

# 1 Spark Streaming with Real Time Data and Kafka

In this part, you will create a Spark Streaming application that will continuously read text data from a real time source, analyze the text for named entities, and send their counts to Apache Kafka. A pipeline using Elasticsearch and Kibana will read the data from Kafka and analyze it visually.

## Setting up your development environment

You can set up everything on your local machine or host them on the cloud. Following are the resources that you will need:

- You will need to get access to a real-time text data source. Some examples are:
  - **Reddit** - You can access the comments on a subreddit using the PRAW Python library: <https://praw.readthedocs.io/en/stable/>  
There is also a streaming API for this at: [https://praw.readthedocs.io/en/latest/code\\_overview/other/subredditstream.html](https://praw.readthedocs.io/en/latest/code_overview/other/subredditstream.html).  
You are free to choose any subreddit.
  - **NewsAPI** - You can download real time news using the News API available at: <https://newsapi.org/>. There is also a Python wrapper for this available at: <https://newsapi.org/docs/client-libraries/python>
  - **Finnhub** For financial news and analysis, the Finnhub library <https://finnhub.io/> can be used. There is a Python wrapper for this available at <https://pypi.org/project/finnhub-python/>.
- Download Apache Kafka and go through the quickstart steps: <https://kafka.apache.org/quickstart>
- Windows users might want to consider WSL, which gives a Unix-like environment on Windows: <https://docs.microsoft.com/en-us/windows/wsl/install-win10>
- You will need a working Spark cluster. It can be setup locally or you can use a cloud based server.
- I have attached a simple project showing how to pass data between Kafka and a Spark Structured Streaming application that performs word count. Instructions for running are also attached.
- Later, you will also need to set up Elasticsearch and Kibana environment to visualize data. You will need to download Elasticsearch, Kibana, and Logstash from <https://www.elastic.co/downloads/>

## Project Steps

For this project, you will need to perform the following steps:

- Create a Python application that reads from a real-time data source, such as the ones mentioned in the previous part. Some libraries have a streaming version that continuously fetch real-time data. For others, you might have to write a loop that gets real time data at periodic intervals.

- This incoming data should continuously be written to a Kafka topic (let's call it topic1 for illustration).
- Create a PySpark structured streaming application that continuously reads data from the Kafka topic (topic1) and keeps a running count of the named entities being mentioned. You should already know how to extract named entities from text. In this part, you will keep a running count of the named entities being mentioned.
- At the trigger time, a message containing the named entities and their counts should be sent to another Kafka topic (let's call it topic2 for illustration).
- Configure Logstash, Elasticsearch, and Kibana to read from the Kafka topic (topic2) and create a bar plot of the top 10 most frequent named entities and their counts.

## Visualizing the data using Elasticsearch and Kibana

You will need Elasticsearch and Kibana installed to get the data from Kafka and visualize it. Details are available at:

<https://www.elastic.co/downloads>

<https://www.elastic.co/start>

You will need to visualize the count of top 10 named entities being mentioned on your chosen data source. You should be able to view this using a barplot in Kibana.

## What to Submit

You are required to submit the following:

- Your project code file.
- Named entity frequency bar plots taken at several intervals. For example, you can show the top 10 named entities being mentioned after 15, 30, 45, and 60 minutes.
- A report explaining your data source and also a small paragraph explaining what your results indicate.
- A README file indicating how to run your code. Please do not use any hard coded paths - all paths should be specified as arguments.