# ▾ Let's have some fun with code

In previous introduction notebook we showed you, how simple is it to use Kafka by yourself.

In order to stream PID data, we used something very similar. Downloading a data, parse them and send to broker server, which is being hosted on AWS. So basicaly we producing data and now you have to consume them.

We will use pyspark and [spark structured streaming API](#).

It is not the only one option, how to consume kafka topics. Kafka itself has APIs to manage it fully, see more [here](#).

First, we have to connect to Kafka's broker. The broker is hosted as [MSK Kafka AWS Service](#).

1. *broker1 b-1-public.bdffelkafka.3jtrac.c19.kafka.us-east-1.amazonaws.com:9196*
2. *broker2 b-2-public.bdffelkafka.3jtrac.c19.kafka.us-east-1.amazonaws.com:9196*

You will need user name and password, provided on previous lectures.

We have 5 topics from which we can read: trams, buses, regbuses, trains and boats, we will start with trams for now.

```
1 # connect to broker
2 JAAS = 'org.apache.kafka.common.security.scram.ScramLoginModule required username="USERXX" password="PWDXX";'
3 tram_stream_topic = spark.readStream \
4   .format("kafka")\
5   .option("kafka.bootstrap.servers", "b-2-public.bdffelkafka.3jtrac.c19.kafka.us-east-1.amazonaws.com:9196," \
6   "b-1-public.bdffelkafka.3jtrac.c19.kafka.us-east-1.amazonaws.com:9196") \
7   .option("kafka.sasl.mechanism", "SCRAM-SHA-512")\
8   .option("kafka.security.protocol", "SASL_SSL") \
9   .option("kafka.sasl.jaas.config", JAAS) \
10  .option("subscribe", "trams") \
11  .load()
```

It is possible to cast json messages directly, using schema, or you can cast it just to string, but it might complicate your work later.

Use schema below. It is possible to save it to other python or dbx notebook and call it externally, see example below for dbx notebook.

```
%run "./pid_schema" # the notebook's name with function in it
schema_pid = get_pid_schema() # use it for casting later
```

To call the schema, it must be wrapped into function with return.

The schema is uploaded on github, copy paste it to separate notebook, or right away in the same notebook.

```
1
2 from pyspark.sql.functions import from_json, col
3 base_trams = tram_stream_topic.select(from_json(col("value").cast("string"), schema_pid).alias("data")).select("data.*") \
```

Let's start the actual spark stream. There are several ways, how to store data, which output mode - what type of sink we will use, documentation is [here](#).

Now we use format `memory` - data will be stored in memory and we will append, ie. we do not wait for complete data (in specified batch or so).

Memory sink is nice for debugging, but you have to be really sure, that data are not big too much. Data are being saved in memory table on Spark's driver.

Supported mode is `Append` or `Complete`.

```
1 #
2 tram_stream_mem_append = base_trams.writeStream \
3         .format("memory")\
4         .queryName("mem_trams")\
5         .outputMode("append")\
6         .start()
```

Using files is useful, when you'd like to append data and save bigger amount and process it later.

You can pull hourly data from stream to your file and analyze it after few days. This can be done by running trigger option and scheduling notebook to run every hour in Workflows menu.

```
1 tram_stream_file = base_trams.writeStream \
```

```
2        .format("parquet")\
3        .option("path", "path/to/destination/dir")
4        .start()
5 #.trigger(once=True).format("delta").queryName(query_name).outputMode("append")
6 # .option("checkpointLocation", "/Filestore/whatever/").toTable("nameoftable")
```

We can try now some basic SQL operations on the mem_trams table.

```
1 %sql
2 select * from mem_trams;
```

We can check how many entries do we have.

```
1 %sql
2 select count(*) from mem_trams;
```

Just remember, if you want to make some transformations on the data, you should not do it directly on stream, since you could corrupt the data in it.

```
1 %sql
2 create table data_trams select * from mem_trams;
```

Now we can do some transformation on the data_trams. For example we can print coordinations of the first 10 trams.

```
1 %sql
2 select geometry.coordinates from data_trams limit 10;
3
```

Ok, let's divide an array with coords.

```
1 %sql
2 SELECT
3   cast(data_trams.geometry.coordinates[0] as double) AS x,
```

```
4    cast(data_trams.geometry.coordinates[1] as double) AS y
5 FROM data_trams;
```

To visualize geopoints, you can use geopy, matplotlib, openstreetmaps and osmnx libs... we will be using some of them. First we need to install libraries.

```
1 %pip install osmnx
```
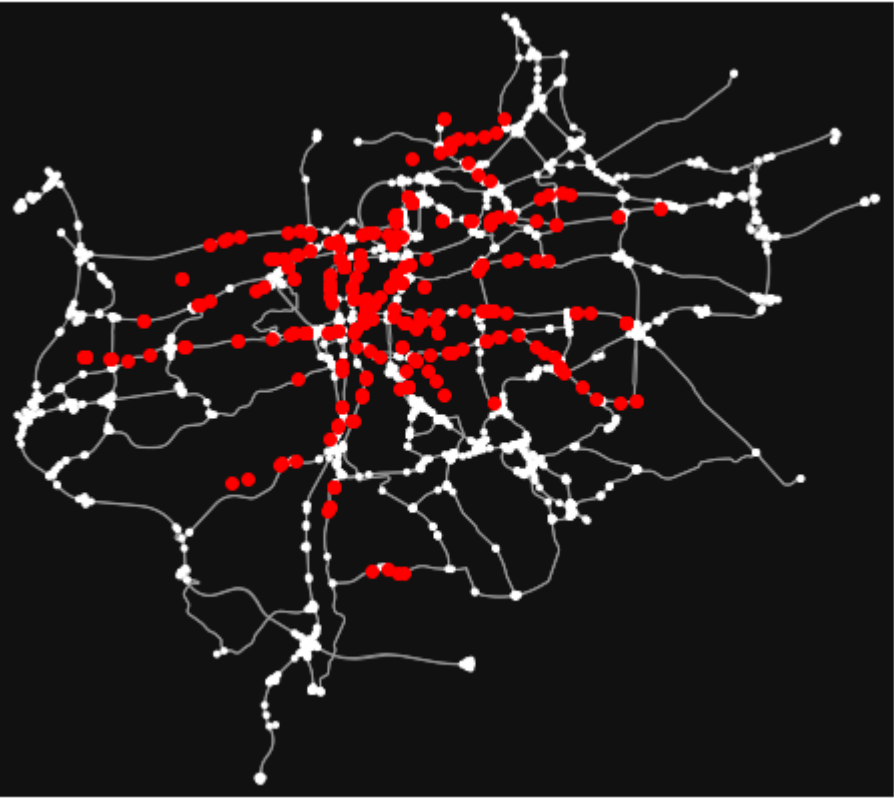
```
1 %pip install numpy==1.23.0
```

Import oxmns and create custom filter (more links you have, longer you will wait). You can create graph by center, by polygon etc. See [OSMNX docs](#)

```
1 import osmnx as ox
2
3 custom_filter='["highway"~"motorway|motorway_link|trunk|trunk_link|primary|primary_link|secondary|secondary_link|road|road_link"]'
4
5 G = ox.graph_from_place("Praha, Czechia", custom_filter=custom_filter)
```

Once you graph is loaded, you can plot it out together with your x and y coordintaes (obtainted via select Spark above).

```
1 import matplotlib.pyplot as plt
2 # this makes your plot wait and not closing
3 fig, ax = ox.plot_graph(G, show=False, close=False)
4 df_geo_p = df_geo.toPandas()
5 # you can plot all, or some subsection for quicker result
6 x = df_geo_p.loc[1:300,'x']
7 y = df_geo_p.loc[1:300,'y']
8
9 ax.scatter(x, y, c='red')
10
11 plt.show()
```

The results is the map below.



Have fun with geo spatial data :)