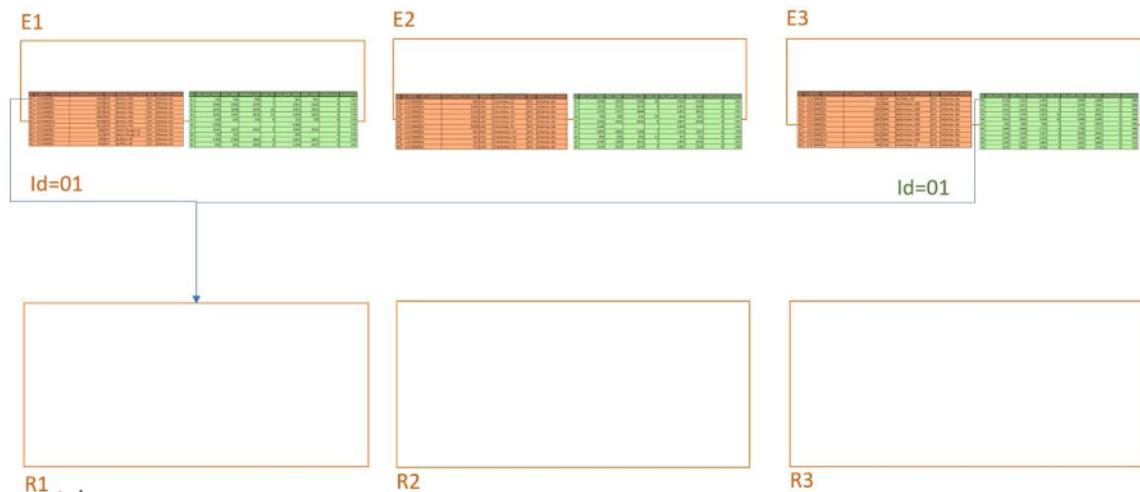
| FL_DATE | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|---------|------------|------------|--------|--------|------------------|------|----------------|
| 31      | 1/1/2000   | DL         | 777    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 32      | 1/1/2000   | DL         | 251    | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 33      | 1/1/2000   | DL         | 1003   | BMI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 34      | 1/1/2000   | DL         | 1501   | BMI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 35      | 1/1/2000   | DL         | 1907   | BMI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 36      | 1/1/2000   | DL         | 2063   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 37      | 1/1/2000   | DL         | 2619   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 38      | 1/1/2000   | DL         | 2821   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 39      | 1/1/2000   | DL         | 346    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 40      | 1/1/2000   | DL         | 412    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 41      | 1/1/2000   | DL         | 299    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 42      | 1/1/2000   | DL         | 495    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |

| FL_DATE | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|---------|------------|------------|--------|--------|------------------|------|----------------|
| 43      | 1/1/2000   | DL         | 1451   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 44      | 1/1/2000   | DL         | 1479   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 45      | 1/1/2000   | DL         | 1857   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 46      | 1/1/2000   | DL         | 1997   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 47      | 1/1/2000   | DL         | 2065   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 48      | 1/1/2000   | US         | 2619   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 49      | 1/1/2000   | US         | 2821   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 50      | 1/1/2000   | DL         | 346    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 51      | 1/1/2000   | DL         | 412    | BTR    | Baton Rouge, LA  | ATL  | Atlanta, GA    |
| 52      | 1/1/2000   | DL         | 299    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 53      | 1/1/2000   | DL         | 495    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |

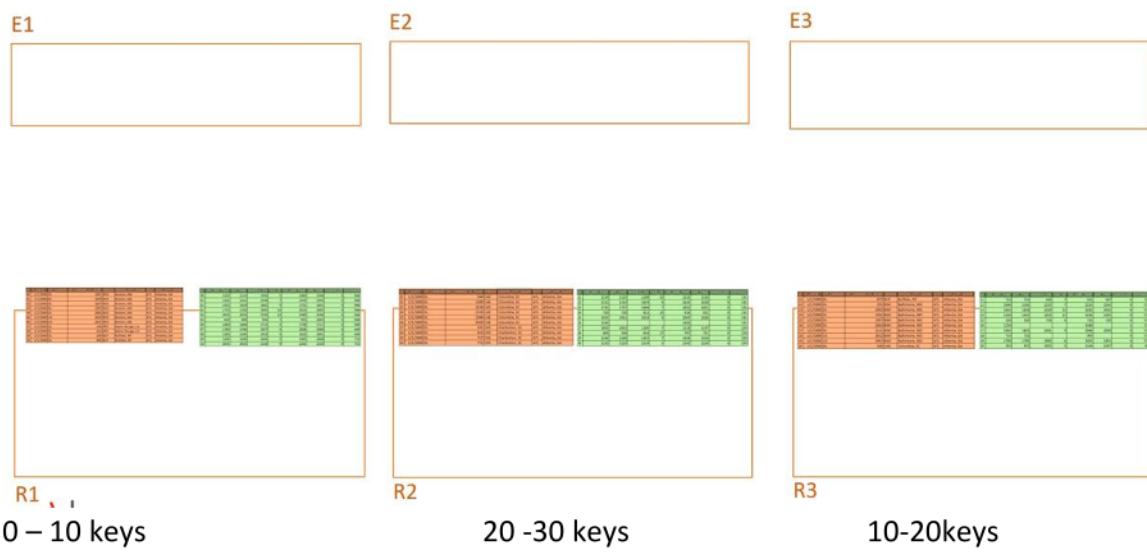
| FL_DATE | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|---------|------------|------------|--------|--------|------------------|------|----------------|
| 54      | 1/1/2000   | DL         | 544    | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 55      | 1/1/2000   | DL         | 1289   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 56      | 1/1/2000   | DL         | 1515   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 57      | 1/1/2000   | DL         | 1519   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 58      | 1/1/2000   | DL         | 1564   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 59      | 1/1/2000   | DL         | 2039   | CAE    | Columbia, SC     | ATL  | Atlanta, GA    |
| 60      | 1/1/2000   | DL         | 391    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 61      | 1/1/2000   | DL         | 423    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 62      | 1/1/2000   | DL         | 717    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |
| 63      | 1/1/2000   | DL         | 771    | CHS    | Charleston, SC   | ATL  | Atlanta, GA    |

| FL_DATE | OP_CARRIER | OP_CARRIER | FL_NUM | ORIGIN | ORIGIN_CITY_NAME | DEST | DEST_CITY_NAME |
|---------|------------|------------|--------|--------|------------------|------|----------------|
| 64      | 1/1/2000   | DL         | 777    | BUF    | Buffalo, NY      | ATL  | Atlanta, GA    |
| 65      | 1/1/2000   | DL         | 251    | BWI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 66      | 1/1/2000   | DL         | 1003   | BMI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 67      | 1/1/2000   | DL         | 1501   | BMI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 68      | 1/1/2000   | DL         | 1907   | BMI    | Baltimore, MD    | ATL  | Atlanta, GA    |
| 69      | 1/1/2000   | DL         | 2063   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 70      | 1/1/2000   | DL         | 2619   | BOS    | Boston, MA       | ATL  | Atlanta, GA    |
| 71</td  |            |            |        |        |                  |      |                |

- Map exchange (map reduce) :: Its like buffer at the executor >> map exchange pick these records and send to map reduce exchange >> Reduce exchange going to collect all record for same key.



here shuffled for 3 partitions and all partitions are processed in parallel.



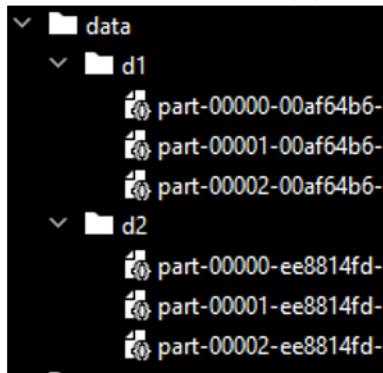
This is shuffle method is main reason that JOIN reduces the performance of app.

Tunning of JOIN operation Is all about optimization of shuffle operation.

| <b>R1</b>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      | <b>R2</b> |            |            |            |         |                  |                  |                |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------|------------|------------|------------|---------|------------------|------------------|----------------|----------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|------|-----|------------|-----|-------------|----|----------|----|--|-----|-----|-----------------|-----|-------------|----|----------|----|--|-----|-----|-----------------|-----|-------------|----|----------|----|--|-----|-----|-------------|-----|-------------|----|----------|----|--|-----|-----|-------------|-----|-------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|----|-----|----------|----------|-----------|---------|--------------|----------|-----------|----------|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|----|------|------|---|---|-----|----|------|------|------|----|------|------|---|---|-----|----|-----|-----|-----|---|-----|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|----|------|------|------|---|------|------|---|---|-----|
| <table border="1"> <thead> <tr> <th>ID</th><th>FL_DATE</th><th>OP_CARRIER</th><th>OP_CARRIER</th><th>FL_NUM</th><th>ORIGIN</th><th>ORIGIN_CITY_NAME</th><th>DEST</th><th>DEST_CITY_NAME</th></tr> </thead> <tbody> <tr><td>00</td><td>1/1/2000</td><td>DL</td><td></td><td>1451</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>01</td><td>1/1/2000</td><td>DL</td><td></td><td>1476</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>02</td><td>1/1/2000</td><td>DL</td><td></td><td>1857</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>03</td><td>1/1/2000</td><td>DL</td><td></td><td>1997</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>04</td><td>1/1/2000</td><td>DL</td><td></td><td>2065</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>05</td><td>1/1/2000</td><td>US</td><td></td><td>2619</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>06</td><td>1/1/2000</td><td>US</td><td></td><td>2621</td><td>BOS</td><td>Boston, MA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>07</td><td>1/1/2000</td><td>DL</td><td></td><td>346</td><td>BTR</td><td>Baton Rouge, LA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>08</td><td>1/1/2000</td><td>DL</td><td></td><td>412</td><td>BTR</td><td>Baton Rouge, LA</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>09</td><td>1/1/2000</td><td>DL</td><td></td><td>299</td><td>BUF</td><td>Buffalo, NY</td><td>ATL</td><td>Atlanta, GA</td></tr> <tr><td>10</td><td>1/1/2000</td><td>DL</td><td></td><td>495</td><td>BUF</td><td>Buffalo, NY</td><td>ATL</td><td>Atlanta, GA</td></tr> </tbody> </table> | ID        | FL_DATE    | OP_CARRIER | OP_CARRIER | FL_NUM  | ORIGIN           | ORIGIN_CITY_NAME | DEST           | DEST_CITY_NAME | 00 | 1/1/2000 | DL |  | 1451 | BOS | Boston, MA | ATL | Atlanta, GA | 01 | 1/1/2000 | DL |  | 1476 | BOS | Boston, MA | ATL | Atlanta, GA | 02 | 1/1/2000 | DL |  | 1857 | BOS | Boston, MA | ATL | Atlanta, GA | 03 | 1/1/2000 | DL |  | 1997 | BOS | Boston, MA | ATL | Atlanta, GA | 04 | 1/1/2000 | DL |  | 2065 | BOS | Boston, MA | ATL | Atlanta, GA | 05 | 1/1/2000 | US |  | 2619 | BOS | Boston, MA | ATL | Atlanta, GA | 06 | 1/1/2000 | US |  | 2621 | BOS | Boston, MA | ATL | Atlanta, GA | 07 | 1/1/2000 | DL |  | 346 | BTR | Baton Rouge, LA | ATL | Atlanta, GA | 08 | 1/1/2000 | DL |  | 412 | BTR | Baton Rouge, LA | ATL | Atlanta, GA | 09 | 1/1/2000 | DL |  | 299 | BUF | Buffalo, NY | ATL | Atlanta, GA | 10 | 1/1/2000 | DL |  | 495 | BUF | Buffalo, NY | ATL | Atlanta, GA | <table border="1"> <thead> <tr> <th>ID</th><th>CRS</th><th>DEP_TIME</th><th>DEP_TIME</th><th>WHEELS_ON</th><th>TAXI_IN</th><th>CRS_ARR_TIME</th><th>ARR_TIME</th><th>CANCELLED</th><th>DISTANCE</th></tr> </thead> <tbody> <tr><td>00</td><td>1115</td><td>1113</td><td>1243</td><td>5</td><td>1400</td><td>1348</td><td>0</td><td>0</td><td>946</td></tr> <tr><td>01</td><td>1315</td><td>1311</td><td>1536</td><td>7</td><td>1559</td><td>1543</td><td>0</td><td>0</td><td>945</td></tr> <tr><td>02</td><td>1415</td><td>1414</td><td>1642</td><td>9</td><td>1721</td><td>1651</td><td>0</td><td>0</td><td>945</td></tr> <tr><td>03</td><td>1715</td><td>1720</td><td>1955</td><td>10</td><td>2013</td><td>2005</td><td>0</td><td>0</td><td>945</td></tr> <tr><td>04</td><td>2015</td><td>2010</td><td>2230</td><td>10</td><td>2800</td><td>2240</td><td>0</td><td>0</td><td>945</td></tr> <tr><td>05</td><td>659</td><td>649</td><td>956</td><td>7</td><td>955</td><td>1003</td><td>0</td><td>0</td><td>946</td></tr> <tr><td>06</td><td>1440</td><td>1446</td><td>1713</td><td>4</td><td>1738</td><td>1717</td><td>0</td><td>0</td><td>946</td></tr> <tr><td>07</td><td>1740</td><td>1744</td><td>1957</td><td>9</td><td>2008</td><td>2006</td><td>0</td><td>0</td><td>449</td></tr> <tr><td>08</td><td>1345</td><td>1345</td><td>1552</td><td>9</td><td>1622</td><td>1601</td><td>0</td><td>0</td><td>449</td></tr> <tr><td>09</td><td>1245</td><td>1245</td><td>1443</td><td>5</td><td>1455</td><td>1448</td><td>0</td><td>0</td><td>712</td></tr> <tr><td>10</td><td>2035</td><td>2035</td><td>2226</td><td>9</td><td>2241</td><td>2235</td><td>0</td><td>0</td><td>712</td></tr> </tbody> </table> | ID | CRS | DEP_TIME | DEP_TIME | WHEELS_ON | TAXI_IN | CRS_ARR_TIME | ARR_TIME | CANCELLED | DISTANCE | 00 | 1115 | 1113 | 1243 | 5 | 1400 | 1348 | 0 | 0 | 946 | 01 | 1315 | 1311 | 1536 | 7 | 1559 | 1543 | 0 | 0 | 945 | 02 | 1415 | 1414 | 1642 | 9 | 1721 | 1651 | 0 | 0 | 945 | 03 | 1715 | 1720 | 1955 | 10 | 2013 | 2005 | 0 | 0 | 945 | 04 | 2015 | 2010 | 2230 | 10 | 2800 | 2240 | 0 | 0 | 945 | 05 | 659 | 649 | 956 | 7 | 955 | 1003 | 0 | 0 | 946 | 06 | 1440 | 1446 | 1713 | 4 | 1738 | 1717 | 0 | 0 | 946 | 07 | 1740 | 1744 | 1957 | 9 | 2008 | 2006 | 0 | 0 | 449 | 08 | 1345 | 1345 | 1552 | 9 | 1622 | 1601 | 0 | 0 | 449 | 09 | 1245 | 1245 | 1443 | 5 | 1455 | 1448 | 0 | 0 | 712 | 10 | 2035 | 2035 | 2226 | 9 | 2241 | 2235 | 0 | 0 | 712 |
| ID                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | FL_DATE   | OP_CARRIER | OP_CARRIER | FL_NUM     | ORIGIN  | ORIGIN_CITY_NAME | DEST             | DEST_CITY_NAME |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 00                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 1451       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 01                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 1476       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 02                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 1857       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 03                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 1997       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 04                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 2065       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 05                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | US         |            | 2619       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 06                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | US         |            | 2621       | BOS     | Boston, MA       | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 07                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 346        | BTR     | Baton Rouge, LA  | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 08                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 412        | BTR     | Baton Rouge, LA  | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 09                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 299        | BUF     | Buffalo, NY      | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 10                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1/1/2000  | DL         |            | 495        | BUF     | Buffalo, NY      | ATL              | Atlanta, GA    |                |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| ID                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | CRS       | DEP_TIME   | DEP_TIME   | WHEELS_ON  | TAXI_IN | CRS_ARR_TIME     | ARR_TIME         | CANCELLED      | DISTANCE       |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 00                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1115      | 1113       | 1243       | 5          | 1400    | 1348             | 0                | 0              | 946            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 01                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1315      | 1311       | 1536       | 7          | 1559    | 1543             | 0                | 0              | 945            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 02                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1415      | 1414       | 1642       | 9          | 1721    | 1651             | 0                | 0              | 945            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 03                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1715      | 1720       | 1955       | 10         | 2013    | 2005             | 0                | 0              | 945            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 04                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 2015      | 2010       | 2230       | 10         | 2800    | 2240             | 0                | 0              | 945            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 05                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 659       | 649        | 956        | 7          | 955     | 1003             | 0                | 0              | 946            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 06                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1440      | 1446       | 1713       | 4          | 1738    | 1717             | 0                | 0              | 946            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 07                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1740      | 1744       | 1957       | 9          | 2008    | 2006             | 0                | 0              | 449            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 08                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1345      | 1345       | 1552       | 9          | 1622    | 1601             | 0                | 0              | 449            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 09                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 1245      | 1245       | 1443       | 5          | 1455    | 1448             | 0                | 0              | 712            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |
| 10                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             | 2035      | 2035       | 2226       | 9          | 2241    | 2235             | 0                | 0              | 712            |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |      |     |            |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |                 |     |             |    |          |    |  |     |     |             |     |             |    |          |    |  |     |     |             |     |             |                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |    |     |          |          |           |         |              |          |           |          |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |      |      |      |    |      |      |   |   |     |    |     |     |     |   |     |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |    |      |      |      |   |      |      |   |   |     |

R1

## Shuffle JOIN [main.py](#)



Here we have 2 partitions each containing 3 JSON files.

```
from pyspark.sql import SparkSession
from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("ShuffleJoin")\
        .getOrCreate()

    logger = Log4j(spark)

    flight_time_df1 = spark.read.json("data/d1/")
    flight_time_df2 = spark.read.json("data/d2/")

    """following conf insures that we having 3 partitions after the shuffle
       which means having 3 exchanges"""
    spark.conf.set("spark.sql.shuffle.partitions", 3)

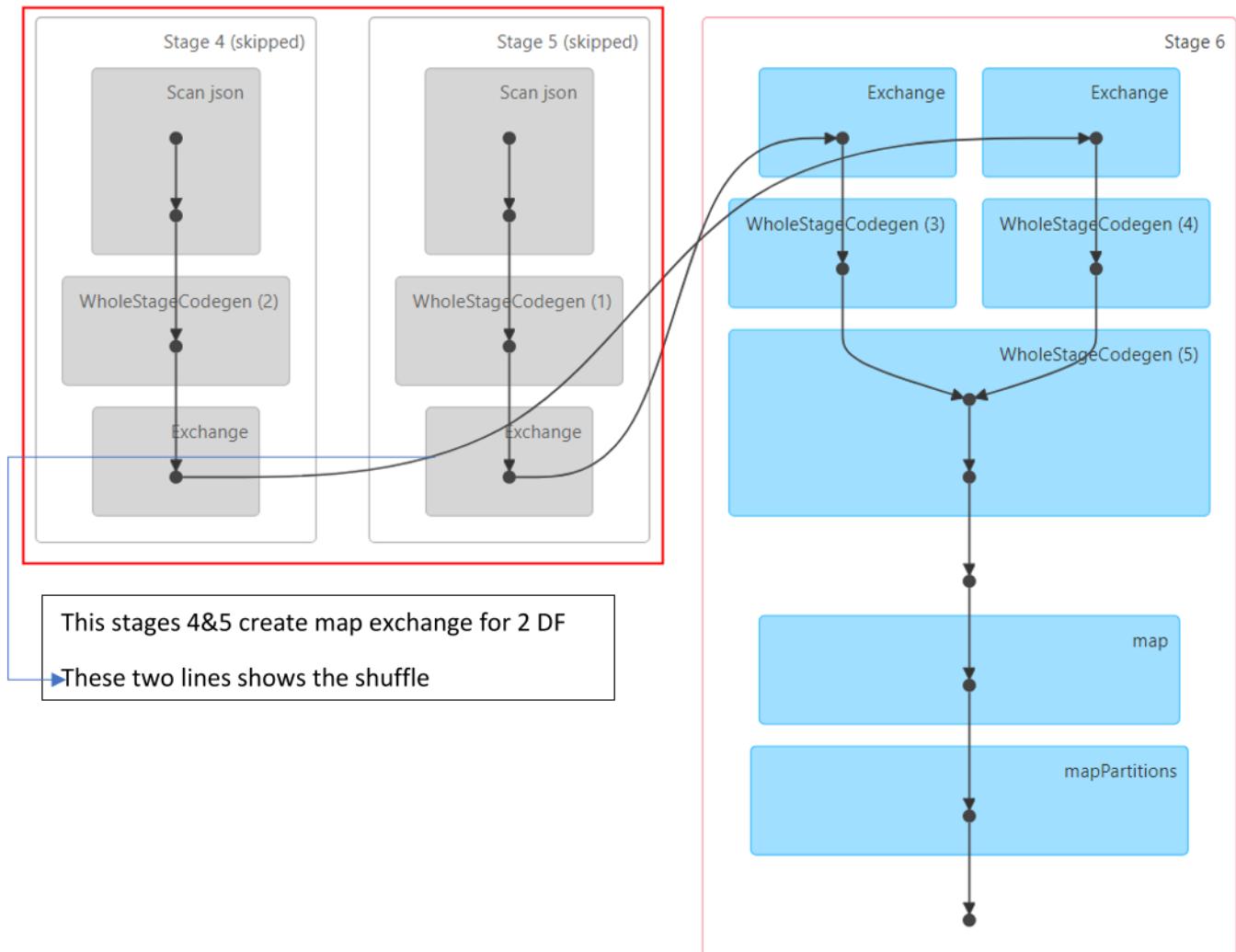
    join_expr = flight_time_df1.id == flight_time_df2.id
    join_df = flight_time_df1.join(flight_time_df2, join_expr, "inner")

    #Dummy action
    join_df.foreach(lambda f: None)
    input("press a key to stop...")
```

localhost:4040/

| Job Id | Description                                                                                                                                                              | Submitted           | Duration | Stages: Succeeded/Total | Tasks (for all stages): Succeeded/Total |
|--------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------|---------------------|----------|-------------------------|-----------------------------------------|
| 4      | foreach at C:\Users\innik\Desktop\Notes\PySpark\Pycharm\ShuffleJoinDemo\main.py:24<br>foreach at C:\Users\innik\Desktop\Notes\PySpark\Pycharm\ShuffleJoinDemo\main.py:24 | 2023/02/28 18:22:50 | 10 s     | 1/1 (2 skipped)         | 3/3 (6 skipped)                         |
| 3      | javaToPython at NativeMethodAccessorsImpl.java:0<br>javaToPython at NativeMethodAccessorsImpl.java:0                                                                     | 2023/02/28 18:22:46 | 4 s      | 1/1                     | 3/3                                     |
| 2      | javaToPython at NativeMethodAccessorsImpl.java:0<br>javaToPython at NativeMethodAccessorsImpl.java:0                                                                     | 2023/02/28 18:22:46 | 3 s      | 1/1                     | 3/3                                     |
| 1      | json at NativeMethodAccessorsImpl.java:0<br>json at NativeMethodAccessorsImpl.java:0                                                                                     | 2023/02/28 18:22:44 | 1 s      | 1/1                     | 3/3                                     |
| 0      | json at NativeMethodAccessorsImpl.java:0<br>json at NativeMethodAccessorsImpl.java:0                                                                                     | 2023/02/28 18:22:41 | 3 s      | 1/1                     | 3/3                                     |

#### DAG Visualization



| Stage Id | Description                                                                        | Submitted           | Duration | Tasks: Succeeded/Total | Input | Output | Shuffle Read | Shuffle Write |
|----------|------------------------------------------------------------------------------------|---------------------|----------|------------------------|-------|--------|--------------|---------------|
| 6        | foreach at C:\Users\innik\Desktop\Notes\PySpark\Pycharm\ShuffleJoinDemo\main.py:24 | 2023/02/28 18:22:50 | 10 s     | 3/3                    |       |        | 24.4 Mill    |               |

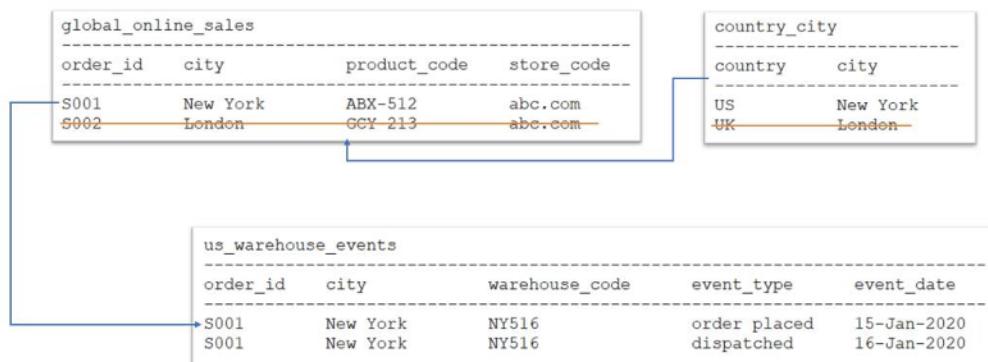
Eliminating this shuffle is best thing in Spark it increases the performance of the app.

## • Optimizing JOINS ref. 61

- There are two JOINING scenarios
  1. Large to Large (This always goes for shuffle operations)
  2. Large to Small (This can be take advantage of broadcast JOIN)

1. Don't code like a novice
2. Shuffle partitions and Parallelism
3. Key Distribution

### 1. Reduce the DF Size before JOIN



1. Don't code like a novice
2. Shuffle partitions and Parallelism
3. Key Distribution

What is the maximum possible parallelism?  
 Executors = 500  
 Shuffle Partitions = 400  
 Unique Keys = 200

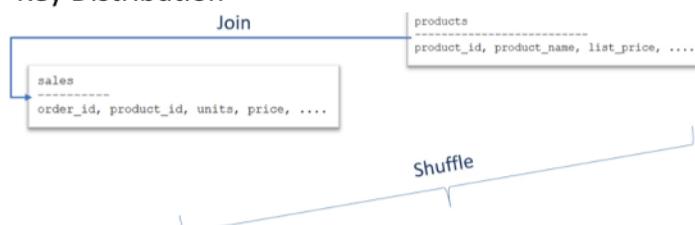
>>

Example,



200 Unique Products  
 Max Shuffle Partitions = 200  
 Max Parallel Tasks = 200

### 3. Key Distribution



Make partitions as per the performance of the operations

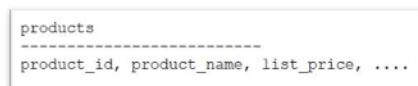


1. Huge Volumes – Filter/Aggregate
2. Parallelism – Shuffles/Executors/Keys
3. Shuffle Distribution – Key Skews

These are the main problems during shuffle joins.

- Large to Small (This can be taken advantage of broadcast JOIN)

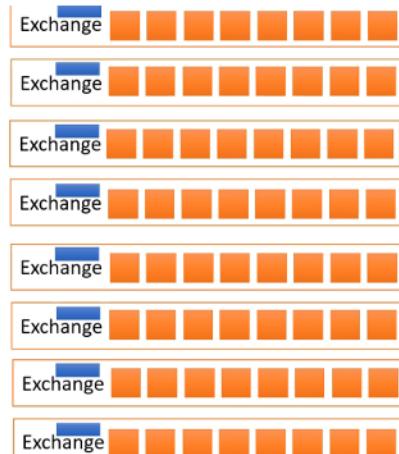
Suppose you have huge DF (sales)



Supposed we used shuffle JOIN



Shuffle



Instead of sending this huge shuffle data to Executors we can dispatch that small size "product" table to all those executors (product table is of 2mb and it will send to 100 executors means now it became 200mb which very small as compared to large which 1 in GB's product also contains all that column which are in large sales DF) this works without shuffle. Smalls means it can me 1GB or 2GB.

Spark-submit options

--driver-memory

--executor-memory

**we know how much memory we have to execute as per that we decide.**

# How to broadcast join?

in many cases spark automatically performed broadcasting

when large DF and Small DF is present.

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import broadcast

from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("ShuffleJoin")\
        .getOrCreate()

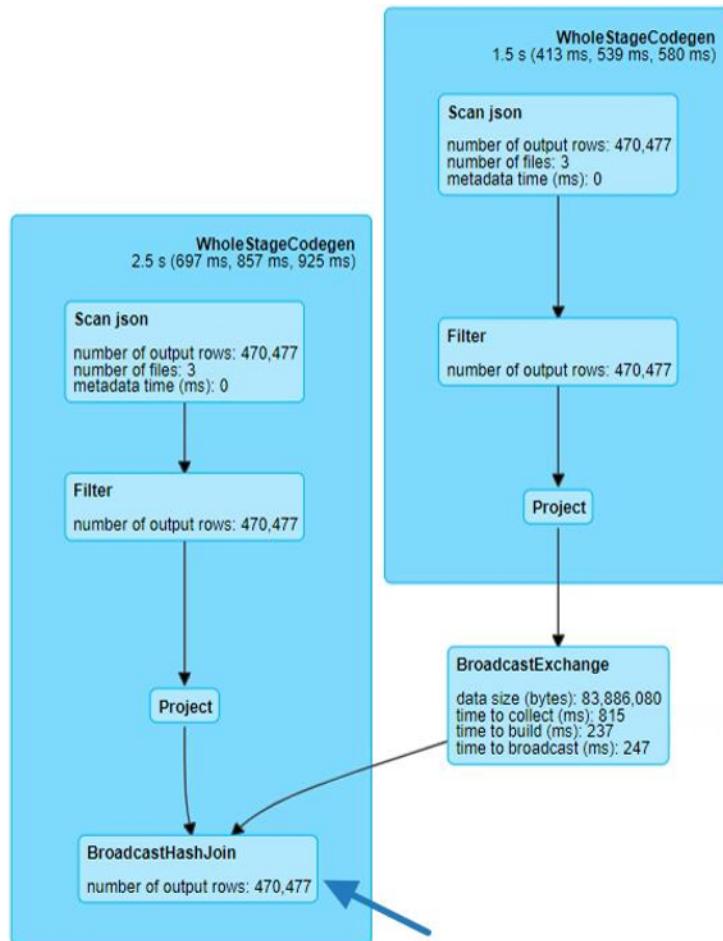
    logger = Log4j(spark)

    flight_time_df1 = spark.read.json("data/d1/")
    flight_time_df2 = spark.read.json("data/d2/")

    """following conf insures that we having 3 partitions after the shuffle
       which means having 3 exchanges"""
    spark.conf.set("spark.sql.shuffle.partitions", 3)

    join_expr = flight_time_df1.id == flight_time_df2.id
    join_df = flight_time_df1.join(broadcast(flight_time_df2), join_expr, "inner")

    #Dummy action
    join_df.foreach(lambda f: None)
    input("press a key to stop...")
```



2 (0c634c72-1c6b-412a-9c03-4461de76676b)

broadcast exchange (runId 0c634c72-1c6b-412a-9c03-4461de76676b)

\$anonfun\$withThreadLocalCaptured\$1 at FutureTask.java:264

- **Implementing Bucket Join**

- Plan in advance | The goal is to avoid shuffle.
- Bucketing is performed only once.
- Bucketing also need shuffle but only when you create a bucket.
- When you have bucket you can join DF **without Shuffle**.
- There no maximum performance is exist. You have to predict it.
- By that prediction you can choice buckets as per the compute power we are selecting 3 buckets in following code of creating **Data Base**.

```
- from pyspark.sql import SparkSession
  from lib.logger import Log4j

  if __name__ == "__main__":
      spark = SparkSession\
          .builder\
          .master("local[3]")\
          .appName("BucketJoin")\
          .enableHiveSupport()\
          .getOrCreate()

      logger = Log4j(spark)

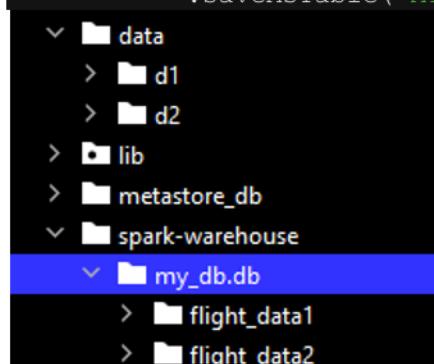
      df1 = spark.read.json("data/d1/")
      df2 = spark.read.json("data/d2/")

      #df1.show()
      #df2.show()

      spark.sql("CREATE DATABASE IF NOT EXISTS MY_DB")
      spark.sql("USE MY_DB")

      df1.coalesce(1).write \
          .bucketBy(3, "id") \
          .mode("overwrite") \
          .saveAsTable("MY_DB.flight_data1")

      df2.coalesce(1).write \
          .bucketBy(3, "id") \
          .mode("overwrite") \
          .saveAsTable("MY_DB.flight_data2")
```



now we have 2 buckets in parquet format.

```
from pyspark.sql import SparkSession
from lib.logger import Log4j

if __name__ == "__main__":
    spark = SparkSession\
        .builder\
        .master("local[3]")\
        .appName("BucketJoin")\
        .enableHiveSupport()\
        .getOrCreate()

    logger = Log4j(spark)

    df1 = spark.read.json("data/d1/")
    df2 = spark.read.json("data/d2/")

    #df1.show()
    #df2.show()

    """spark.sql("CREATE DATABASE IF NOT EXISTS MY_DB")
    spark.sql("USE MY_DB")

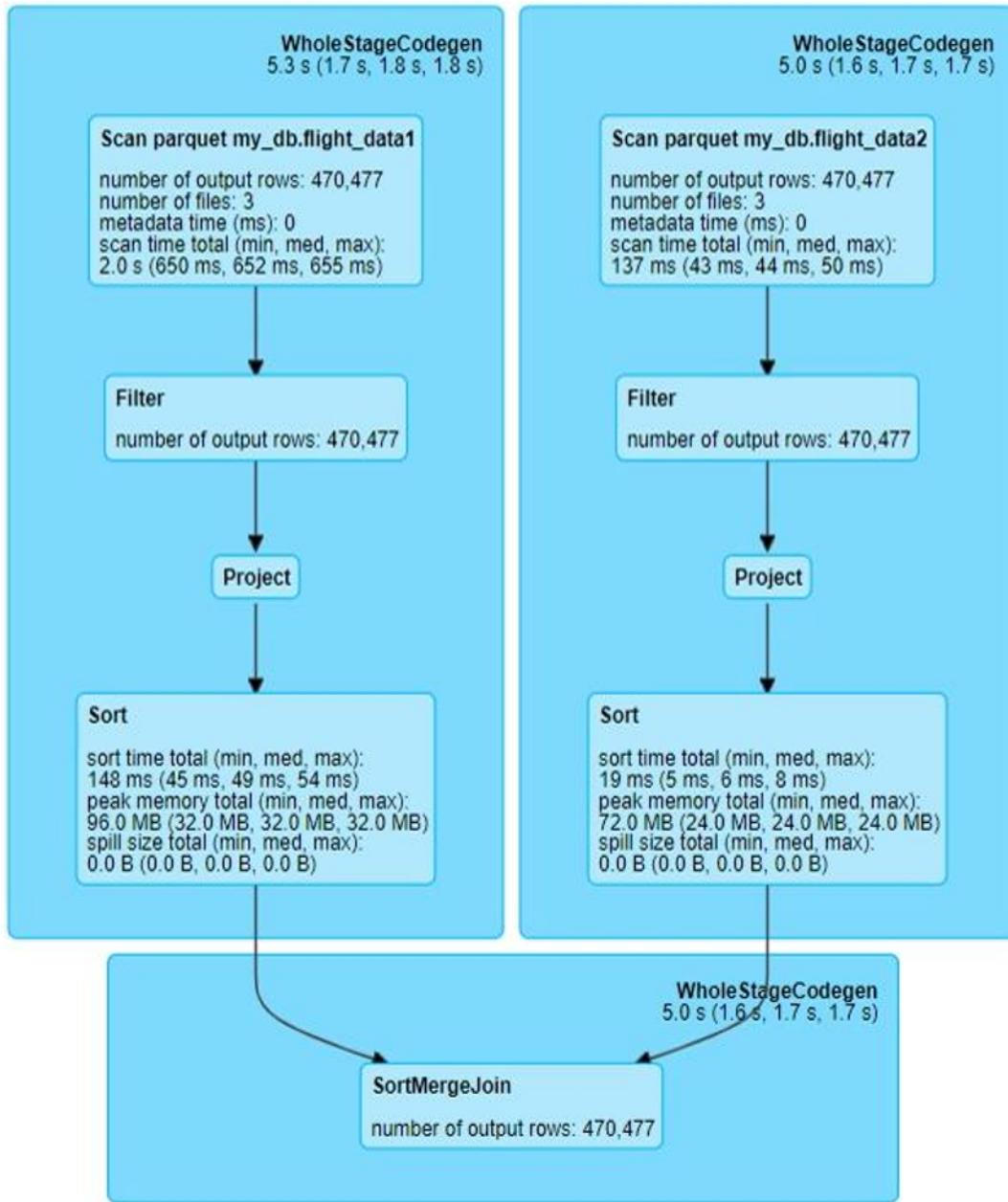
    df1.coalesce(1).write \
        .bucketBy(3, "id") \
        .mode("overwrite") \
        .saveAsTable("MY_DB.flight_data1")

    df2.coalesce(1).write \
        .bucketBy(3, "id") \
        .mode("overwrite") \
        .saveAsTable("MY_DB.flight_data2")"""

    df3 = spark.read.table("MY_DB.flight_data1")
    df4 = spark.read.table("MY_DB.flight_data2")

    spark.conf.set("spark.sql.autoBroadcastJoinThreshold", -1) # tables are small so
    #the spark automatically broadcast it so by this line of code we disable broadcast.
    join_expr = df3.id == df4.id
    join_df = df3.join(df4, join_expr, "inner")

    join_df.collect()
    input("press a key to stop...")
```



Here you can see we performed join with out Shuffle with the help of bucket.

Its possible only when you preplanned and understanding the data and also understand what and how to use it.