# Projekat 2 – Big data

Nikola Đorđević 1567

# Kafka producer

```python
producer = KafkaProducer(bootstrap_servers='kafka:9092',
                         value_serializer=lambda x: x.encode('utf-8'))

with open('brightkite_small.csv') as file:
    reader = csv.reader(file)
    for row in reader:
        message = ','.join(row)
        timestamp = str(datetime.now())
        message_with_timestamp = f"{message},{timestamp}"
        producer.send('topic3', value=message_with_timestamp)

producer.flush()
```

- Kafka producer čita red po red iz fajla u kome se nalaze podaci, i pročitani red šalje na odgovarajući topic
- Polja jednog reda se šalju odvojena zarezom
- Za svaki poslati red se dodaje i vremenska oznaka kada je poslat

Spark streaming aplikacija

# Čitanje i parsiranje podataka sa Kafka topic-a

```python
df = spark \
  .readStream \
  .format("kafka") \
  .option("kafka.bootstrap.servers", "kafka:9092") \
  .option("subscribe", "topic3") \
  .option("startingOffsets", "earliest") \
  .option("maxOffsetsPerTrigger", 1) \
  .load()

parsed_df =  df.selectExpr("CAST(value AS STRING)")

split_col = split(parsed_df['value'], ',')
parsed_df = parsed_df.withColumn("user", split_col.getItem(0))
parsed_df = parsed_df.withColumn("check_in_time", split_col.getItem(1))
parsed_df = parsed_df.withColumn("latitude", split_col.getItem(2))
parsed_df = parsed_df.withColumn("longitude", split_col.getItem(3))
parsed_df = parsed_df.withColumn("location_id", split_col.getItem(4))
parsed_df = parsed_df.withColumn("time_spent", split_col.getItem(5))
parsed_df = parsed_df.withColumn("time_stamp", split_col.getItem(6))

parsed_df = parsed_df.where(parsed_df.user != 'user')
parsed_df = parsed_df.withColumn("check_in_time", to_timestamp(parsed_df["check_in_time"], "yyyy-MM-dd'T'HH:mm:ss.SSSXXX"))
parsed_df = parsed_df.withColumn("time_stamp", from_utc_timestamp(parsed_df["time_stamp"], "UTC"))
```

# Statistiški proračuni za provedeno vreme na određenoj lokaciji, za određenog korisnika

```python
# Prosecno, maksimalno i minimalno vreme koje je svaki korisnik proveo na odredjenoj lokaciji
windowed_df_1 = parsed_df.groupBy(window(parsed_df.time_stamp, "10 seconds"), parsed_df.user, parsed_df.location_id) \
    .agg(avg(parsed_df.time_spent).alias("avg_time_spent"),
        min(parsed_df.time_spent).alias("min_time_spent"),
        max(parsed_df.time_spent).alias("max_time_spent"))
```

```python
def write_to_influxdb(df, epoch_id):
    # Inicijalizacija InfluxDB klijenta
    token = "Xl99hGe7tyPW-wWrgpZRo8vOxA6bK2nR-X3MEoqkigZqnSG1vSqpKOoBmLZdWpWbYKKMKNEfqAAX4FMoKhd5ug=="
    org = "brightkite-org"
    bucket = "brightkite-bucket"
    client = InfluxDBClient(url="http://influxdb:8086", token=token)

    # Kreiranje instance WriteApi klase
    write_api = client.write_api(write_options=SYNCHRONOUS)

    df = df\
        .withColumn("avg_time_spent", col("avg_time_spent").cast("double")) \
        .withColumn("min_time_spent", col("min_time_spent").cast("double")) \
        .withColumn("max_time_spent", col("max_time_spent").cast("double"))

    for row in df.collect():
        point = Point("statistics_3") \
            .tag("window", row.window) \
            .tag("location_id", row.location_id) \
            .tag("user", row.user) \
            .field("avg_time_spent", row.avg_time_spent) \
            .field("min_time_spent", row.min_time_spent) \
            .field("max_time_spent", row.max_time_spent)
        write_api.write(bucket=bucket, org=org, record=point)
```
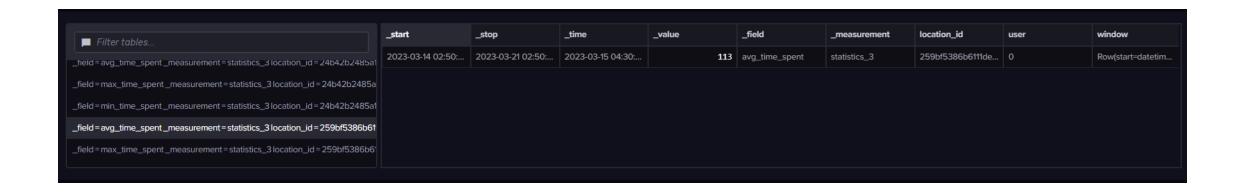
- Za vremenski prozor od 10 sekundi, podaci se grupišu po korisnicima i lokacijama
- Računa se prosečno, minimalno i maksimalno vreme koje je korisnik proveo na lokaciji

# Rezultat u InfluxDB bazi podataka

# Broj korisnika koji je posetio svaku od lokacija

```python
# Broj korisnika koji je posetio svaku od lokacija
windowed_df_2 = parsed_df.groupBy(window(parsed_df.time_stamp, "10 seconds"), parsed_df.location_id) \
    .agg(approx_count_distinct(parsed_df.user).alias("distinct_users"))
```

```python
def write_to_influxdb_2(df, epoch_id):
    # Inicijalizacija InfluxDB klijenta
    token = "Xl99hGe7tyPW-wWrgpZRo8vOxA6bK2nR-X3MEoqkigZqnSG1vSqpKOoBmLZdWpWbYKKMKNEfqAAX4FMoKhd5ug=="
    org = "brightkite-org"
    bucket = "brightkite-bucket"
    client = InfluxDBClient(url="http://influxdb:8086", token=token)

    # Kreiranje instance WriteApi klase
    write_api = client.write_api(write_options=SYNCHRONOUS)

    for row in df.collect():
        point = Point("dist_users") \
            .tag("window", row.window) \
            .tag("location_id", row.location_id) \
            .field("distinct_users", row.distinct_users)
        write_api.write(bucket=bucket, org=org, record=point)
```

- Određivanje broja korisnika koji je posetio svaku od lokacija za vremenski prozor od 10 sekundi, i upis rezultata u InfluxDB bazu podataka
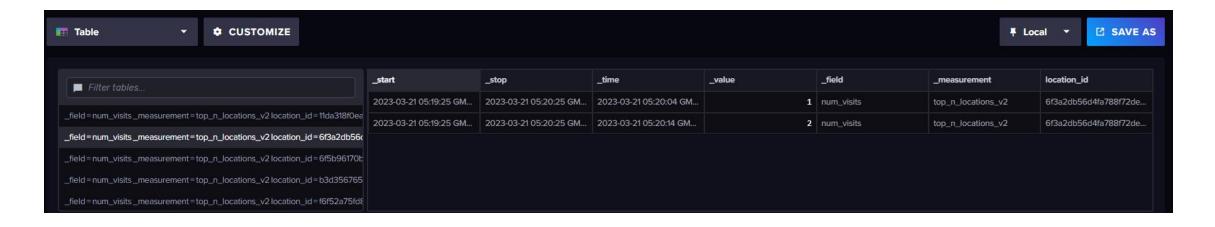
# Rezultat u InfluxDB bazi podataka

# Određivanje top N lokacija i upis u InfluxDB

```python
windowed_df_3 = parsed_df.groupBy(window(parsed_df.time_stamp, "10 seconds"), parsed_df.location_id) \
    .agg(count("*").alias("num_visits"))
```

```python
def write_to_influxdb_3(df, epoch_id):
    N = 2
    # Inicijalizacija InfluxDB klijenta
    token = "Xl99hGe7tyPW-wWrgpZRo8vOxA6bK2nR-X3MEoqkigZqnSG1vSqpKOoBmLZdWpWbYKKMKNEfqAAX4FMoKhd5ug=="
    org = "brightkite-org"
    bucket = "brightkite-bucket"
    client = InfluxDBClient(url="http://influxdb:8086", token=token)

    # Kreiranje instance WriteApi klase
    write_api = client.write_api(write_options=SYNCHRONOUS)

    df = df.orderBy(col("num_visits").desc()).limit(N)

    for row in df.collect():
        point = Point("top_n_locations_v2") \
            .tag("location_id", row.location_id) \
            .field("num_visits", row.num_visits)
        write_api.write(bucket=bucket, org=org, record=point)
```

Na osnovu broja poseta za određenu lokaciju, u vremenskom prozoru od 10 sekundi, vrši se sortiranje lokacija i određivanje top N lokacija, koje se upisuju u InfluxDB bazu

# Rezultat u InfluxDB bazi podataka

| | _start | _stop | _time | _value | _field | _measurement | location_id |
|---|---|---|---|---|---|---|---|
| Filter tables... | | | | | | | |
| _field=num_visits_measurement=top_n_locations_v2 location_id=11da318f0ea | 2023-03-21 05:19:25 GM... | 2023-03-21 05:20:25 GM... | 2023-03-21 05:20:04 GM... | 1 | num_visits | top_n_locations_v2 | 6f3a2db56d4fa788f72de... |
| _field=num_visits_measurement=top_n_locations_v2 location_id=6f3a2db56c | 2023-03-21 05:19:25 GM... | 2023-03-21 05:20:25 GM... | 2023-03-21 05:20:14 GM... | 2 | num_visits | top_n_locations_v2 | 6f3a2db56d4fa788f72de... |
| _field=num_visits_measurement=top_n_locations_v2 location_id=6f5b96170k | | | | | | | |
| _field=num_visits_measurement=top_n_locations_v2 location_id=b3d356765 | | | | | | | |
| _field=num_visits_measurement=top_n_locations_v2 location_id=f6f52a75fd8 | | | | | | | |

# Skripta za pokretanje aplikacije

```bash
#!/bin/bash

spark/bin/spark-submit --master spark://spark-master:7077 --packages org.apache.spark:spark-sql-kafka-0-10_2.12:3.1.2 app_2.py
```

# Web UI

**Application: PySpark_Kafka_Consumer**

*Spark 3.1.2*

**ID:** app-20230321150800-0000
**Name:** PySpark_Kafka_Consumer
**User:** root
**Cores:** Unlimited (16 granted)
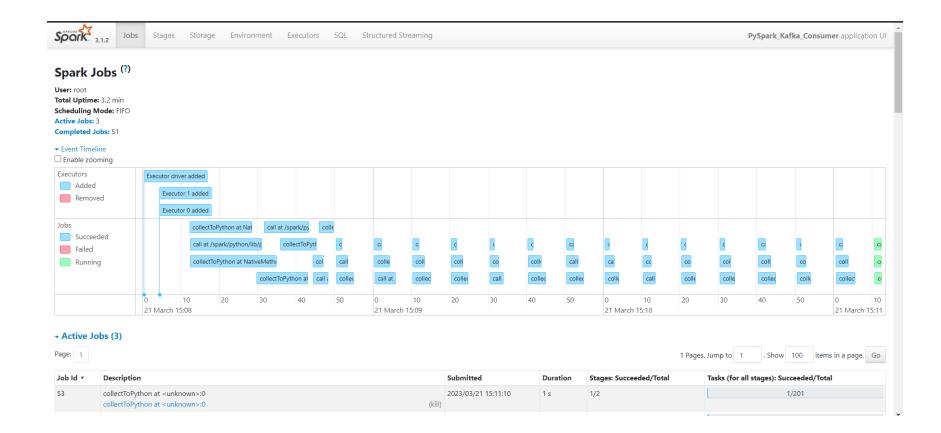**Executor Limit:** Unlimited (2 granted)
**Executor Memory:** 1024.0 MiB
**Executor Resources:**
**Submit Date:** 2023/03/21 15:08:00
**State:** RUNNING
**Application Detail UI**

### ▾ Executor Summary (2)

| ExecutorID | Worker | Cores | Memory | Resources | State | Logs |
|---|---|---|---|---|---|---|
| 1 | worker-20230321150626-172.18.0.7-37795 | 8 | 1024 | | RUNNING | stdout stderr |
| 0 | worker-20230321150627-172.18.0.8-46755 | 8 | 1024 | | RUNNING | stdout stderr |

# Web UI

# Flink aplikacija

# Čitanje i parsiranje podataka sa Kafka topic-a

```java
public static DataStream<CheckIn> StreamConsumer(String inputTopic, String server, StreamExecutionEnvironment environment) throws Exception {
    FlinkKafkaConsumer<String> flinkKafkaConsumer = createStringConsumerForTopic(inputTopic, server);
    DataStream<String> stringInputStream = environment.addSource(flinkKafkaConsumer);

    return stringInputStream.map(new MapFunction<String, CheckIn>() {
        private static final long serialVersionUID = -999736771747691234L;

        @Override
        public CheckIn map(String value) throws Exception {
            String[] split = value.split(regex:",");
            return new CheckIn(
                    Integer.parseInt(split[0]),
                    split[1],
                    Double.parseDouble(split[2]),
                    Double.parseDouble(split[3]),
                    split[4],
                    Integer.parseInt(split[5])
            );
        }
    });
}
```

# Statistički proračuni za provedeno vreme na određenoj lokaciji

```java
SingleOutputStreamOperator groupedStream = dataStream
.keyBy(CheckIn::getLocation_id)
.window(SlidingProcessingTimeWindows.of(Time.seconds(10), Time.seconds(10)))
.aggregate(new MinMaxAggregateFunction());


groupedStream.print();


CassandraSink.addSink(groupedStream)
        .setHost("cassandra")
        .setQuery("INSERT INTO flink.min_max(location_id, max, min) values (?, ?, ?);")
        .build();
```

- Određivanje minimalnog i maksimalnog provedenog vremena na svakoj od lokacija i upis rezultata u Cassandra bazu podataka
- Podaci su grupisani u vremenski (po tipu "Sliding") prozor veličine 10 sekundi i koraka 10 sekundi

# Rezultat u Cassandra DB



| location_id | max | min |
|---|---|---|
| a8d770948720fd154ccfe093a21b73e0 | 201 | 201 |
| 4dfb8a7526182e70b53ea6a7ee670b5a3d694375 | 277 | 277 |
| 424eb3dd143292f9e013efa00486c907 | 354 | 340 |
| e653dbdc794b11ddbd6a0030487eb504 | 347 | 347 |
| dd7cd3d264c2d063832db506fba8bf79 | 217 | 107 |
| b3d356765cc8a4aa7ac5cd18caafd393 | 200 | 3 |
| 3a955c5c8336ae8040fa61ceece8a0c0 | 30 | 30 |
| 6f5b96170b7744af3c7577fa35ed0b8f | 188 | 188 |
| 08a35293e09f508494096c1c1b3819edb9df50db | 267 | 8 |
| dd98ca10f19211dd948d003048c10834 | 81 | 81 |
| 6ccba23d2e036ef1bcff1feb17443d6c7ef8f579 | 210 | 14 |
| 6f3a2db56d4fa788f72def616f79b7a4 | 254 | 20 |
| d2f35ed2724211deafe6003048c0801e | 184 | 184 |
| 36740e782a458ebcab9f14ae3a8a5f19d04d0e41 | 151 | 151 |
| 828fb63e770c11ddb8450030487eb504 | 261 | 261 |
| 85cd3b30938d11dd9bf0003048c10834 | 331 | 331 |
| e2460b3f5040fc666cd99128a0c133b7 | 333 | 333 |
| 098e8e9f4b606ce204a76dc504b61697 | 284 | 284 |
| b3fe66207ee111dd89eb0030487eb504 | 352 | 352 |

# Statistički proračuni za provedeno vreme na određenoj lokaciji

```
SingleOutputStreamOperator meanStream = dataStream
    .keyBy(CheckIn::getLocation_id)
    .window(SlidingProcessingTimeWindows.of(Time.seconds(10), Time.seconds(10)))
    .aggregate(new MeanAggregate());


meanStream.print();


CassandraSink.addSink(meanStream)
        .setHost("cassandra")
        .setQuery("INSERT INTO flink.mean(location_id, mean) values (?, ?);")
        .build();
```

- Određivanje prosečnog provedenog vremena na svakoj od lokacija i upis rezultata u Cassandra bazu podataka
- Podaci su grupisani u vremenski (po tipu "Sliding") prozor veličine 10 sekundi i koraka 10 sekundi

# Rezultat u Cassandra DB

# Broj korisnika koji je posetio svaku od lokacija

```java
DataStream<Tuple2<String, Integer>> location_user_counts = dataStream
.map(new MapFunction<CheckIn, String>() {
    @Override
    public String map(CheckIn value) throws Exception {
        return value.getLocation_id() + "-" + value.getUser();
    }
})
.map(new MapFunction<String, Tuple2<String, Integer>>() {
    @Override
    public Tuple2<String, Integer> map(String value) throws Exception {
        return new Tuple2<String, Integer>(value, 1);
    }
})
.keyBy(new KeySelector<Tuple2<String, Integer>, String>() {
    @Override
    public String getKey(Tuple2<String, Integer> value) throws Exception {
        return value.f0;
    }
})
//.window(SlidingProcessingTimeWindows.of(Time.seconds(10), Time.seconds(10)))
.window(TumblingProcessingTimeWindows.of(Time.seconds(10)))
.sum(1);

location_user_counts.print();

CassandraSink.addSink(location_user_counts)
        .setHost("cassandra")
        .setQuery("INSERT INTO flink.location_user_counts(user_location, counts) values (?, ?);")
        .build();
```

- Određivanje broja korisnika koji je posetio svaku od lokacija i upis rezultata u Cassandra bazu podataka
- Podaci su grupisani u vremenski (po tipu "Tumbling") prozor veličine 10 sekundi

# Rezultat u Cassndra DB



| user_location | counts |
| --- | --- |
| 64b925364ac71005af756eabb94096e2-0 | 1 |
| 9848afcc62e500a01cf6fbf24b797732f8963683-0 | 2 |
| a2451a4a10d0061e87ebcd21b422f44b-0 | 1 |
| 7a0f88982aa015062b95e3b4843f9ca2-0 | 5 |
| e2460b3f5040fc666cd99128a0c133b7-0 | 1 |
| 974b77e179e872596500b2d22c38bf26-0 | 1 |
| c83ac485e066dbf68b8619c7ab5b2579-0 | 1 |
| 8f060f74f59a02df0993b24964334eea-0 | 2 |
| d2f35ed2724211deafe6003048c0801e-0 | 1 |
| 828fb63e770c11ddb8450030487eb504-0 | 1 |
| 115b5206a20011dd9f01003048c10834-0 | 1 |
| dcc06bf19e775f436c2225be50e14922-0 | 2 |
| c69bd906c5210caee0719913f836303c-0 | 1 |
| e63d729e56954aeb23ba669d2c7a2805-0 | 1 |
| 6f5b96170b7744af3c7577fa35ed0b8f-0 | 1 |

# Top N lokacija

```java
DataStream<Tuple2<String, Integer>> topN = dataStream
.map(new MapFunction<CheckIn, String>() {
    @Override
    public String map(CheckIn value) throws Exception {
        return value.getLocation_id();
    }
})
.map(new MapFunction<String, Tuple2<String, Integer>>() {
    @Override
    public Tuple2<String, Integer> map(String value) throws Exception {
        return new Tuple2<String, Integer>(value, 1);
    }
})
.keyBy(new KeySelector<Tuple2<String, Integer>, String>() {
    @Override
    public String getKey(Tuple2<String, Integer> value) throws Exception {
        return value.f0;
    }
})
//.window(SlidingProcessingTimeWindows.of(Time.seconds(10), Time.seconds(10)))
.window(TumblingProcessingTimeWindows.of(Time.seconds(10)))
.sum(1)
//.windowAll(SlidingProcessingTimeWindows.of(Time.seconds(10), Time.seconds(10)))
.windowAll(TumblingProcessingTimeWindows.of(Time.seconds(10)))
.process(new ProcessAllWindowFunction<Tuple2<String, Integer>, Tuple2<String, Integer>, TimeWindow>() {
    @Override
    public void process(Context context, Iterable<Tuple2<String, Integer>> iterable, Collector<Tuple2<String, Integer>> collector) throws Exception {
        PriorityQueue<Tuple2<String, Integer>> queue = new PriorityQueue<>(Comparator.comparingInt(o -> o.f1));
        for (Tuple2<String, Integer> t : iterable) {
            queue.offer(t);
            if (queue.size() > N) {
                queue.poll();
            }
        }
        List<Tuple2<String, Integer>> topN = new ArrayList<>(queue);
        topN.sort(Comparator.comparingInt(o -> -o.f1));
        for (Tuple2<String, Integer> t : topN) {
            collector.collect(t);
        }
    }
});
```
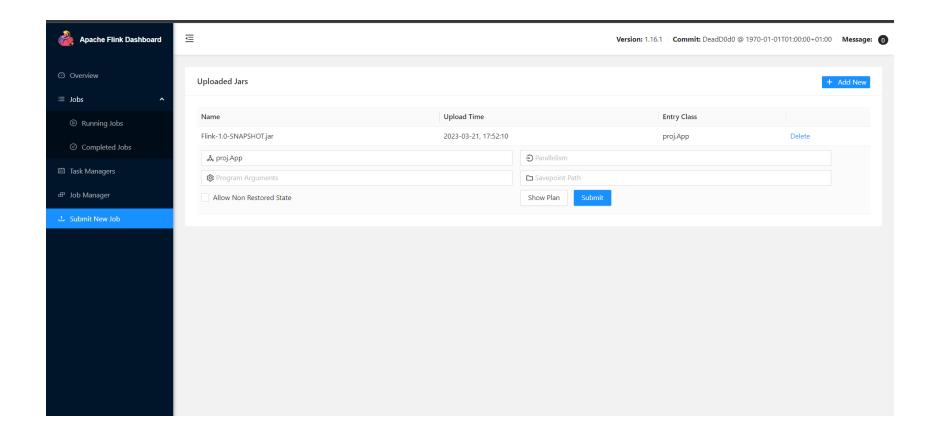
- Pronalazak Top N lokacija u odnosu na broj poseta
- Podaci su grupisani u vremenski (po tipu "Tumbling") prozor veličine 10 sekundi

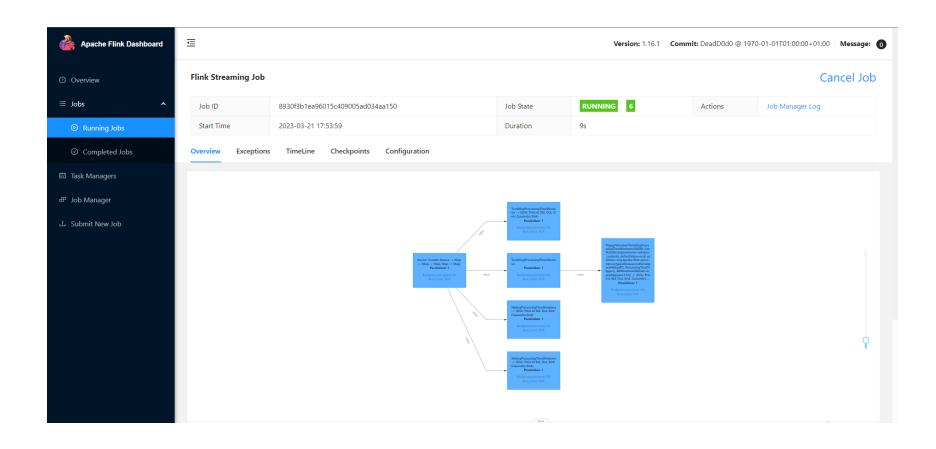# Rezultat u Cassndra DB

```
location_id                      | counts
---------------------------------+--------
dd7cd3d264c2d063832db506fba8bf79 |      6
6f3a2db56d4fa788f72def616f79b7a4 |      5
2ef143e12038c870038df53e0478cefc |      9
f6f52a75fd80e27e3770cd3a87054f27 |      7
7a0f88982aa015062b95e3b4843f9ca2 |      5
```

# Top N lokacija  - Cassandra DB

```
CassandraSink.addSink(topN)
        .setHost("cassandra")
        .setQuery("INSERT INTO flink.top_n_locations(location_id, counts) values (?, ?);")
        .build();
```

# Web UI



- Aplikacija se startuje tako što se upload-uje željeni JAR fajl preko web UI - a

# Web UI

# Web UI

| Name | Status | Bytes Received | Records Received | Bytes Sent | Records Sent | Parallelism | Start Time | Tasks |
|------|--------|----------------|------------------|------------|--------------|-------------|------------|-------|
| Source: Custom Source -> Map -> (Map -> Map, Map -> Map) | RUNNING | 0 B | 0 | 0 B | 198 | 1 | 2023-03-21 17: | 1 |
| TumblingProcessingTimeWindows -> (Sink: Print to Std. Out, Sink: Ca... | RUNNING | 5.15 KB | 99 | 0 B | 0 | 1 | 2023-03-21 17: | 1 |
| TumblingProcessingTimeWindows | RUNNING | 4.96 KB | 99 | 0 B | 0 | 1 | 2023-03-21 17: | 1 |
| TriggerWindow(TumblingProcessingTimeWindows(10000), ListStateD... | RUNNING | 4 B | 0 | 0 B | 0 | 1 | 2023-03-21 17: | 1 |
| SlidingProcessingTimeWindows -> (Sink: Print to Std. Out, Sink: Cassa... | RUNNING | 10.5 KB | 99 | 0 B | 0 | 1 | 2023-03-21 17: | 1 |
| SlidingProcessingTimeWindows -> (Sink: Print to Std. Out, Sink: Cassa... | RUNNING | 10.5 KB | 99 | 0 B | 0 | 1 | 2023-03-21 17: | 1 |