# Lab 3 - Data Intensive Computing

Stream Processing with Flink

## Windows and Watermarks

In an ideal stream, every time the end of a window is hit, one can assume that all the records belonging to that window have arrived and hence the processing of that window can be started immediately. However, due to issues such as network latency and processing power, it can happen that records arrive past the closing time of a window. To solve this issue, the sources of these events can emit watermarks: when a watermark is received we can assume that all the records before the watermark have been observed. With this indication, it is enough to wait for the right watermark before processing the contents of a window.

## Preprocessing

With this preprocessing, the latitude and longitude of each ride are converted into cell ids, then only the rides that started or ended in an terminal  of the airport are kept.

```scala
    private val extractLatLon = (ride: TaxiRide) =>
     if (ride.isStart)
       (ride.startLon, ride.startLat)
     else
       (ride.endLon, ride.endLat)


    private val extractTerminal = (lon: Float, lat: Float) =>
     TerminalUtils.gridToTerminal(GeoUtils.mapToGridCell(lon, lat))

    val maxDelay = 60 // events are out of order by max 60 seconds
    val speed = 600 // events of 10 minutes are served in 1 second

    val env = StreamExecutionEnvironment.getExecutionEnvironment
    env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
    val rides = env.addSource(new TaxiRideSource("nycTaxiRides.gz", maxDelay,
    speed))

    val jfkRides = rides
     .map(extractLatLon)
     .map(extractTerminal.tupled)
     .filter(_ != Terminal_404)
```

# Problem 1 - Number of taxis in each terminal by hour

The stream of ride events is grouped by terminal and windowed with a tumbling window of one hour. Then, every terminal-window pair is processed to output a tuple with the terminal, the ride events count and the start of the window in hours.

```scala
val windowedCountsByTerminal = jfkRides
  .keyBy(terminal => terminal)
  .window(TumblingEventTimeWindows.of(Time.hours(1)))
  .apply((terminal: Terminal, window, records, out: Collector[(Terminal, Int, Int)]) => {
    val cal: Calendar = Calendar.getInstance()
    cal.setTimeZone(TimeZone.getTimeZone("America/New_York"))
    cal.setTimeInMillis(window.getStart)
    val hour = cal.get(Calendar.HOUR_OF_DAY)
    out.collect((terminal, records.size, hour))
  })
windowedCountsByTerminal.print()
```

# Problem 2 - Busiest terminal by hour

In order to find the busiest terminal for every hour, the stream of rides is windowed without grouping on the terminal id. Then, for every window a count of rides per terminal is computed. Finally, only the terminal with the highest rides count is emitted.

```scala
val windowedBusiestTerminal = jfkRides
  .timeWindowAll(Time.hours(1))
  .apply((window, records, out: Collector[(Terminal, Int, Int)]) => {
    val cal: Calendar = Calendar.getInstance()
    cal.setTimeZone(TimeZone.getTimeZone("America/New_York"))
    cal.setTimeInMillis(window.getStart)
    val hour = cal.get(Calendar.HOUR_OF_DAY)
    val counts = records
      .foldLeft(HashMap.empty[Terminal, Int].withDefaultValue(0))((acc, terminal) => {
        acc(terminal) += 1
        acc
      })
    val busiestTerminal = counts.maxBy(_._2)
    out.collect((busiestTerminal._1, busiestTerminal._2, hour))
  })
windowedBusiestTerminal.print()
```

```scala
1   package lab3.problems
2
3   import java.util.{Calendar, TimeZone}
4
5   import com.dataartisans.flinktraining.exercises.datastream_java.datatypes.TaxiRide
6   import com.dataartisans.flinktraining.exercises.datastream_java.sources.TaxiRideSource
7   import com.dataartisans.flinktraining.exercises.datastream_java.utils.GeoUtils
8   import org.apache.flink.streaming.api.TimeCharacteristic
9   import org.apache.flink.streaming.api.scala._
10  import org.apache.flink.streaming.api.windowing.assigners.TumblingEventTimeWindows
11  import org.apache.flink.streaming.api.windowing.time.Time
12  import org.apache.flink.util.Collector
13
14  import scala.collection.mutable.HashMap
15
16  object JfkTerminals {
17    private val extractLatLon = (ride: TaxiRide) =>
18      if (ride.isStart)
19        (ride.startLon, ride.startLat)
20      else
21        (ride.endLon, ride.endLat)
22
23    private val extractTerminal = (lon: Float, lat: Float) =>
24      TerminalUtils.gridToTerminal(GeoUtils.mapToGridCell(lon, lat))
25
26    def main(args: Array[String]) {
27      val maxDelay = 60 // events are out of order by max 60 seconds
28      val speed = 600 // events of 10 minutes are served in 1 second
29
30      val env = StreamExecutionEnvironment.getExecutionEnvironment
31      env.setStreamTimeCharacteristic(TimeCharacteristic.EventTime)
32      val rides = env.addSource(new TaxiRideSource("nycTaxiRides.gz", maxDelay, speed))
33
34      val jfkRides = rides
35        .map(extractLatLon)
36        .map(extractTerminal.tupled)
37        .filter(_ != Terminal_404)
38
39      val windowedCountsByTerminal = jfkRides
40        .keyBy(terminal => terminal)
41        .window(TumblingEventTimeWindows.of(Time.hours(1)))
42        .apply((terminal: Terminal, window, records, out: Collector[(Terminal, Int, Int)]) => {
43          val cal: Calendar = Calendar.getInstance()
44          cal.setTimeZone(TimeZone.getTimeZone("America/New_York"))
45          cal.setTimeInMillis(window.getStart)
46          val hour = cal.get(Calendar.HOUR_OF_DAY)
47          out.collect((terminal, records.size, hour))
48        })
49      windowedCountsByTerminal.print()
50
51      val windowedBusiestTerminal = jfkRides
52        .timeWindowAll(Time.hours(1))
53        .apply((window, records, out: Collector[(Terminal, Int, Int)]) => {
54          val cal: Calendar = Calendar.getInstance()
55          cal.setTimeZone(TimeZone.getTimeZone("America/New_York"))
56          cal.setTimeInMillis(window.getStart)
57          val hour = cal.get(Calendar.HOUR_OF_DAY)
58
59          val counts = records
60            .foldLeft(HashMap.empty[Terminal, Int].withDefaultValue(0))((acc, terminal) => {
61              acc(terminal) += 1
62              acc
63            })
64          val busiestTerminal = counts.maxBy(_._2)
65          out.collect((busiestTerminal._1, busiestTerminal._2, hour))
66        })
67      windowedBusiestTerminal.print()
68
69      // execute program
70      env.execute("Airport Terminals")
71    }
72  }
```

```scala
package lab3.problems

sealed trait Terminal{def grid: Int}
case object Terminal_1 extends Terminal {val grid = 71436}
case object Terminal_2 extends Terminal {val grid = 71688}
case object Terminal_3 extends Terminal {val grid = 71191}
case object Terminal_4 extends Terminal {val grid = 70945}
case object Terminal_5 extends Terminal {val grid = 70190}
case object Terminal_6 extends Terminal {val grid = 70686}
case object Terminal_404 extends Terminal {val grid = -1}

object TerminalUtils {
  val terminals: Set[Terminal] = Set(Terminal_1, Terminal_2, Terminal_3, Terminal_4, Terminal_5, Terminal_6)

  def gridToTerminal(gridCell: Int): Terminal = {
    terminals.find(t => t.grid == gridCell) match {
      case Some(terminal) => terminal;
      case None => Terminal_404;
    }
  }

}
```