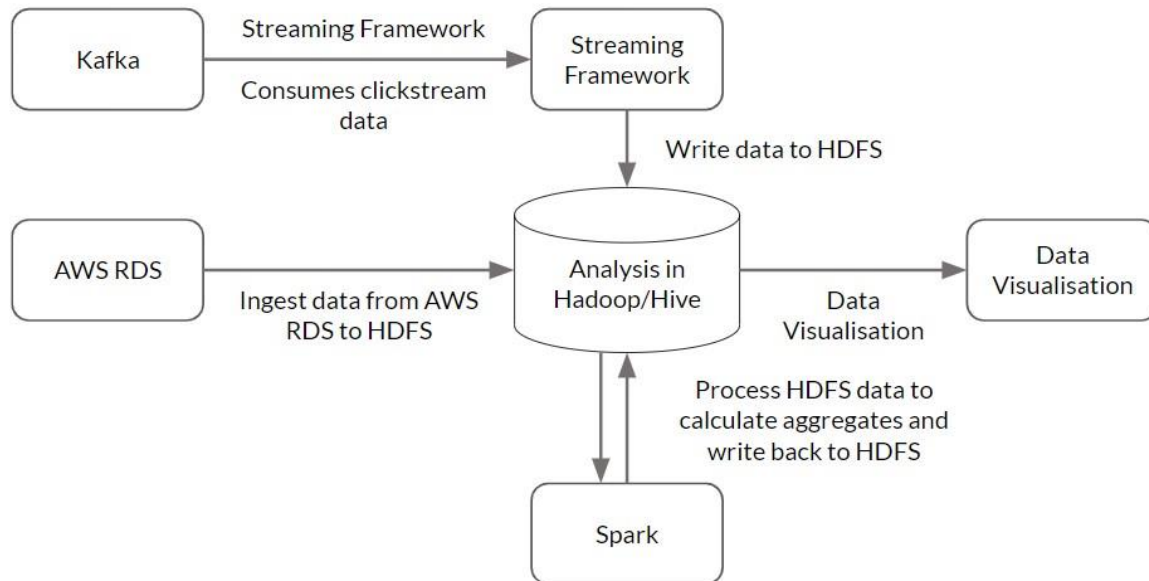# Logic For First Submission

**Overall Architecture:**



According to the given architecture, we have performed the following steps:

In the first step, we ingest stream data from Kafka into a stream processing framework. Then, we consume batch data from AWS RDS into HDFS. Where data aggregation is necessary, we read and write our data from Spark to HDFS. The final aggregated data is stored in Hive tables. The data in Hive tables is later queried to calculate different metrics.

For the first submission (MID), we consume data from Kafka and RDS and load them into Hive managed tables. This process is split into four tasks:

**Task 1: Write a job to consume clickstream data from Kafka and ingest to Hadoop**

Please find commented file "**spark_kafka_to_local.py**"

1) At first, we create a Spark session with application name "Kafka to HDFS"

```python
spark = SparkSession.builder \
    .appName("Kafka To HDFS") \
    .getOrCreate()
spark.sparkContext.setLogLevel('WARN')
```

2) Then we read the incoming clickstream data from Kafka and storing the data in the hdfs. To do this, we use the credentials provided on the platform.

Details of the Kafka Broker are as follows:

- **Bootstrap-server** - 18.211.252.152
- **Port** - 9092
- **Topic** - de-capstone3

The code to read the data is shown below:

```python
df = spark \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", "18.211.252.152:9092") \
    .option("subscribe","de-capstone3") \
    .option("startingOffsets","earliest") \
    .load()
df.printSchema()
```

3) Then we transform the data by dropping a few columns and changing **value** column data type

```python
df = df.withColumn('value_str', df['value'].cast('string').alias('key_str')).drop('value') \
        .drop('key','topic','partition','offset','timestamp','timestampType')
```

4) We write the data to hdfs directory and keep it running until terminated (Write stream in json format).

```python
df1=df.writeStream \
    .outputMode("append") \
    .format("json") \
    .option("truncate","false") \
    .option("path", "/home/hadoop/clickstream_data") \
    .option("checkpointLocation","/home/hadoop/clickstream_checkpoint") \
    .start()

console_df = df \
    .writeStream \
    .format("console") \
    .outputMode("append") \
```

Code execution output:
After running the code, we get the following output

9492-17ae-11eb-adc1-0242ac120002", "is_button_click": "Yes", "is_page_view": "No", "is_scroll_up": "Yes", "is_scroll_down": "No", "timestamp\n": "2020-07-23 23:59:19\n"} |
|{"customer_id": "16929816", "app_version": "1.2.25", "OS_version": "Android", "lat": "-34.8576385", "lon": "-136.639035", "page_id": "e7bc5fb2-1231-11eb-adc1-0242ac120002", "button_id": "f
cba68aa-1231-11eb-adc1-0242ac120002", "is_button_click": "No", "is_page_view": "No", "is_scroll_up": "No", "is_scroll_down": "Yes", "timestamp\n": "2020-05-13 16:35:17\n"} |
|{"customer_id": "41929865", "app_version": "1.1.14", "OS_version": "iOS", "lat": "-34.3613155", "lon": "-137.467833", "page_id": "b328829e-17ae-11eb-adc1-0242ac120002", "button_id": "a95dd
57b-779f-49db-819d-b6960483e554", "is_button_click": "No", "is_page_view": "Yes", "is_scroll_up": "No", "is_scroll_down": "No", "timestamp\n": "2020-04-13 07:58:26\n"}
|{"customer_id": "68940320", "app_version": "1.1.17", "OS_version": "Android", "lat": "81.738239", "lon": "-104.207763", "page_id": "de545711-3914-4450-8c11-b17b8dabb5e1", "button_id": "fcb
a68aa-1231-11eb-adc1-0242ac120002", "is_button_click": "Yes", "is_page_view": "No", "is_scroll_up": "Yes", "is_scroll_down": "No", "timestamp\n": "2020-10-14 00:39:05\n"}
|{"customer_id": "95405928", "app_version": "1.2.10", "OS_version": "Android", "lat": "44.185627", "lon": "78.125488", "page_id": "e7bc5fb2-1231-11eb-adc1-0242ac120002", "button_id": "e1e99
492-17ae-11eb-adc1-0242ac120002", "is_button_click": "No", "is_page_view": "No", "is_scroll_up": "No", "is_scroll_down": "No", "timestamp\n": "2020-07-03 19:15:34\n"}
|{"customer_id": "56235860", "app_version": "1.4.16", "OS_version": "Android", "lat": "-61.178407", "lon": "8.534126", "page_id": "de545711-3914-4450-8c11-b17b8dabb5e1", "button_id": "a95dd
57b-779f-49db-819d-b6960483e554", "is_button_click": "Yes", "is_page_view": "Yes", "is_scroll_up": "Yes", "is_scroll_down": "Yes", "timestamp\n": "2020-06-02 21:32:44\n"}
|{"customer_id": "81758845", "app_version": "4.4.30", "OS_version": "Android", "lat": "78.904984", "lon": "167.315596", "page_id": "e7bc5fb2-1231-11eb-adc1-0242ac120002", "button_id": "a95d
d57b-779f-49db-819d-b6960483e554", "is_button_click": "No", "is_page_view": "No", "is_scroll_up": "No", "is_scroll_down": "No", "timestamp\n": "2020-04-17 19:52:56\n"}
|{"customer_id": "46822067", "app_version": "1.2.26", "OS_version": "Android", "lat": "45.8993985", "lon": "135.895118", "page_id": "de545711-3914-4450-8c11-b17b8dabb5e1", "button_id": "fcb
a68aa-1231-11eb-adc1-0242ac120002", "is_button_click": "Yes", "is_page_view": "Yes", "is_scroll_up": "Yes", "is_scroll_down": "Yes", "timestamp\n": "2020-08-09 12:02:09\n"}|
|{"customer_id": "37623067", "app_version": "3.2.27", "OS_version": "iOS", "lat": "-74.1131505", "lon": "-170.715421", "page_id": "b328829e-17ae-11eb-adc1-0242ac120002", "button_id": "a95dd
57b-779f-49db-819d-b6960483e554", "is_button_click": "No", "is_page_view": "No", "is_scroll_up": "No", "is_scroll_down": "No", "timestamp\n": "2020-01-20 00:06:06\n"}
|{"customer_id": "14571693", "app_version": "3.4.21", "OS_version": "iOS", "lat": "-28.8329055", "lon": "161.132369", "page_id": "e7bc5fb2-1231-11eb-adc1-0242ac120002", "button_id": "fcba68
aa-1231-11eb-adc1-0242ac120002", "is_button_click": "No", "is_page_view": "Yes", "is_scroll_up": "No", "is_scroll_down": "No", "timestamp\n": "2020-08-16 11:20:40\n"}
|{"customer_id": "62605529", "app_version": "2.1.36", "OS_version": "Android", "lat": "44.196239", "lon": "-15.354515", "page_id": "e7bc5fb2-1231-11eb-adc1-0242ac120002", "button_id": "a95d
d57b-779f-49db-819d-b6960483e554", "is_button_click": "Yes", "is_page_view": "Yes", "is_scroll_up": "No", "is_scroll_down": "No", "timestamp\n": "2020-06-24 01:55:32\n"}
+-----------------------------------------------------------------------------------------------------------------------------------+

only showing top 20 rows

^CTraceback (most recent call last):
  File "/home/hadoop/spark_kafka_to_local.py", line 44, in <module>
    console_df.awaitTermination()
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/sql/streaming.py", line 103, in awaitTermination
  File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/java_gateway.py", line 1255, in __call__
  File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/java_gateway.py", line 985, in send_command
  File "/usr/lib/spark/python/lib/py4j-0.10.7-src.zip/py4j/java_gateway.py", line 1152, in send_command
  File "/usr/lib64/python3.7/socket.py", line 589, in readinto
    return self._sock.recv_into(b)
  File "/usr/lib/spark/python/lib/pyspark.zip/pyspark/context.py", line 278, in signal_handler
KeyboardInterrupt
[hadoop@ip-172-31-32-203 ~]$ hadoop fs -ls /home/hadoop/
Found 2 items
drwxr-xr-x   - hadoop hadoop          0 2023-10-23 20:55 /home/hadoop/clickstream_checkpoint
drwxr-xr-x   - hadoop hadoop          0 2023-10-23 20:55 /home/hadoop/clickstream_data
[hadoop@ip-172-31-32-203 ~]$ hadoop fs -ls /home/hadoop/clickstream_data
Found 2 items
drwxr-xr-x   - hadoop hadoop          0 2023-10-23 20:55 /home/hadoop/clickstream_data/_spark_metadata
-rw-r--r--   1 hadoop hadoop    1267706 2023-10-23 20:55 /home/hadoop/clickstream_data/part-00000-e65e868d-b14a-44d5-b5e3-bd36bafbbdd6-c000.json

We now clean the loaded Kafka data to a more structured format and save it to a csv file.

Please find commented file "**spark_local_flatten.py**"

1) Firstly, we get the existing spark session "Kafka to HDFS"

```
spark = SparkSession.builder \
    .appName("Kafka To HDFS") \
    .getOrCreate()
spark.sparkContext.setLogLevel('WARN')
```

2) Then we read the clickstream data that we stored in json format into the dataframe

```
df = spark.read.json('/home/hadoop/clickstream_data/part*')
```

3) We extract the columns from json **value_str** in dataframe and create a new dataframe with new columns

```
df1 = df.select(
    get_json_object(df["value_str"],"$.customer_id").alias("customer_id"),
    get_json_object(df["value_str"],"$.app_version").alias("app_version"),
    get_json_object(df["value_str"],"$.OS_version").alias("OS_version"),
    get_json_object(df["value_str"],"$.lat").alias("lat"),
    get_json_object(df["value_str"],"$.lon").alias("lon"),
    get_json_object(df["value_str"],"$.page_id").alias("page_id"),
    get_json_object(df["value_str"],"$.button_id").alias("button_id"),
    get_json_object(df["value_str"],"$.is_button_click").alias("is_button_click"),
    get_json_object(df["value_str"],"$.is_page_view").alias("is_page_view"),
    get_json_object(df["value_str"],"$.is_scroll_up").alias("is_scroll_up"),
    get_json_object(df["value_str"],"$.is_scroll_down").alias("is_scroll_down"),
    get_json_object(df["value_str"],"$.timestamp\n").alias("click_timestamp")
)
```

4) We then save the new dataframe to a csv file in hdfs directory

```
df1.coalesce(1).write.format('com.databricks.spark.csv').mode('overwrite').save('/home/hadoop/clickstream_data_flatten', header = 'false')
```

The steps to run the python files and the steps is provided in commented file "**step1.pdf**"

Code execution output:
After running the code, we get the following output



**Task 2**: Write a job to ingest the bookings data from AWS RDS to Hadoop

Please find commented file "**step2.pdf**" for the sqoop command
- To do this, we write a sqoop job using the credentials provided on the platform

The RDS connection string and credentials are as follows:

- RDS **Connection String** -

  jdbc:mysql://upgraddetest.cyaielc9bmnf.us-east-1.rds.amazonaws.com/testdatabase

- **Username** - student

- **Password** - STUDENT123
- **Table Name** - bookings

```
sqoop import \
--connect jdbc:mysql://upgraddetest.cyaielc9bmnf.us-east-1.rds.amazonaws.com/testdatabase \
--table bookings \
--username student --password STUDENT123 \
--target-dir /home/hadoop/bookings_data \
-m1
```

Code execution output:
After running the code, we get the following output



**Task 3**: **Create aggregates for finding total bookings date-wise using Spark script**

To prepare aggregates, data is read from HDFS, processed by a processing framework such as Spark and written back to HDFS to create a Hive table for the aggregated data.

Please find commented file "**datewise_bookings_aggregates_spark.py**"
1) To perform aggregation of data, we first create a Spark session with application name "datewise_booking"

```
spark = SparkSession.builder.appName("datewise_booking").getOrCreate()
sc = spark.sparkContext
sc
```

2) Then we read the bookings data into the dataframe

```
df = spark.read.csv("/home/hadoop/bookings_data/part*",inferSchema = True)
```

3) After the above step, we create a new dataframe with renamed columns

```python
df1 = df.withColumnRenamed("_c0","booking_id")\
        .withColumnRenamed("_c1","customer_id") \
        .withColumnRenamed("_c2","driver_id") \
        .withColumnRenamed("_c3","customer_app_version")  \
        .withColumnRenamed("_c4","customer_phone_os_version") \
        .withColumnRenamed("_c5","pickup_lat") \
        .withColumnRenamed("_c6","pickup_lon") \
        .withColumnRenamed("_c7","drop_lat") \
        .withColumnRenamed("_c8","drop_lon") \
        .withColumnRenamed("_c9","pickup_timestamp")  \
        .withColumnRenamed("_c10","drop_timestamp")  \
        .withColumnRenamed("_c11","trip_fare") \
        .withColumnRenamed("_c12","tip_amount")  \
        .withColumnRenamed("_c13","currency_code") \
        .withColumnRenamed("_c14","cab_color")  \
        .withColumnRenamed("_c15","cab_registration_no") \
        .withColumnRenamed("_c16","customer_rating_by_driver")  \
        .withColumnRenamed("_c17","rating_by_customer")  \
        .withColumnRenamed("_c18","passenger_count")
```

4) Then we create aggregations based on date and store the transformed data in a csv file

```python
df2 =  df1.withColumn("booking_date", date_format('pickup_timestamp', "yyyy-MM-dd"))

df2.show(5)
# Date wise aggregation stored in date.
date = df2.select('booking_date').groupBy('booking_date').count()

date.show(10)
date.count()
date.coalesce(1).write.format('com.databricks.spark.csv').mode('overwrite').save('/home/hadoop/datewise_aggreation', header = 'false')
```

The steps to run the python file and the steps are provided in commented file "**step3.pdf**"

Code execution output:
After running the code, we get the following output

## Task 4: Hive managed Table Creation

Please find commented file "**step4.pdf**"
1) Create a Hive-managed table for clickstream data
2) Create a Hive-managed table for bookings data
3) Create a Hive-managed table for aggregated data in Task 3

```
hadoop@ip-172-31-33-247:~

[hadoop@ip-172-31-33-247 ~]$ hive

Logging initialized using configuration in file:/etc/hive/conf.dist/hive-log4j2.properties Async: false
hive> create database if not exists cab_booking ;
OK
Time taken: 0.91 seconds
hive> use cab_booking ;
OK
Time taken: 0.073 seconds
hive> create table if not exists clickstream_data
    > (customer_id string ,app_version string, os_version string,lat string ,lon string ,page_id
    > string,button_id string , is_button_click varchar(3) ,is_page_view varchar(3) ,is_scroll_up
    > varchar(3) ,is_scroll_down varchar(3), click_timestamp string )
    > row format delimited fields terminated by ","
    > location '/home/hadoop/clickstream_data_flatten/';
OK
Time taken: 0.422 seconds
hive> create table if not exists booking_data
    > (booking_id string ,customer_id string ,driver_id string , customer_app_version string,
    > customer_phone_os_version string , pickup_lat double , pickup_lon double, drop_lat double,
    > drop_lon double, pickup_timestamp string , drop_timestamp string ,trip_fare int,
    > tip_amount int, currency_code string ,cab_color string, cab_registration_no string ,
    > customer_rating_by_driver int, rating_by_customer int ,passenger_count int )
    > row format delimited fields terminated by ","
    > location '/home/hadoop/booking_data_csv/';
OK
Time taken: 0.07 seconds
hive> create table if not exists datewise_data
    > (booking_date string , count int)
    > row format delimited fields terminated by ","
    > location '/home/hadoop/datewise_aggreation/';
OK
Time taken: 0.05 seconds
hive>
```

Code execution output:

After creating the Hive tables, we check if the data is properly loaded into them