# PROGRAMAÇÃO DE ENTRADAS STREAMING EM CLUSTERS

Segundo projeto de pesquisa

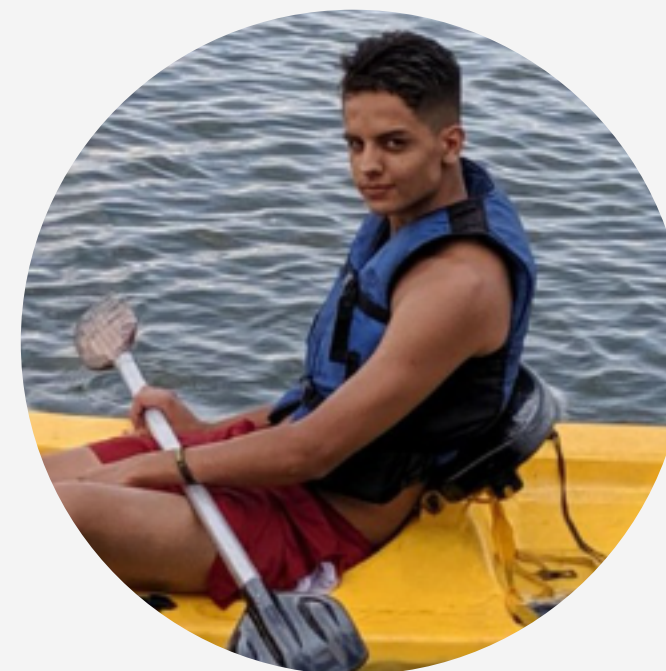# Equipe

**Eduardo Afonso**
19/0012307

**Thiago Paiva**
19/0020377

**Rafael Cleydson**
19/0019085

# Problema

Contar palavras
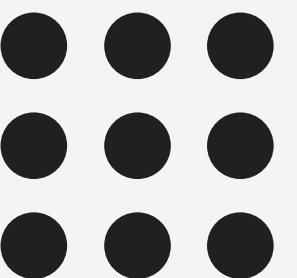Contar palavras que começam com S
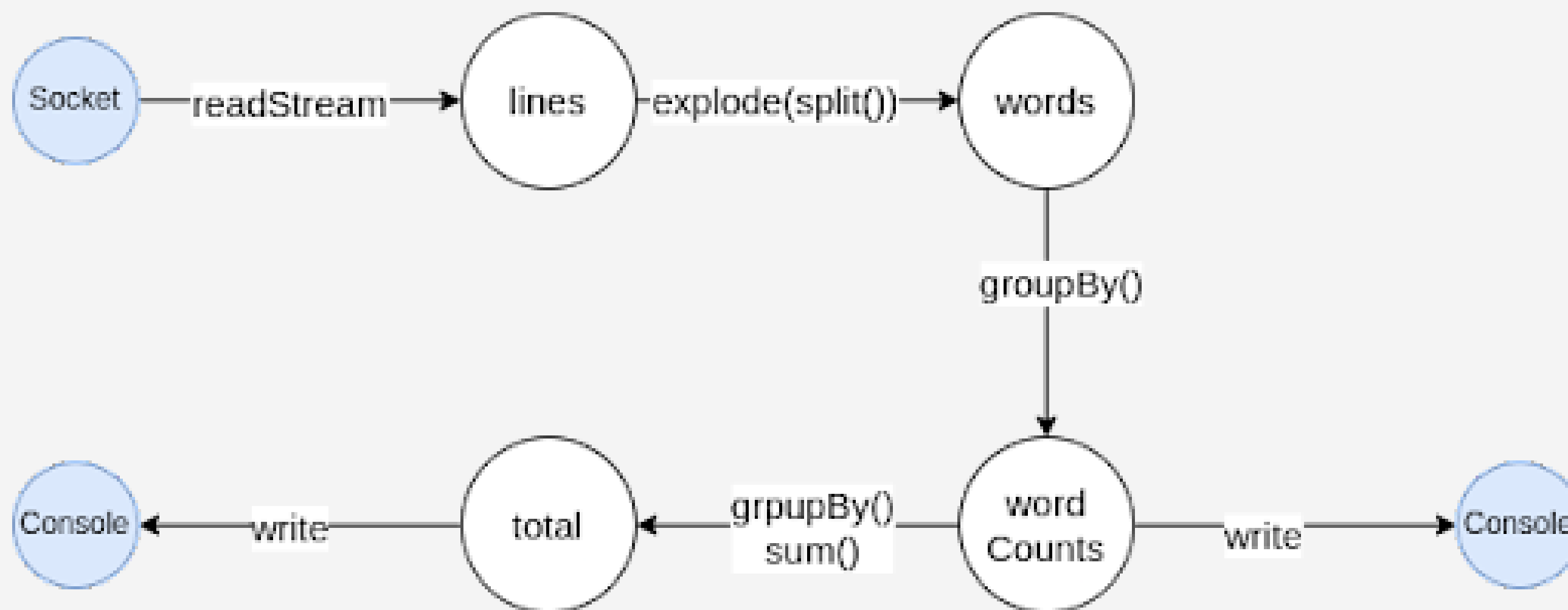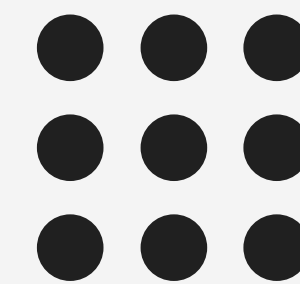Contar palavras que começam com P
Contar palavras que começam com R
Contar palavras com 6 caracteres
Contar palavras com 8 caracteres
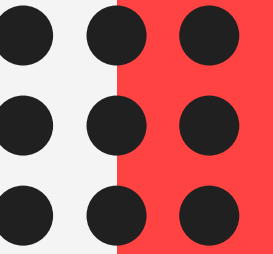Contar palavras com 11 caracteres

# Solução Socket - Parte 1

Socket —readStream→ lines —explode(split())→ words

words —groupBy()→ word Counts

word Counts —grpupBy() sum()→ total —write→ Console

word Counts —write→ Console

# Solução Socket - Conexão

```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import split, explode, lit, col, upper

SOCKET_HOST = "localhost"
SOCKET_PORT = "9999"

spark = SparkSession \
    .builder \
    .appName("P2 - PSPD - Socket") \
    .getOrCreate()

lines = spark \
    .readStream \
    .format("socket") \
    .option("host", SOCKET_HOST) \
    .option("port", SOCKET_PORT) \
    .load()
```
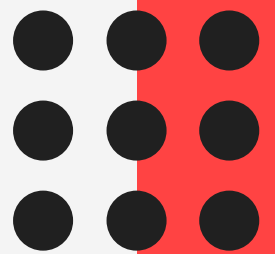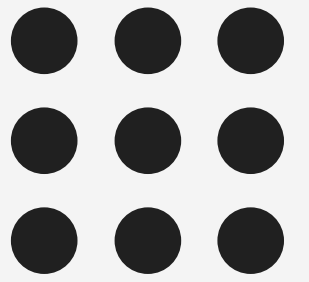
# Solução Socket - Contagem

```python
19  # Split the lines into words
20  words = lines.select(
21      explode(
22          split(lines.value, "\s+")
23      ).alias('word')
24  )
25  words = words.select(upper(words.word).alias('word'))
26
27  # Generate running word count
28  wordCounts = words.groupBy("word").count()
29
30  def foreach_batch_func(df, _):
31      """Find the total number of words and write it and wordsCount in console"""
32      total = df \
33          .groupBy() \
34          .sum() \
35          .select(lit('TOTAL').alias('key'), col('sum(count)').alias('value'))
36
37      df.write.format('console').save()
38      total.write.format('console').save()
39
40  # Sink
41  query = wordCounts \
42      .writeStream \
43      .outputMode("complete") \
44      .foreachBatch(foreach_batch_func) \
45      .start()
46
47  query.awaitTermination()
```
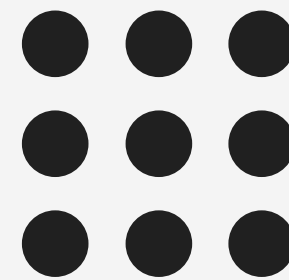
# Comandos Exec Socket

```
nc -lk 9999 < {filePath}
```
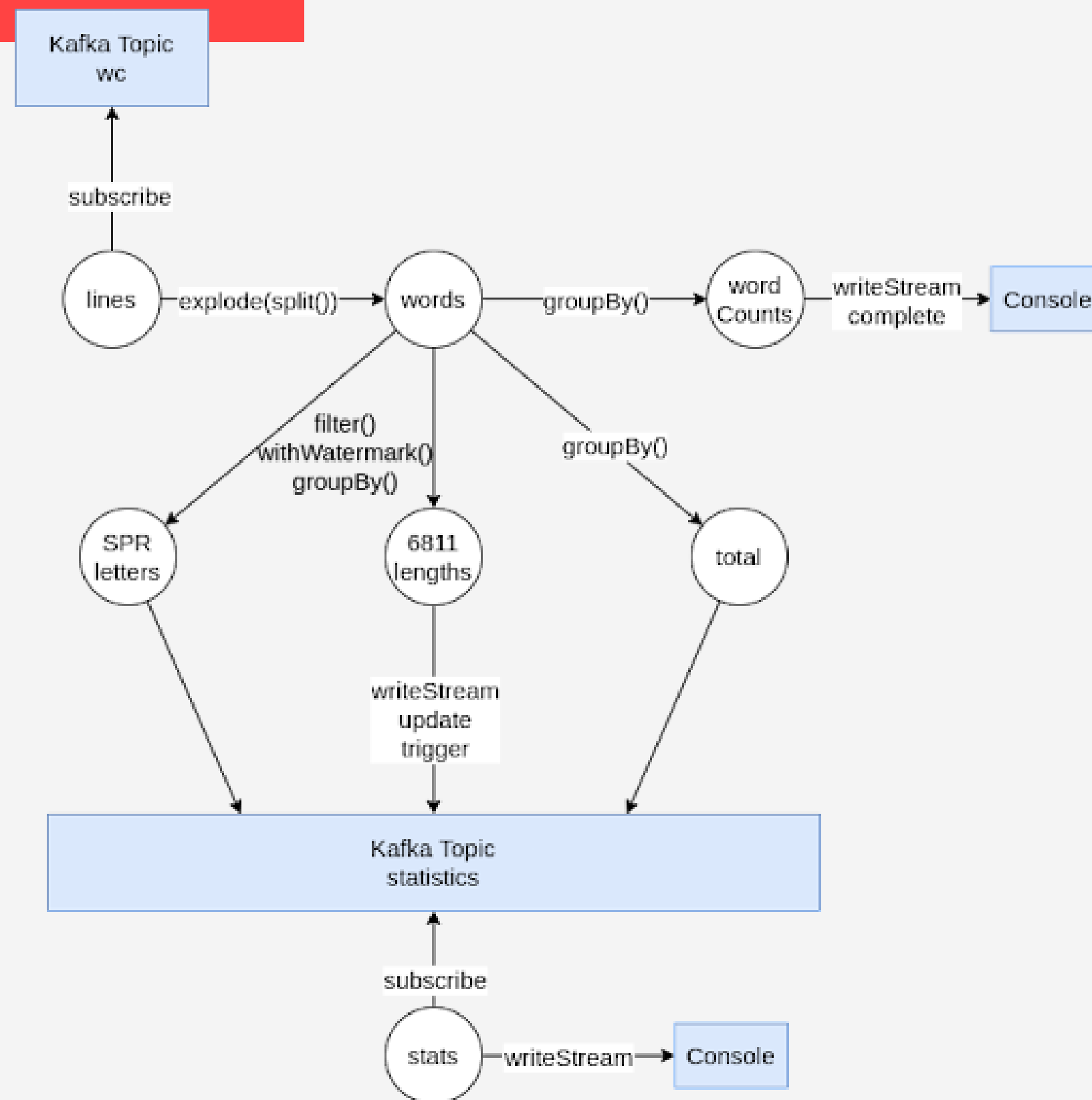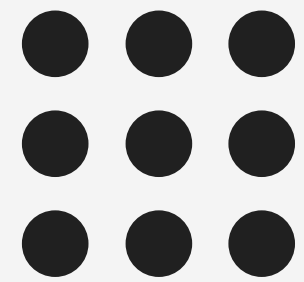
```
bin/pyspark
```

# Solução Socket Saida

```
+----------+-----+          +-----+--------+
|      word|count|          | key|   value|
+----------+-----+          +-----+--------+
|  YTSGWVEV|    1|          |TOTAL|19462001|
|        DZ|    1|          |    S|  170323|
|     VDXZW|    1|          |    R|  169696|
| WNZPUDOPOF|   1|          |    P|  170119|
|   FEKZUVF|    1|          +-----+--------+
|        LT|    1|
|EVVHRWKETXG|   1|
| DAAPTIXRX|    1|          ---------------------
| GFSZXJVEE|    1|          Batch: 8
| VWEXCUXOQW|   1|          ---------------------
|       LEB|    1|
|      NBKP|    1|          +---+------+
|         K|   17|          |key| value|
| ATPPYCWGX|    1|          +---+------+
|   PZKVVNI|    1|          |  8|368409|
|   FMXBZKE|    1|          | 11|367776|
|        MM|    1|          |  6|368207|
|KOTWSZIBIDT|   1|          +---+------+
|      RDJY|    1|
| BRKCICSACE|   1|          ---------------------
+----------+-----+          Batch: 9
only showing top 20 rows    ---------------------

                            +-----+---------+
                            | key|    value|
                            +-----+---------+
                            |TOTAL|100000001|
                            +-----+---------+
```

# Solução Kafka - Parte 2

## Solução Kafka - Conexao

```python
from pyspark.sql import SparkSession
from pyspark.sql.functions import length, explode, split, substring, upper, window

INTERVAL = '3 seconds'
KAFKA_SERVER = 'localhost:9092'
WORDS_TOPIC = 'wc'
STATS_TOPIC = 'statistics'

spark = SparkSession \
    .builder \
    .appName("P2 - PSPD - Transformer") \
    .getOrCreate()

# Create DataFrame representing the stream of input lines and subscribe it in kafka topic
lines = spark \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", KAFKA_SERVER) \
    .option("subscribe", WORDS_TOPIC) \
    .option('includeTimestamp', 'true') \
    .load()
```
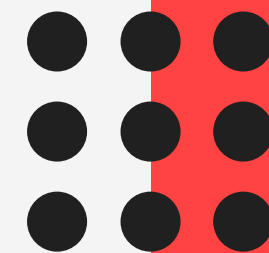
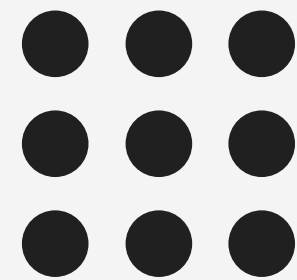# Solução Kafka - Contagem

```python
23  # Split the lines into words
24  words = lines.select(
25      explode(
26          split(lines.value, "\s+")).alias("word"),
27          lines.timestamp
28      )
29  words = words.select(upper(words.word).alias('word'), words.timestamp)
30
31  # Group words
32  wordCounts = words.groupBy("word").count()
33
34  # Count the total of words readed
35  total = words \
36      .groupBy() \
37      .count() \
38      .selectExpr("'TOTAL' as key", "CAST(count AS STR:
39
```

```python
40  # Count the words that startswith S, P and R
41  letters = words \
42      .filter(upper(substring(words.word, 0, 1)).isin(["S", "P", "R"])) \
43      .withWatermark("timestamp", INTERVAL) \
44      .groupBy(
45          window(words.timestamp, INTERVAL, INTERVAL),
46          upper(substring(words.word, 0, 1)).alias("key"),
47      ) \
48      .count() \
49      .selectExpr("key", "CAST(count AS STRING) as value")
50
51  # Count the words that has length 6, 8 and 11
52  lengths = words \
53      .filter(length(words.word).isin([6, 8, 11])) \
54      .withWatermark("timestamp", INTERVAL) \
55      .groupBy(
56          window(words.timestamp, INTERVAL, INTERVAL),
57          length(words.word).alias("key")
58      ) \
59      .count() \
60      .selectExpr("CAST(key AS STRING)", "CAST(count AS STRING) as value")
61
62  # Sinks
```
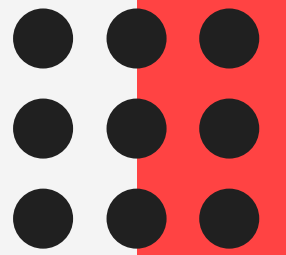
# Solução Kafka - Sink

```python
# Sinks
qW = wordCounts \
    .writeStream \
    .outputMode("complete") \
    .format("console") \
    .start()

qT = total \
    .writeStream \
    .outputMode("complete") \
    .format("kafka") \
    .option("kafka.bootstrap.servers", KAFKA_SERVER) \
    .option('topic', STATS_TOPIC) \
    .option('checkpointLocation', '/tmp/spark/total-stats') \
    .start()

qLen = lengths \
    .writeStream \
    .outputMode("update") \
    .format("kafka") \
    .option("kafka.bootstrap.servers", KAFKA_SERVER) \
    .option('topic', STATS_TOPIC) \
    .option('checkpointLocation', '/tmp/spark/len-stats') \
    .trigger(processingTime=INTERVAL) \
    .start()

qLet = letters \
    .writeStream \
    .outputMode("update") \
    .format("kafka") \
    .option("kafka.bootstrap.servers", KAFKA_SERVER) \
    .option('topic', STATS_TOPIC) \
    .option('checkpointLocation', '/tmp/spark/let-stats') \
    .trigger(processingTime=INTERVAL) \
    .start()
```
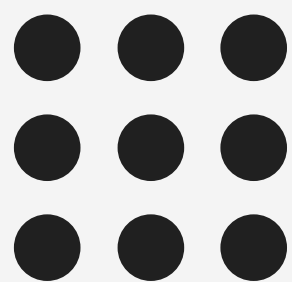
## Solução Kafka - Visualização das Estatisticas

```python
from pyspark.sql import SparkSession

KAFKA_SERVER = 'localhost:9092'
STATS_TOPIC = 'statistics'

spark = SparkSession \
    .builder \
    .appName("P2 - PSPD - Stats Consumer") \
    .getOrCreate()

stats = spark \
    .readStream \
    .format("kafka") \
    .option("kafka.bootstrap.servers", KAFKA_SERVER) \
    .option("subscribe", "statistics") \
    .load()

q = stats \
    .selectExpr("CAST(key AS STRING)", "CAST(value AS STRING)") \
    .writeStream \
    .format('console') \
    .outputMode('append') \
    .trigger(processingTime='3 seconds')\
    .start()

q.awaitTermination()
```

# Solução Kafka Saida

```
+----------+-----+
|      word|count|
+----------+-----+
|  YTSGWVEV|    1|
|        DZ|    1|
|     VDXZW|    1|
| WNZPUDOPOF|   1|
|   FEKZUVF|    1|
|        LT|    1|
|EVVHRWKETXG|   1|
| DAAPTIXRX|    1|
| GFSZXJVEE|    1|
| VWEXCUXOQW|   1|
|       LEB|    1|
|      NBKP|    1|
|         K|   17|
|  ATPPYCWGX|   1|
|   PZKVVNI|    1|
|   FMXBZKE|    1|
|        MM|    1|
|KOTWSZIBIDT|   1|
|      RDJY|    1|
| BRKCICSACE|   1|
+----------+-----+
only showing top 20 rows
```

```
+-----+--------+
|  key|   value|
+-----+--------+
|TOTAL|19462001|
|    S|  170323|
|    R|  169696|
|    P|  170119|
+-----+--------+


------------------------
Batch: 8
------------------------

+---+------+
|key| value|
+---+------+
|  8|368409|
| 11|367776|
|  6|368207|
+---+------+


------------------------
Batch: 9
------------------------

+-----+--------+
|  key|   value|
+-----+--------+
|TOTAL|100000001|
+-----+--------+
```

# Obrigado pela atenção