# CS-523DE-Big-Data-Technology-Project

1. Iyosiyas Workie Mitiku - 614166
2. El Mehdi Korda - 614134

# Agenda

**1**   **Project Description - Spark Streaming  - Data Flow, Kafka Consumer, Kafka Producer**

**2**   **Spark Streaming - Integration with Hive**

**3**   **Further Results**
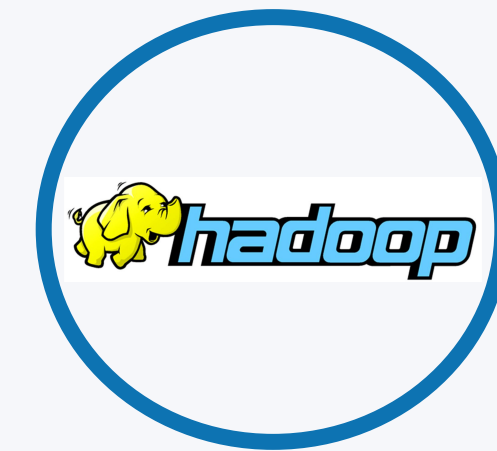
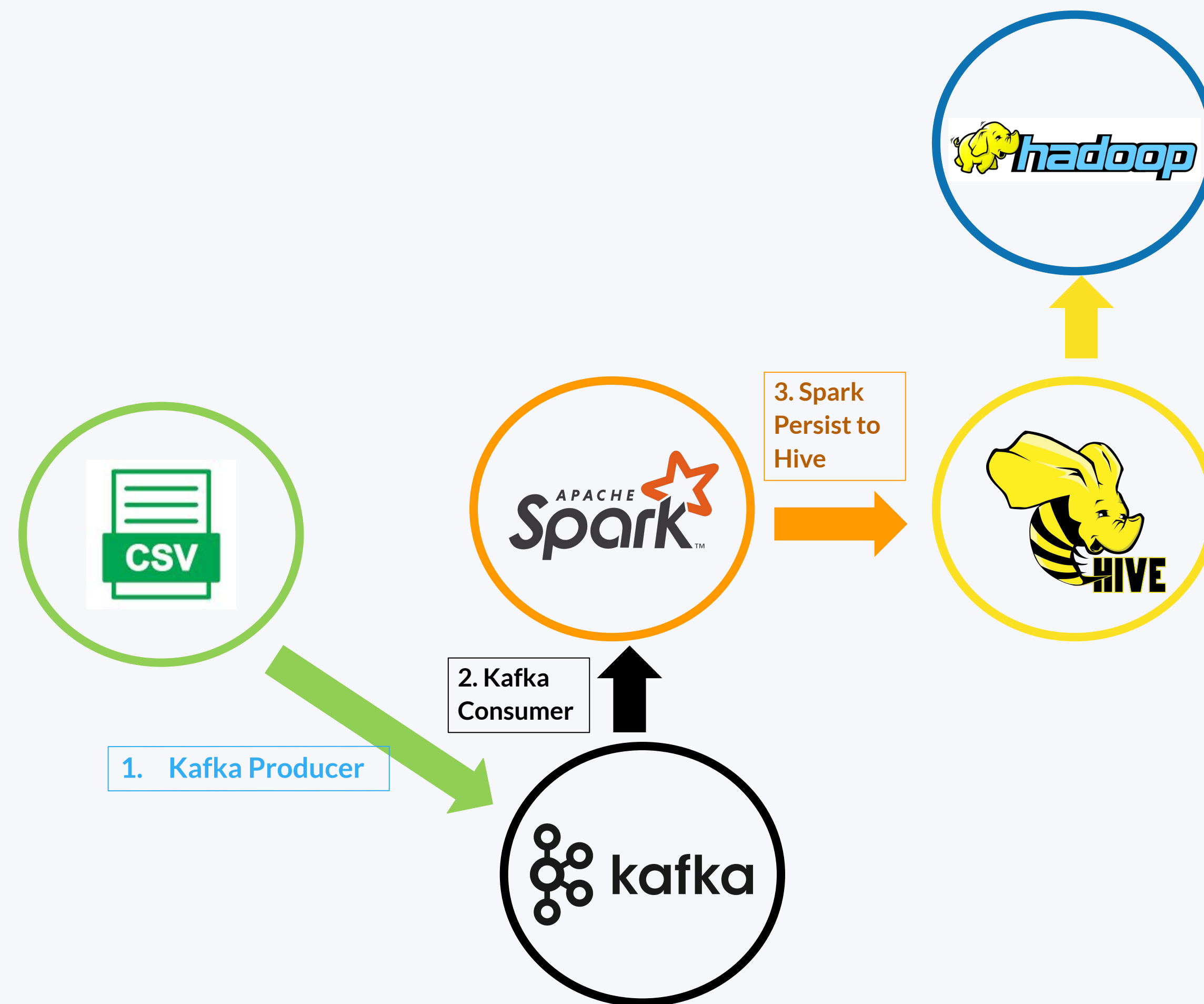# 1.Project Description

## Data

**Technologies**

- player_country
- first_name
- String last_name
- goals_scored
- champions_league_matches_played
- english_league_matches_played
- minutes_played
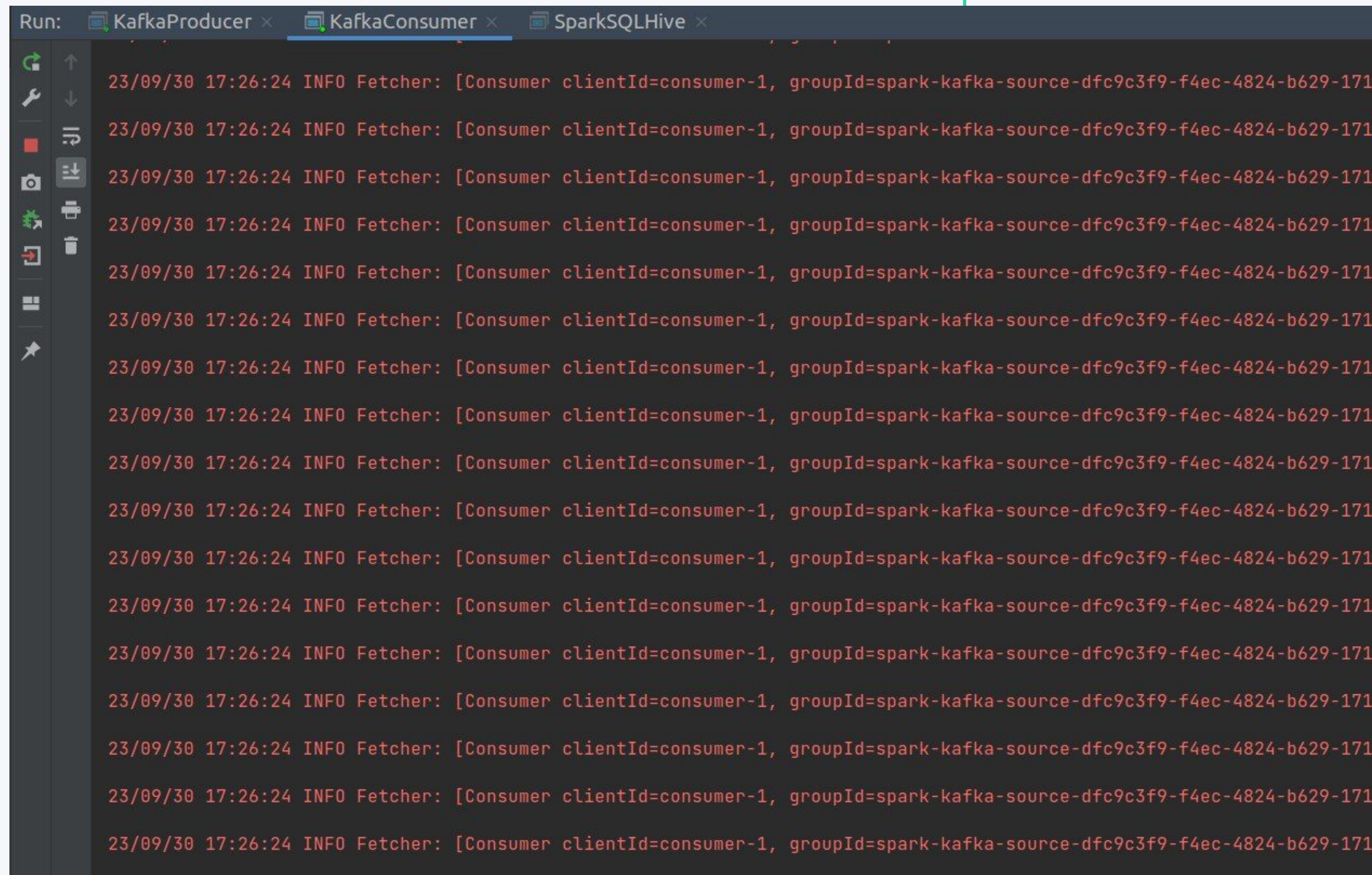- assists
- tackles
- age

# 1. Spark Streaming  - Data Flow



3. Spark Persist to Hive

2. Kafka Consumer

1.    Kafka Producer

# 1.1 Kafka Producer

kafka

- **We have used Kafka as a broker to produce messages that we get from our data source**

- **As seen in the screenshot the Footballer data we get is being produced to a kafka topic we specified**



Run: KafkaProducer ×    KafkaConsumer ×    SparkSQLHive ×

```
message(318, China,Philip,Wyson,2,35,33,735,12,450,20) sent to partition(0), offset(2721) in 2 ms
China,Ferd,Lampert,4,38,24,680,6,406,20
message(320, China,Ferd,Lampert,4,38,24,680,6,406,20) sent to partition(0), offset(2723) in 2 ms
Portugal,Desirae,Reedman,2,25,15,633,78,279,25
message(322, Portugal,Desirae,Reedman,2,25,15,633,78,279,25) sent to partition(0), offset(2725) in 1 ms
Bangladesh,Borg,Dust,9,12,23,540,8,305,28
message(324, Bangladesh,Borg,Dust,9,12,23,540,8,305,28) sent to partition(0), offset(2727) in 2 ms
Brazil,Carolyn,Ovize,13,6,30,645,40,335,26
message(326, Brazil,Carolyn,Ovize,13,6,30,645,40,335,26) sent to partition(0), offset(2729) in 2 ms
Philippines,Whitman,Sudworth,14,14,39,492,63,393,19
message(328, Philippines,Whitman,Sudworth,14,14,39,492,63,393,19) sent to partition(0), offset(2731) in 1 ms
Peru,Ashton,Sabbin,16,17,30,894,11,479,28
message(330, Peru,Ashton,Sabbin,16,17,30,894,11,479,28) sent to partition(0), offset(2733) in 2 ms
Armenia,Felike,Hughs,16,34,30,910,95,266,23
message(332, Armenia,Felike,Hughs,16,34,30,910,95,266,23) sent to partition(0), offset(2735) in 2 ms
Russia,Dick,Casali,19,29,35,156,87,187,23
message(334, Russia,Dick,Casali,19,29,35,156,87,187,23) sent to partition(0), offset(2737) in 2 ms
```

# 1.1 Kafka Consumer

- **As seen in the image we are listening to what we produced**

- **After listening from a topic we are write it to spark**

# 2. Spark Streaming - Integration with Hive



- **As seen in the screenshot we are using Spark to process data stored in Hive.**

- **This integration allows users to take advantage of Spark's performance and scalability for Hive workloads.**

# 2. Spark Streaming - Integration with Hive

- **As seen in the screenshot we are using spark to write sql messages which is then converted to mapreduce then gives the result on the right**

# 3. Hadoop - Further Results Captured



- **The image on the right shows our hadoop which is the underlying storage used to store our data**

# Thank you