```python
In [1]:   1  from pyspark.sql import SparkSession
          2  import pyspark.sql.functions as F
          3  from pyspark.sql.types import *
          4
          5  spark = SparkSession\
          6      .builder\
          7      .appName("chapter-21-stream-kafka")\
          8      .getOrCreate()
          9
         10  import os
         11  SPARK_BOOK_DATA_PATH = os.environ['SPARK_BOOK_DATA_PATH']
```

# Kafka

https://spark.apache.org/docs/2.4.0/structured-streaming-kafka-integration.html#deploying
(https://spark.apache.org/docs/2.4.0/structured-streaming-kafka-integration.html#deploying)

```
 ./bin/spark-submit --packages org.apache.spark:spark-sql-kafka-0-
10_2.11:2.4.0
```

## setup ¶

steps to run kafka and create topic (https://github.com/wgong/py4kids/blob/master/lesson-71-
kafka/Calories-Alert-Kafka/kafka.README.md)

```python
In [2]:   1  # Subscribe to 1 topic
          2  streaming = spark.readStream.format("kafka")\
          3    .option("kafka.bootstrap.servers", "localhost:9092")\
          4    .option("subscribe", "Hello-Kafka")\
          5    .load()
```

```python
In [3]:   1  streaming.printSchema()
```

```
root
 |-- key: binary (nullable = true)
 |-- value: binary (nullable = true)
 |-- topic: string (nullable = true)
 |-- partition: integer (nullable = true)
 |-- offset: long (nullable = true)
 |-- timestamp: timestamp (nullable = true)
 |-- timestampType: integer (nullable = true)
```

### write to memory for test

```
In [4]:  1  streaming.selectExpr("CAST(key AS STRING)", "CAST(value AS STRING)")
         2    .writeStream\
         3    .queryName("test_transform")\
         4    .format("memory")\
         5    .outputMode("append")\
         6    .start()
```

Out[4]: `<pyspark.sql.streaming.StreamingQuery at 0x7fe08da2cac8>`

```
In [8]:  1  spark.sql("select * from test_transform").show(truncate=False)
```

```
+----+--------------------------------------------+
|key |value                                       |
+----+--------------------------------------------+
|null|kafka is a distributed pub/sub message broker|
|null|spark is a distributed big-data platform    |
+----+--------------------------------------------+
```

```
In [9]:  1  spark.sql("select * from test_transform").show(truncate=False)
```

```
+----+--------------------------------------------+
|key |value                                       |
+----+--------------------------------------------+
|null|kafka is a distributed pub/sub message broker|
|null|spark is a distributed big-data platform    |
|null|spark also has mllib for machine learning   |
|null|databricks is the company behind spark      |
+----+--------------------------------------------+
```

```
In [10]:  1  spark.sql("select * from test_transform").show(truncate=False)
```

```
+----+--------------------------------------------+
|key |value                                       |
+----+--------------------------------------------+
|null|kafka is a distributed pub/sub message broker|
|null|spark is a distributed big-data platform    |
|null|spark also has mllib for machine learning   |
|null|databricks is the company behind spark      |
|null|tensorflow 2.0 was released last week       |
|null|tensorflow.js is very interesting           |
+----+--------------------------------------------+
```
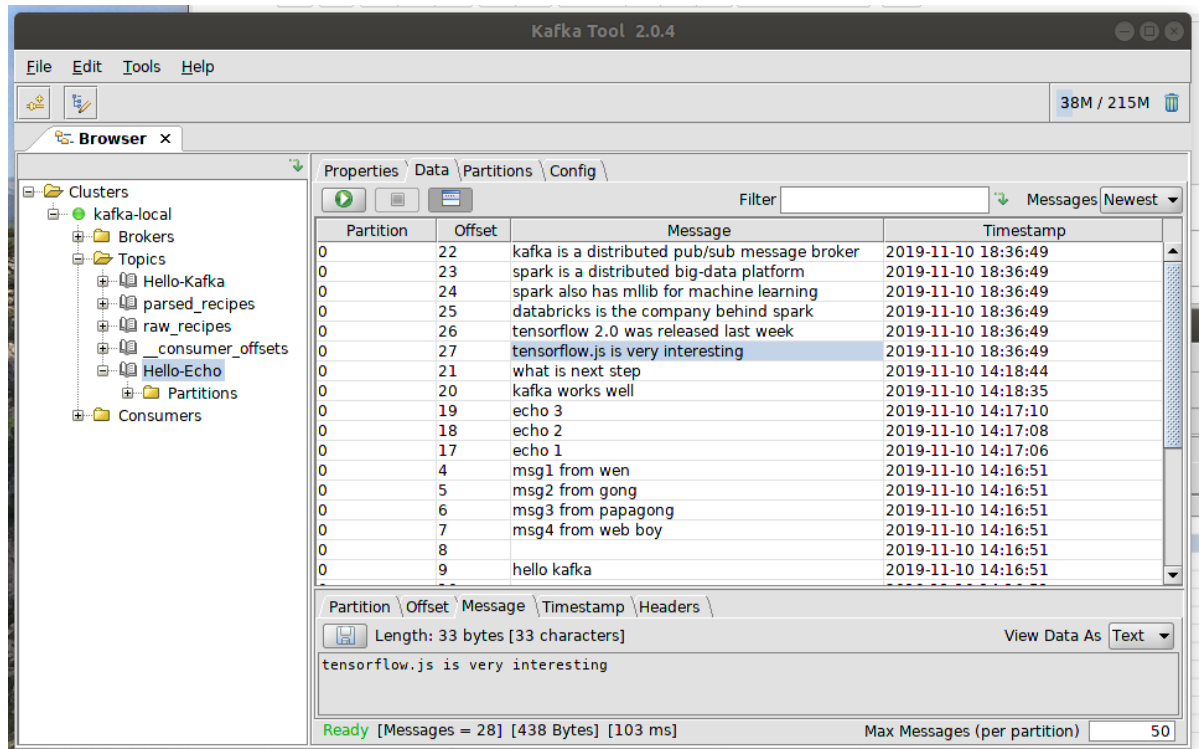
**write to another topic**

```
In [11]:  1  streaming.selectExpr("CAST(key AS STRING)", "CAST(value AS STRING)")
          2    .writeStream\
          3    .format("kafka")\
          4    .option("kafka.bootstrap.servers", "localhost:9092")\
          5    .option("checkpointLocation", "/tmp/kafka-checkpoint")\
          6    .option("topic", "Hello-Echo")\
          7    .start()
```

Out[11]: `<pyspark.sql.streaming.StreamingQuery at 0x7fe08da2cf60>`

Check in Kafkatool to see messages are echoed to the new topic = "Hello-Echo"



```
In [ ]:  1
```

```
In [ ]:  1
```

below codes are not tested

```
In [ ]:  1  # Subscribe to 1 topic
         2  df1 = spark.readStream.format("kafka")\
         3    .option("kafka.bootstrap.servers", "host1:port1,host2:port2")\
         4    .option("subscribe", "topic1")\
         5    .load()
```

```
In [ ]:  1  # Subscribe to multiple topics
         2  df2 = spark.readStream.format("kafka")\
         3    .option("kafka.bootstrap.servers", "host1:port1,host2:port2")\
         4    .option("subscribe", "topic1,topic2")\
         5    .load()
```

```
In [ ]:    1  # Subscribe to a pattern
           2  df3 = spark.readStream.format("kafka")\
           3    .option("kafka.bootstrap.servers", "host1:port1,host2:port2")\
           4    .option("subscribePattern", "topic.*")\
           5    .load()
```

```
In [ ]:    1  # COMMAND ----------
           2
           3  df1.selectExpr("topic", "CAST(key AS STRING)", "CAST(value AS STRING
           4    .writeStream\
           5    .format("kafka")\
           6    .option("kafka.bootstrap.servers", "host1:port1,host2:port2")\
           7    .option("checkpointLocation", "/to/HDFS-compatible/dir")\
           8    .start()
           9
          10  df1.selectExpr("CAST(key AS STRING)", "CAST(value AS STRING)")\
          11    .writeStream\
          12    .format("kafka")\
          13    .option("kafka.bootstrap.servers", "host1:port1,host2:port2")\
          14    .option("checkpointLocation", "/to/HDFS-compatible/dir")\
          15    .option("topic", "topic1")\
          16    .start()
          17
          18
          19  # COMMAND ----------
          20
          21  socketDF = spark.readStream.format("socket")\
          22    .option("host", "localhost").option("port", 9999).load()
```

```
In [ ]:    1  # COMMAND ----------
           2
           3  activityCounts.writeStream.trigger(processingTime='5 seconds')\
           4    .format("console").outputMode("complete").start()
```

```
In [ ]:    1  # COMMAND ----------
           2
           3  activityCounts.writeStream.trigger(once=True)\
           4    .format("console").outputMode("complete").start()
```

```
In [ ]:    1  # COMMAND ----------
```