# ID2221 - Lab 2

Daniel Sharifi          Antonios Mantzaris
dsharifi@kth.se          mantza@kth.se

## General Comments

Our delivery is structured within 3 separate folders, one for each task in the lab. Task 1 and 2 can be run by running "sbt run" inside the task folder. We completed task3 with a DataBricks notebook, and thus we simply saved the notebook as a HTML file to display the results.

For task1 and task2 we see that the average values for all keys are converging towards 12.5, which is the expected value, considering each value is randomly sampled with equal probability in the range [1, 26].

## Task 1

**Comments:**
We process the stream in task1 by processing the records in a stateful manner, with mapWithState(). To achieve this we created a case class Keeper, which keeps tracks of the cumulative sum and number of records. Each time a new record is seen in the stream the state object will be updated by increasing its count by one, and adding the new record value to the cumulative sum. Once the state is updated the mapping function returns the average by dividing the cumulative sum with the number of records.

**Results:**

```
Time: 1633898985000 ms
-----------------------------------
(d,3.0)
(h,8.0)
(p,6.0)
(h,4.5)
(r,7.0)
(v,17.0)
(p,6.5)
(l,14.0)
(t,3.0)
(j,18.0)
...
```

```
Time: 1633898987000 ms
-----------------------------------
(z,12.561371841155236)
(b,12.423779277491068)
(t,12.561000406669377)
(j,12.427586206896551)
(l,12.647649658497388)
(l,12.649126330588471)
(b,12.42627505457432)
(f,12.431905049285858)
(l,12.649196787148595)
(r,12.705641864268193)
...
```

```
Time: 1633898986000 ms
-----------------------------------
(d,12.20854453294714)
(h,12.331252295262578)
(r,12.754104341481211)
(v,12.331755280407865)
(f,12.435916002896452)
(n,12.50266240681576)
(x,12.692035398230088)
(b,12.320307048150733)
(l,12.619748653500897)
(t,12.540805604203152)
...
```
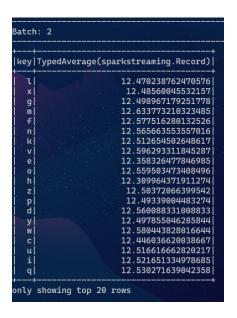
# Task2

**Comments:**

After parsing the key value pairs into typed classes, we simply computed average by grouping records by key, and using the building aggregate tool to compute averages:

averages = ds.groupByKey(_.key).agg(typed.avg(_.value))

Aggregations such as avg are stateful by default in spark structured streaming[0].

**Results:**

# Task 3

Open the html file attached in the task3 folder.

## References:

[0]
http://spark.apache.org/docs/latest/structured-streaming-programming-guide.html#arbitrary-stateful-operations