

Experiment report: Stream Processing and Flink

Done by student: Bylkova Kristina (伯汀娜), 1820249049

Course: Big Data Analysis Technology

Date: October 2024

Summary

This report shows the MindMap of stream processing and flink. Here is also the experiment screenshots.

Stream Processing MindMap



Figure 1: Stream Processing Mindmap

Flink MindMap

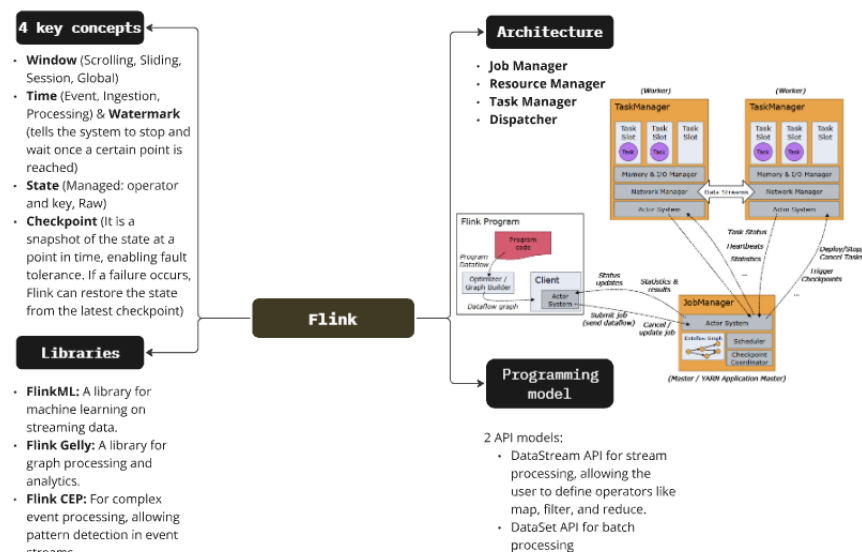


Figure 2: Flink MindMap

Task

1. Explain and illustrate the concepts of State and Checkpoint, as well as their working principles, and how the fault recovery can be achieved through State and Checkpoint. (using text and the diagram)
2. With reference to the Flink architecture diagram, explain the components and working principles of two architectures in Flink's YARN mode: the task submission architecture and the runtime architecture. (Flink runtime architectures include Local mode, Standalone cluster mode, and YARN mode).

Answer:

1. State is updated as events are processed. Each operator maintains some internal state as it processes each incoming event. Checkpointing occurs at regular intervals (configured by the user), where the system records the state of each operator. This means that at each checkpoint, a snapshot of the current state is written to a durable store.

If a failure occurs, the system will recover by using the most recent checkpoint (rolling back the system to the last successful checkpoint).

The fault recovery process works:

- Failure Detection (when a failure occurs, the system detects the failure)

- State Restoration (the system will attempt to recover by restoring the state from the last successful checkpoint)
- Reprocessing (after restoring the state, the system can resume processing from where it left off, reprocessing the unprocessed records)

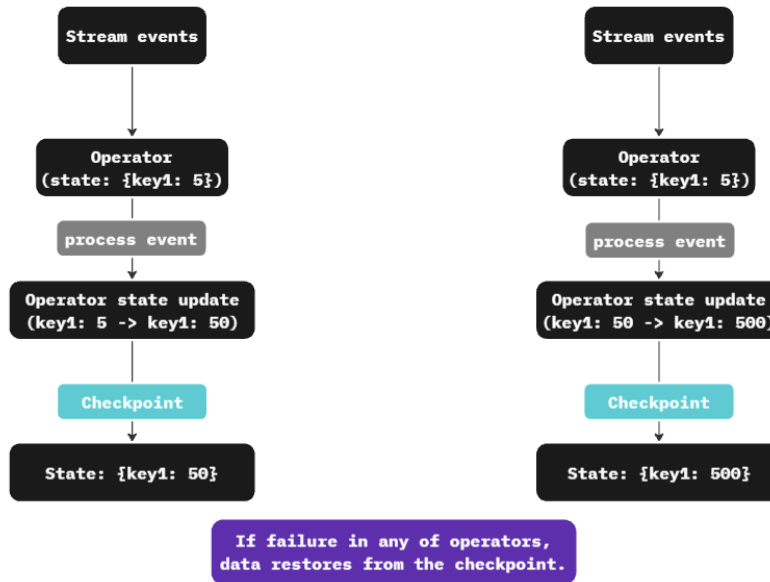


Figure 3: State and Checkpoint

2. Task Submission Architecture:

Components:

- Client (submit the Flink job to the YARN cluster)
- ResourceManager (manages resources across all available nodes in the cluster and allocates them for various applications)
- ApplicationMaster (manage the life cycle of the Flink job, including resource allocation and monitoring)
- JobManager (schedule tasks, distribute them to the TaskManagers, and manage job execution)

The client submits a Flink job to YARN using the JobManager. The ResourceManager assigns an ApplicationMaster to the job. The ApplicationMaster requests resources (containers) for the TaskManagers from the ResourceManager. Once resources are allocated, the JobManager coordinates the execution of tasks on the available TaskManagers.

Flink runtime architecture:

Components:

- JobManager coordinates the job execution. It is responsible for scheduling tasks, managing state, and recovering from failures
- TaskManagers are the worker nodes responsible for executing the tasks assigned by the JobManager
- The ResourceManager is responsible for managing the resources in the YARN cluster

In Local Mode, Flink runs on a single machine for development and testing purposes. In Standalone Cluster Mode, Flink is deployed across multiple machines or containers, where it runs in a distributed mode, but without any external resource manager like YARN or Mesos. In YARN Mode, Flink runs as an application within a YARN (Yet Another Resource Negotiator) cluster. YARN manages the resources across the cluster, and Flink leverages it to allocate resources for job execution.

Running

```
kristina@kristina-VirtualBox:~/flink-1.19.1$ ./bin/start-cluster.sh
Starting cluster.
Starting standalone session daemon on host kristina-VirtualBox.
Starting taskexecutor daemon on host kristina-VirtualBox.
```

Figure 4: In terminal

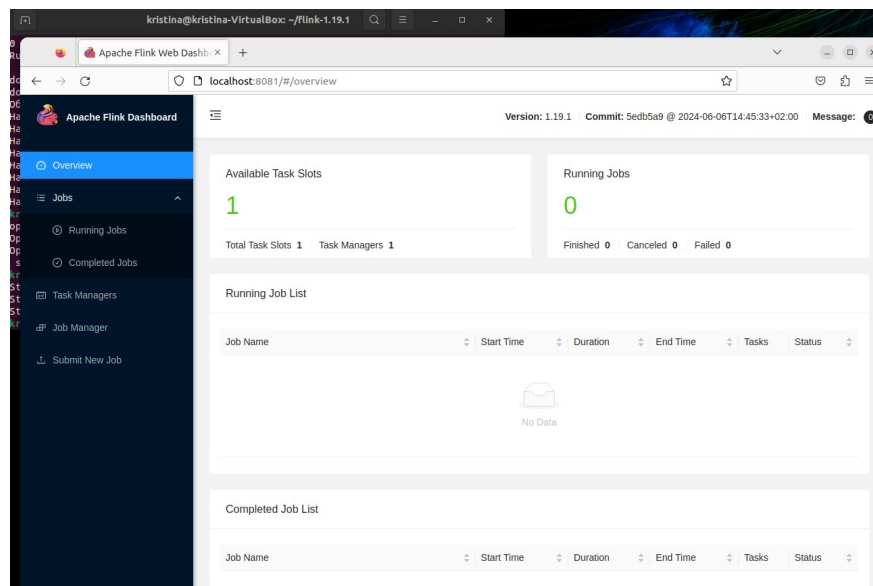


Figure 5: Web page

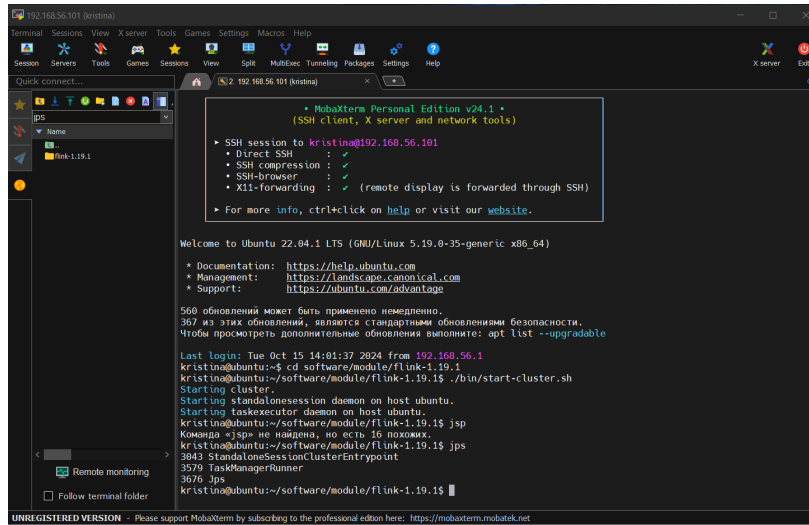


Figure 6: Connection via MobaXTerm

Flink operations

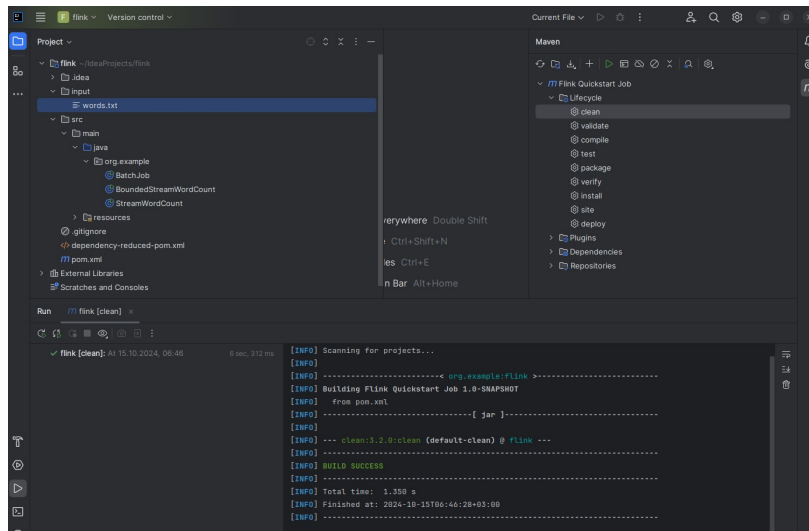


Figure 7: Clean

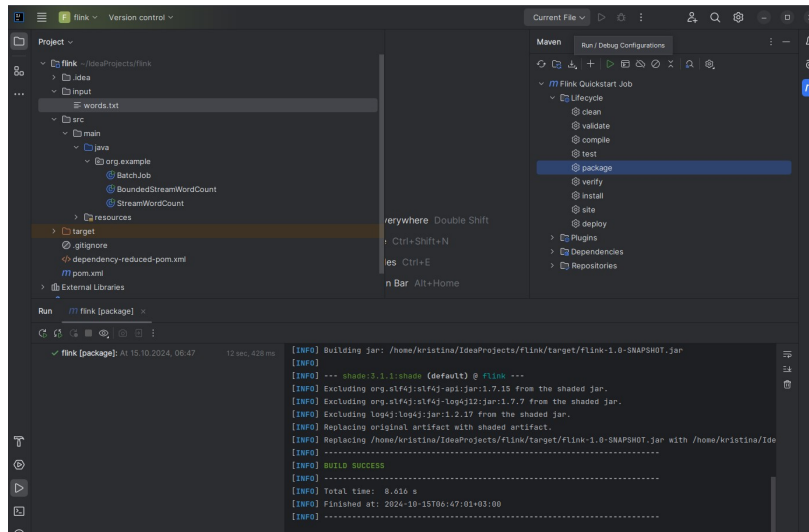


Figure 8: Package

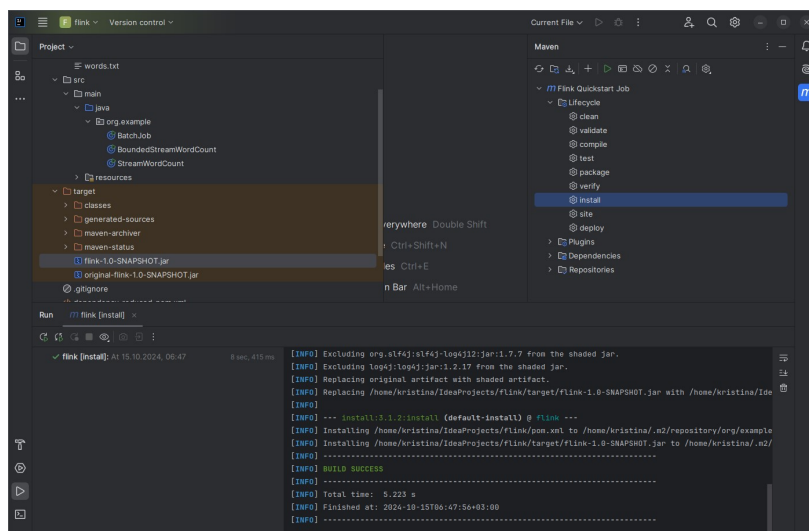


Figure 9: Install

Experiment

BatchJob experiment

```
1 package org.example;
2
3 import org.apache.flink.api.common.functions.FlatMapFunction;
4 import org.apache.flink.api.java.DataSet;
5 import org.apache.flink.api.java.ExecutionEnvironment;
6 import org.apache.flink.api.java.operators.DataSource;
7 import org.apache.flink.api.java.tuple.Tuple2;
8 import org.apache.flink.util.Collector;
9
10
11 public class BatchJob {
12
13     public static void main(String[] args) throws Exception {
14         final ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();
15         DataSource<String> stringDataSource = env.readTextFile("input/words.txt");
16
17         DataSet<Tuple2<String, Integer>> counts =
18             stringDataSource.flatMap(new Tokenizer())
19                 .groupBy(0)
20                 .sum(1);
21         counts.print();
22     }
23
24     public static class Tokenizer implements FlatMapFunction<String, Tuple2<String, Integer>> {
25
26         @Override
27         public void flatMap(String value, Collector<Tuple2<String, Integer>> out) {
28             String[] tokens = value.toLowerCase().split("\\W+");
29
30             for (String token : tokens) {
31                 if (token.length() > 0) {
32                     out.collect(new Tuple2<>(token, 1));
33                 }
34             }
35         }
36     }
37 }
38 }
```

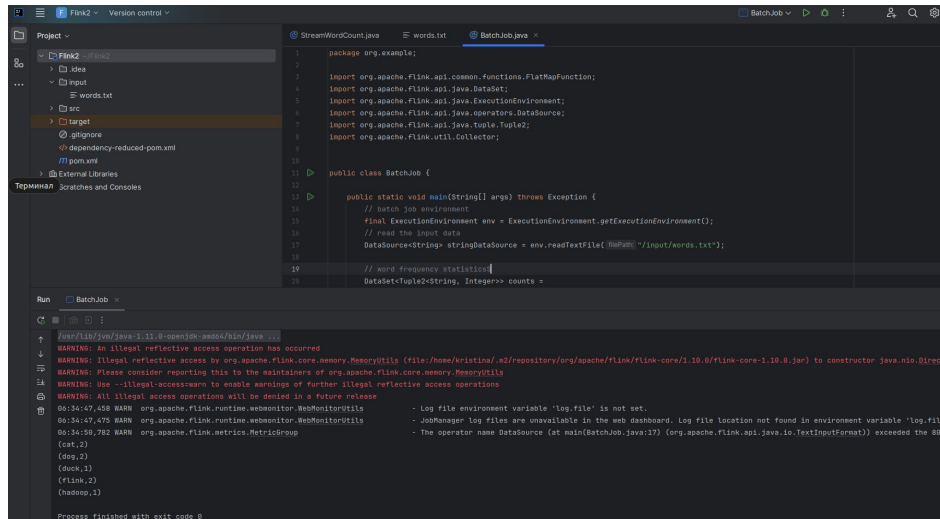


Figure 10: BatchJob execution

BoundedStreamWordCount experiment

```

1 package org.example;
2
3 import org.apache.flink.api.common.typeinfo.Types;
4 import org.apache.flink.api.java.tuple.Tuple2;
5 import org.apache.flink.streaming.api.datastream.DataStreamSource;
6 import org.apache.flink.streaming.api.datastream.KeyedStream;
7 import org.apache.flink.streaming.api.datastream.SingleOutputStreamOperator;
8 import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment;
9 import org.apache.flink.util.Collector;
10
11 import java.lang.reflect.Type;
12
13
14 public class BoundedStreamWordCount {
15
16     public static void main(String[] args) throws Exception {
17         final StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();
18
19         DataStreamSource<String> lineDataStreamSource = env.readTextFile("input/words.txt");
20
21         SingleOutputStreamOperator<Tuple2<String, Long>> wordAndOneTuple = lineDataStreamSource.flatMap((
22             String line, Collector<Tuple2<String, Long>> out) -> {
23             String[] words = line.split(" ");
24             for (String word: words){
25                 out.collect(Tuple2.of(word, 1L));
26             }
27         })
28         .returns(Types.TUPLE(Types.STRING, Types.LONG));
29

```



```

30    KeyedStream<Tuple2<String, Long>, String> wordAndOneKeyedStream = wordAndOneTuple.keyBy(data ->
    data.f0);
31
32    SingleOutputStreamOperator<Tuple2<String, Long>> sum = wordAndOneKeyedStream.sum(1);
33
34    sum.print();
35
36    env.execute();
37 }
38 }

```

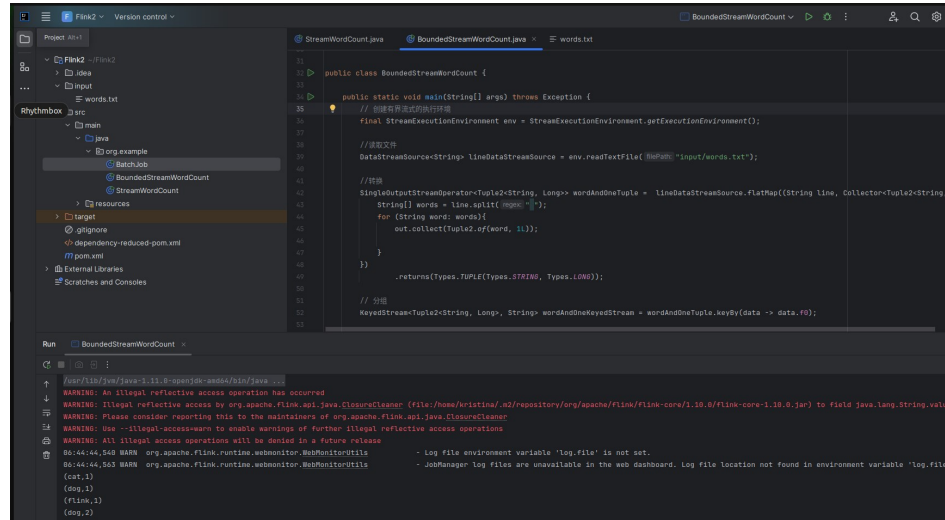


Figure 11: BoundedStreamWordCount execution

InfiniteStreamingJob (StreamWordCount) experiment

In this experiment after running the code we can see the process of the real-time streaming data statistics updates. In the end we get output file: "output1.csv".

```
1 package org.example;
2
3 import org.apache.flink.api.common.functions.FlatMapFunction;
4 import org.apache.flink.streaming.api.datastream.DataStream;
5 import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment;
6 import org.apache.flink.api.java.tuple.Tuple2;
7 import org.apache.flink.util.Collector;
8
9 public class InfiniteStreamingJob {
10
11     public static void main(String[] args) throws Exception {
12         final StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();
13
14         if (args.length != 1) {
15             System.out.println("Usage: StreamingJob <output path>");
16             return;
17         }
18         String outputPath = args[0];
19
20         DataStream<String> textStream = env.socketTextStream("localhost", 9999);
21
22         DataStream<Tuple2<String, Integer>> counts = textStream
23             .flatMap(new Tokenizer())
24             .keyBy(value -> value.f0)
25             .sum(1);
26
27         counts.print();
28
29         counts.writeAsText(outputPath).setParallelism(1);
30
31         env.execute("Infinite Streaming Word Count Job");
32     }
33
34     public static class Tokenizer implements FlatMapFunction<String, Tuple2<String, Integer>> {
35         @Override
36         public void flatMap(String value, Collector<Tuple2<String, Integer>> out) {
37             String[] tokens = value.toLowerCase().split("\\W+");
38             for (String token : tokens) {
39                 if (token.length() > 0) {
40                     out.collect(new Tuple2<>(token, 1));
41                 }
42             }
43         }
44     }
45 }
```

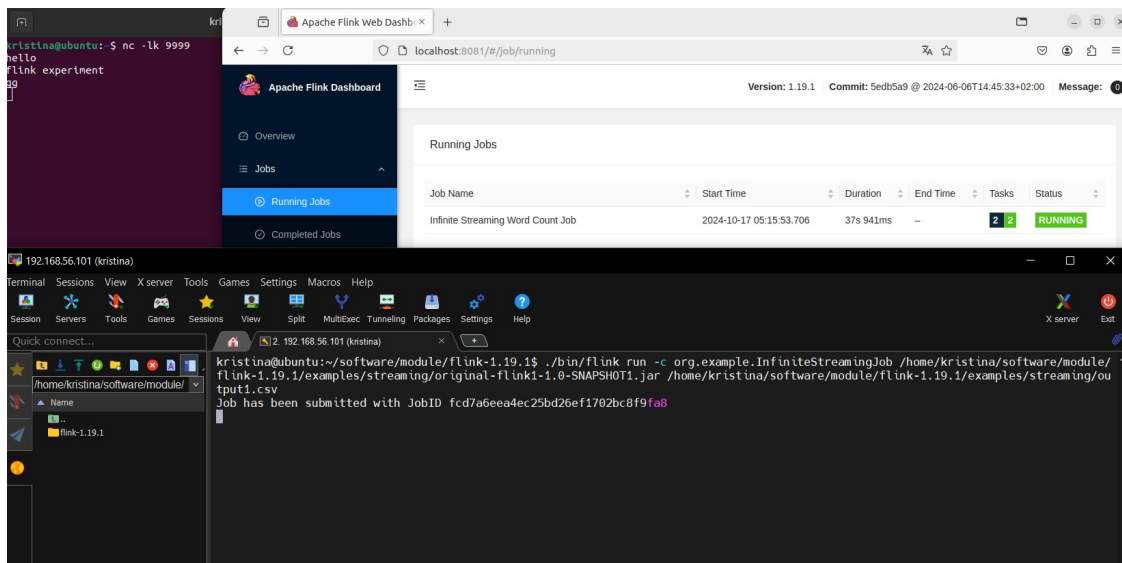


Figure 12: Streaming job execution

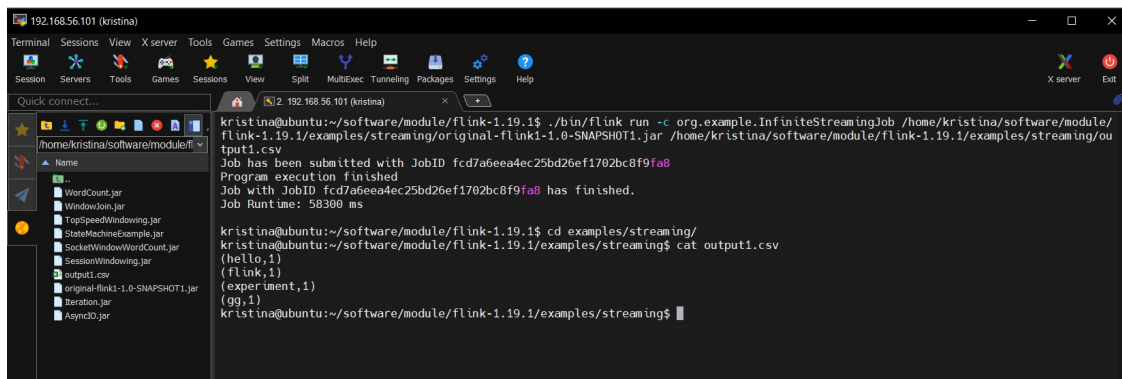


Figure 13: Streaming job result

Optional experiment

TaxiTripDuration experiment

Here we read the train.csv file and output the duration of each trip

```

1 package org.example;
2
3 import org.apache.flink.api.common.functions.MapFunction;
4 import org.apache.flink.api.java.DataSet;
5 import org.apache.flink.api.java.ExecutionEnvironment;
6 import org.apache.flink.api.java.tuple.Tuple2;
7 import org.apache.flink.api.java.io.TextOutputFormat;
8
9 public class TaxiTripDuration {
10     public static void main(String[] args) throws Exception {

```

```

11 final ExecutionEnvironment env = ExecutionEnvironment.getExecutionEnvironment();
12
13 String inputPath = args[0];
14 String outputPath = args[1];
15
16 DataSet<String> source = env.readTextFile(inputPath);
17
18 DataSet<Tuple2<String, Integer>> result = source
19     .map(new MapFunction<String, Tuple2<String, Integer>>() {
20         @Override
21         public Tuple2<String, Integer> map(String line) throws Exception {
22             String[] fields = line.split(",");
23
24             if (fields[0].equals("id")) {
25                 return null;
26
27             try {
28                 String id = fields[0].trim();
29                 Integer tripDuration = Integer.parseInt(fields[10].trim());
30                 System.out.println("Processing ID: " + id + ", Trip Duration: " + tripDuration);
31                 return new Tuple2<>(id, tripDuration);
32             } catch (NumberFormatException | ArrayIndexOutOfBoundsException e) {
33                 System.err.println("Error processing line: " + line + " - " + e.getMessage());
34                 return null;
35             }
36         }
37     })
38     .filter(tuple -> tuple != null);
39
40 DataSet<String> header = env.fromElements("id,trip_duration");
41 DataSet<String> output = header.union(result.map(tuple -> tuple.f0 + "," + tuple.f1));
42
43 output.writeAsText(outputPath);
44
45 env.execute("Taxi Trip Duration");
46 }
47 }
48 }

```

```

... 9 more
kristina@ubuntu:~/software/module/flink-1.19.1$ ./bin/flink run -c org.example.TaxiTripDuration /home/kristina/software/module/flink-1.19.1/examples/streaming/original-flink1-1.0-SNAPSHOT.jar /home/kristina/software/module/train.csv /home/kristina/software/module/flink-1.19.1/examples/streaming/output2.csv
Job has been submitted with JobID 5bf863d4eac925f3db7d021f00bcb867
Program execution finished
Job with JobID 5bf863d4eac925f3db7d021f00bcb867 has finished.
Job Runtime: 13949 ms
kristina@ubuntu:~/software/module/flink-1.19.1$

```

Figure 14: Taxi trip Duration execution

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P
1	id	trip duration														
2	id2875421	455														
3	id2377394	663														
4	id3858529	2124														
5	id3504673	429														
6	id2181028	435														
7	id0801584	443														
8	id1813257	341														
9	id1324603	1551														
10	id1301050	255														
11	id0012891	1225														
12	id1436371	1274														
13	id1299289	1128														
14	id1187965	1114														
15	id0799785	260														
16	id2900608	1414														
17	id3319787	211														
18	id3379579	2316														
19	id1154431	731														
20	id3552682	1317														
21	id3390316	251														
22	id2070428	486														
23	id0809232	652														
24	id2352683	423														
25	id1603037	1163														
26	id3321406	2485														
27	id0129640	1283														
28	id3587298	1130														
29	id2104175	694														
30	id3973319	892														

Figure 15: Taxi trip duration result

PassengerCount experiment

In this code we read train.csv file and output passengers' count.

```

1 package org.example;
2
3 import org.apache.flink.api.common.functions.MapFunction;
4 import org.apache.flink.api.common.functions.ReduceFunction;
5 import org.apache.flink.streaming.api.datastream.DataStream;
6 import org.apache.flink.streaming.api.environment.StreamExecutionEnvironment;
7 import org.apache.flink.api.java.tuple.Tuple2;
8
9 public class PassengerCount {
10     public static void main(String[] args) throws Exception {
11         final StreamExecutionEnvironment env = StreamExecutionEnvironment.getExecutionEnvironment();
12
13         String inputPath = args[0];
14         String outputPath = args[1];
15
16         DataStream<String> source = env.readTextFile(inputPath);
17
18         DataStream<Tuple2<Integer, Integer>> result = source
19             .map(new MapFunction<String, Tuple2<Integer, Integer>>() {
20                 @Override
21                 public Tuple2<Integer, Integer> map(String line) throws Exception {
22                     String[] fields = line.split("\\s*,\\s*");
23                     if (fields.length < 6) {
24                         return null;
25                     }

```

```

26
27     if (fields[0].equals("id")) {
28         return null;
29     }
30
31     try {
32         int passengerCount = Integer.parseInt(fields[4].trim());
33         System.out.println("Processing Passenger Count: " + passengerCount);
34         return new Tuple2<>(passengerCount, 1);
35     } catch (NumberFormatException e) {
36         System.err.println("Error processing line: " + line + " - " + e.getMessage());
37         return null;
38     }
39 }
40 })
41 .filter(tuple -> tuple != null)
42 .keyBy(tuple -> tuple.f0)
43 .reduce(new ReduceFunction<Tuple2<Integer, Integer>>() {
44     @Override
45     public Tuple2<Integer, Integer> reduce(Tuple2<Integer, Integer> value1, Tuple2<Integer,
46     Integer> value2) throws Exception {
47         return new Tuple2<>(value1.f0, value1.f1 + value2.f1);
48     }
49 });
50
51 result.map(tuple -> "passenger_count,count\n" + tuple.f0 + "," + tuple.f1)
52     .writeAsText(outputPath);
53
54 env.execute("Passenger Count");
55 }

```

```

kristina@ubuntu:~/software/module/flink-1.19.1$ ./bin/flink run -c org.example.PassengerCount /home/kristina/software/module/flink-
1.19.1/examples/streaming/flink1-1.0-SNAPSHOT.jar /home/kristina/software/module/train.csv /home/kristina/software/module/flink-1.1
9.1/examples/streaming/output3.csv
Job has been submitted with JobID 506b8b2e1baa64d8dfc910de9f95876e
Program execution finished
Job with JobID 506b8b2e1baa64d8dfc910de9f95876e has finished.
Job Runtime: 19166 ms

```

Figure 16: Passenger Count execution

output3.csv - LibreOffice Calc

Файл Правка Вид Вставка Формат Стили Лист Данные Сервис Окно Справка

Liberation Sans 10 пт Ж К Ч . A . % 74 00 00

A1 fx S

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O
1	S	count													
2	1	1													
3	passenger_count	count													
4	1	2													
5	passenger_count	count													
6	1	3													
7	passenger_count	count													
8	1	4													
9	passenger_count	count													
10	1	5													
11	passenger_count	count													
12	6	1													
13	passenger_count	count													
14	4	1													
15	passenger_count	count													
16	1	6													
17	passenger_count	count													
18	1	7													
19	passenger_count	count													
20	1	8													
21	passenger_count	count													
22	1	9													
23	passenger_count	count													
24	4	2													
25	passenger_count	count													
26	2	1													
27	passenger_count	count													
28	1	10													
29	passenger_count	count													
30	1	11													
31	passenger_count	count													
32	1	12													
33	passenger_count	count													
34	1	13													
35	passenger_count	count													
36	1	14													
37	passenger_count	count													
38	1	15													
39	passenger_count	count													
40	1	16													

Figure 17: Passenger Count result