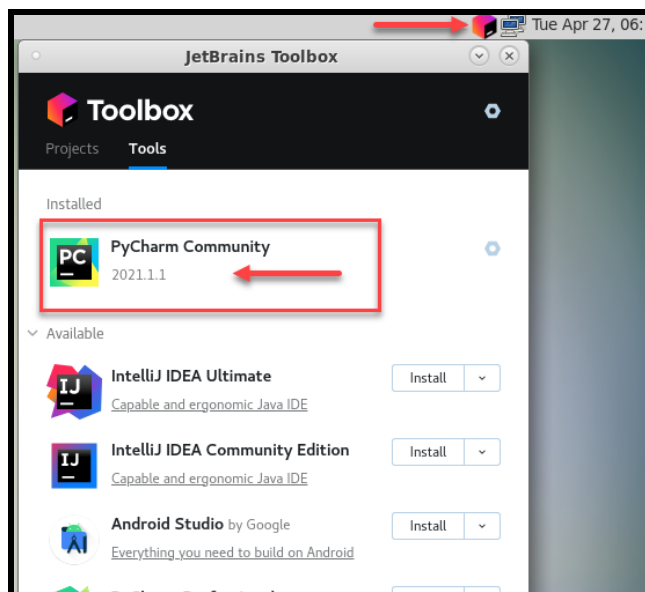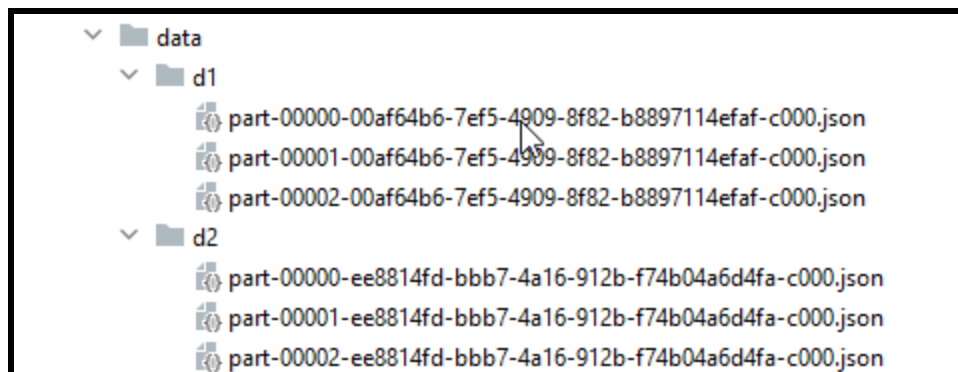# Lab: Shuffle Join

## Introduction

This exercise would help you to explore the internals of the Spark Join so that you can get a better understanding of the plans generated by Catalyst Optimizer for executing Joins in Spark.

## Let's get Started

**Run Pycharm using below screenshot**



We have two data sets named as d1 and d2. Each dataset contains three data files.

**Note:** We created three files purposely because we wanted to make sure we get three partitions when we read them.

Here is the code snippet:

```python
from pyspark.sql import SparkSession

from lib.logger import Log4j

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

    logger = Log4j(spark)

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

Here you can observe that, we are creating a spark session with 3 three parallel threads.

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

    logger = Log4j(spark)

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

We would be reading these two datasets in two different data frames.

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

Now, set the shuffle partition configuration to get three partitions after the shuffle, which means having three reduced exchanges:

```
spark.conf.set("spark.sql.shuffle.partitions", 3)
```

Let's define the join condition.

After performing the join operation we are going to do an inner join. But, Join is a transformation, so nothing is going to actually happen until we take action.

So, let's add a dummy action here. Then we hold the job to look at the Spark UI and understand what's going on there.

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

Let's run it:

```
JoinDemo.py ×
1      from pyspark.sql import SparkSession
2
3      from lib.logger import Log4j
4
5      if    name    == "  main   ":
6    ▶ Run 'JoinDemo'        Ctrl+Shift+F10
7    🐞 Debug 'JoinDemo'
8    🔧 Modify Run Configuration...        10") \
9                  .master("local[3]") \
10                 .getOrCreate()
11
12         logger = Log4j(spark)
13
14         flight_time_df1 = spark.read.json("data/d1/")
15         flight_time_df2 = spark.read.json("data/d2/")
16
17         spark.conf.set("spark.sql.shuffle.partitions", 3)
18
19         join_expr = flight_time_df1.id == flight_time_df2.id
20         join_df = flight_time_df1.join(flight_time_df2, join_expr, "inner")
21
22         join_df.collect()
23         input("press a key to stop...")
```

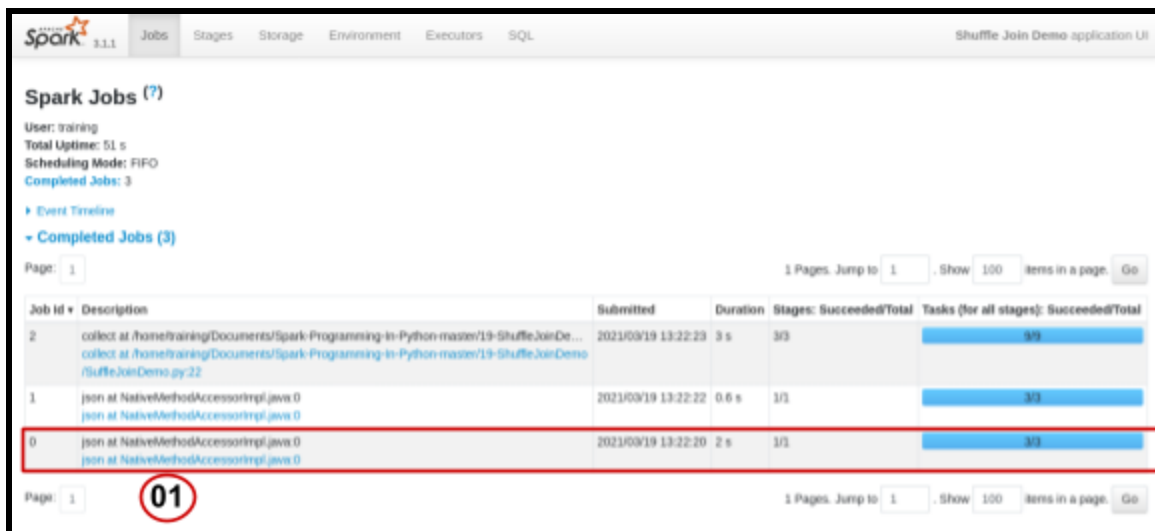# Start the Spark UI.

**Open the browser:** http://localhost:4040

# Jobs

The jobs page shows all the jobs and their run times. Additionally, the page shows how many stages each job contains.

The first job is to read the first data frame.

Reading a data frame is a single-stage and straightforward operation.

So, we have a single stage for this job.



The same happened for the second data frame.

The Join operation happened here.



The join operation is accomplished in three stages.

These two stages are to create a map exchange for the two data frames and lines are indicating the shuffle.

So data moves from the map exchange to the reduce-exchange and rest all is a simple sort-merge join.



The final stage was doing a shuffle read, that thing also happened in three parallel tasks because we configured the shuffle partitions to 3.

Great! So now you know how this Shuffle sort-merge join works under the hood.

Knowing internals can take you too far in investigating join performance, tuning, and improving them.

**Voila!!** We have successfully completed this lab.