

1. Spark Streaming with Real Time Data and Kafka and ELK stack

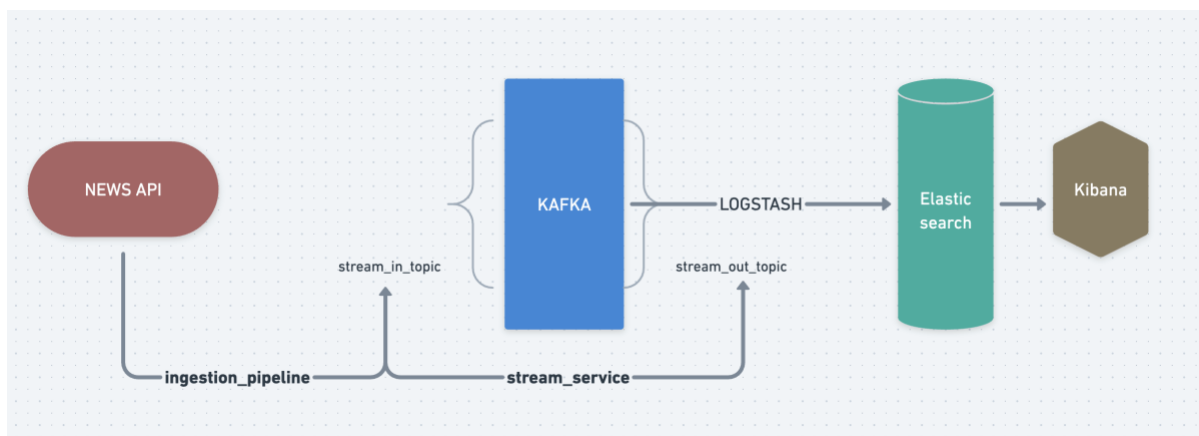
We are performing 'Named Entity Recognition (NER) count' on the live streaming data from [newsapi](https://newsapi.org/) [<https://newsapi.org/>].

Data source -

<https://newsapi.org/docs/endpoints/everything>

This API searches through millions of articles from over 80,000 large and small news sources and blogs. This endpoint suits article discovery and analysis.

Below is the pipeline architecture –



Application configuration -

Project configurations are defined in file **app-config.yml** and can be changed as per need –

```
bootstrap_servers: kafka:29092
input_topic: stream_in_topic
output_topic: stream_out_topic
api_key: f3abe15c03e04488b3ade0a1e5eb1167
news_api_url:
https://newsapi.org/v2/everything?q=the&apiKey={api_key}&sortBy=publishedAt&language=en
checkpoint_location: '/tmp/checkpoint'
news_api_request_delay: 20 # seconds
```

Running steps –

The project is setup in docker-containerized environment and can be easily launched using 'docker-compose up'.

It will automatically download all the base images and will launch the necessary services. Docker daemon has to be running for it. Please see here or more details - <https://docs.docker.com/compose/>

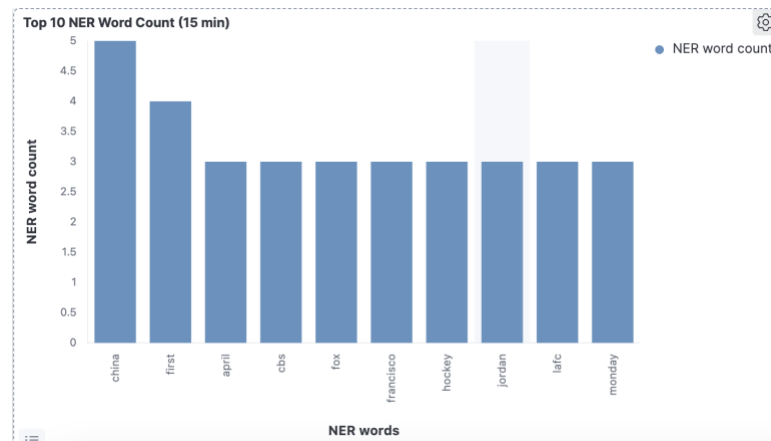
Respective launched applications can be found at –

- Kafka - <localhost:9092> [Inside containers it will be accessible as <kafka:29092>]
- Elasticsearch - <http://localhost:9200/>
- Kibana - <http://localhost:5601/> [Please import [*'export.ndjson'*] file in kibana to import dashboard analysis.]

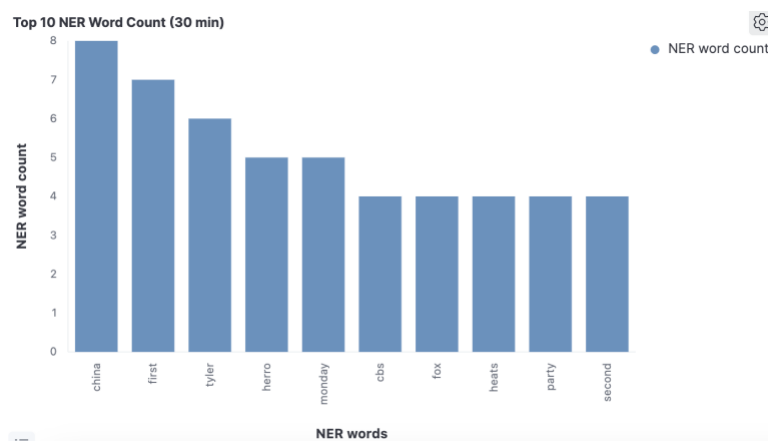
Results –

NER word count analysis has been demonstrated using various plots. Here is the description for them –

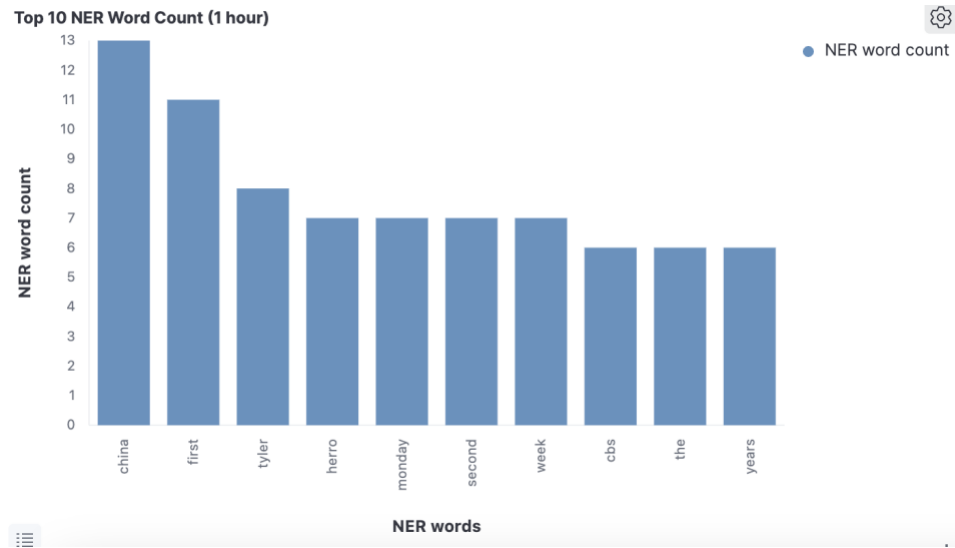
1. Top 10 NER Word Count (All time)



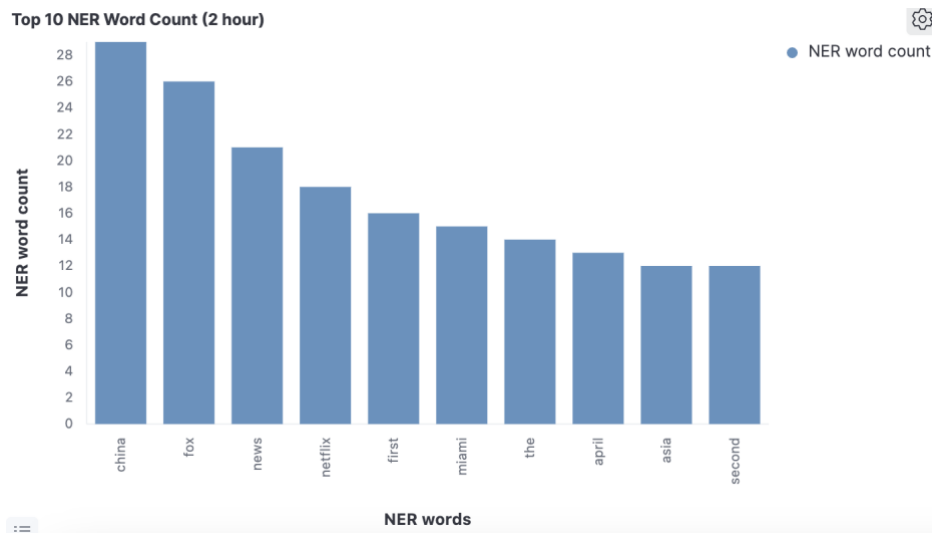
Top 10 NER Word Count (15 min)



Top 10 NER Word Count (30 min)



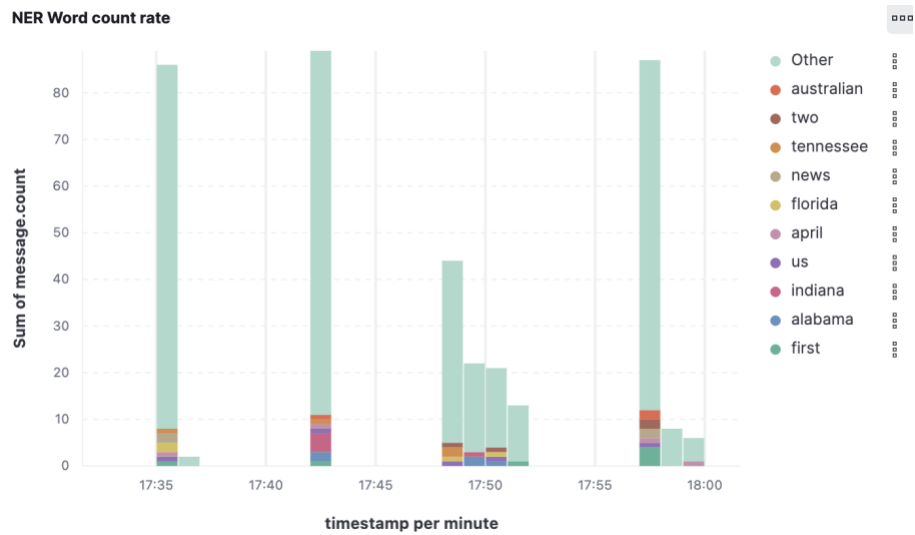
Top 10 NER Word Count (1 hour)



Top 10 NER Word Count (2 hour)

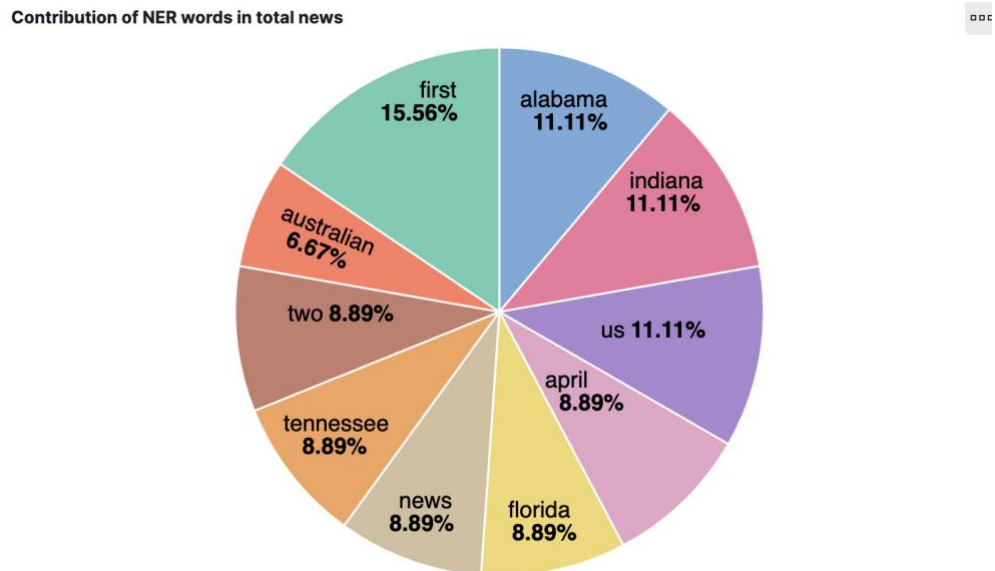
The above graph shows the top 10 NER words that have occurred most frequently in all time.

2. NER Word count rate



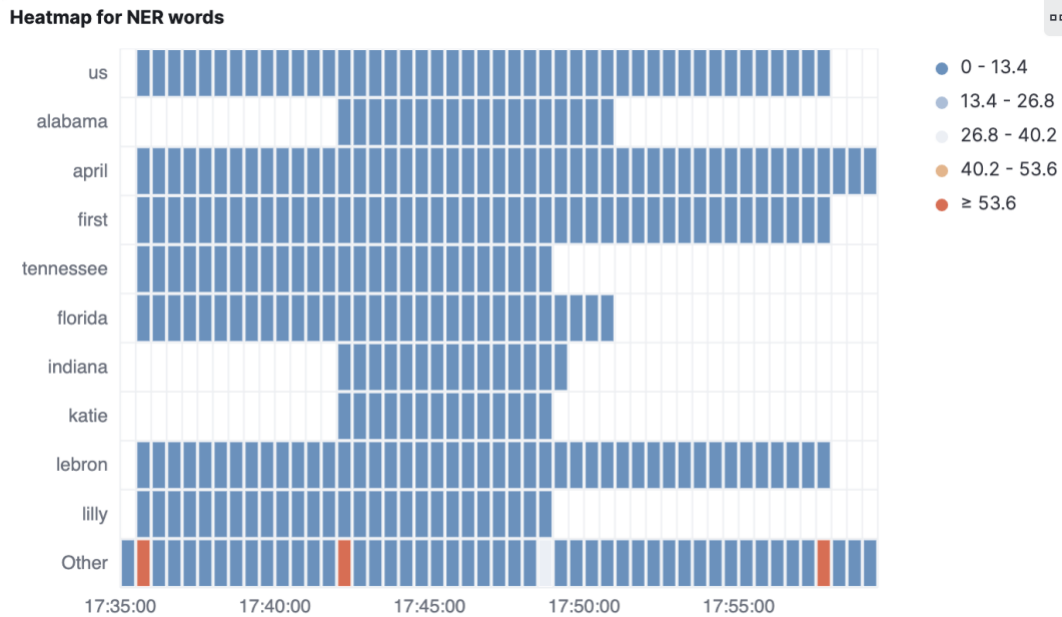
The stacked barplot graph shows the rate of NER words every minute.

3. Contribution of NER words in total news



The pie chart shows the individual contribution of top 10 NER words in all time news.

4. Heatmap for NER words



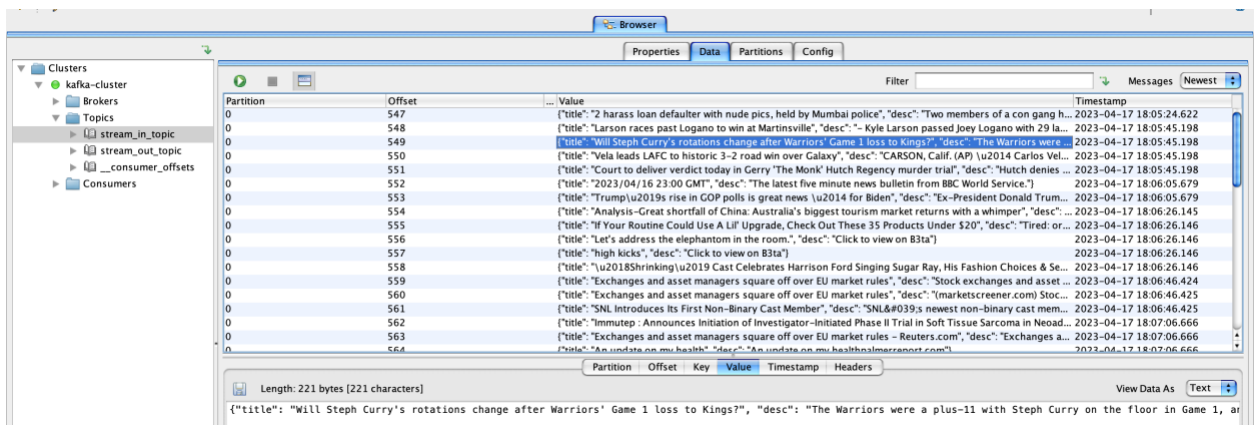
The heatmap shows the frequency of NER words at various timestamps. Heat is represented using color scales.

Result Analysis

The graphs shows the running frequency count of different NER words extracted from real time news. We observe that NER words like 'US', 'florida', 'indiana', 'katie' have been observed in the data at different times.

Here are further screenshots of data at various stages –

1. At input Kafka topic – (stream_in_topic)



2. At output Kafka topic – (stream_out_topic)

Partition	Offset	Key	Value	Timestamp
0	521		["word":"ferrari","count":1]	2023-04-17 18:06:00.444
0	522		["word":"lana","count":1]	2023-04-17 18:06:01.310
0	523		["word":"seattle","count":1]	2023-04-17 18:06:01.580
0	524		["word":"mumbai","count":1]	2023-04-17 18:06:01.580
0	525		["word":"margaret","count":1]	2023-04-17 18:06:01.656
0	526		["word":"rey","count":1]	2023-04-17 18:06:02.441
0	527		["word":"leads","count":1]	2023-04-17 18:06:20.910
0	528		["word":"tafc","count":1]	2023-04-17 18:06:21.385
0	529		["word":"vela","count":1]	2023-04-17 18:06:21.861
0	530		["word":"today","count":1]	2023-04-17 18:06:22.261
0	531		["word":"gop","count":1]	2023-04-17 18:06:40.649
0	532		["word":"trumps","count":1]	2023-04-17 18:06:40.994
0	533		["word":"china","count":1]	2023-04-17 18:07:02.337
0	534		["word":"ford","count":1]	2023-04-17 18:07:03.545
0	535		["word":"first","count":1]	2023-04-17 18:07:20.269
0	536		["word":"snf","count":1]	2023-04-17 18:07:22.539
0	537		["word":"jason","count":1]	2023-04-17 18:07:40.638
0	538		["word":"sant","count":1]	2023-04-17 18:07:40.803

Partition: 0
Offset: 521
Offset Hex: 209
Message Length: 28

3. Inside Elasticsearch –

Time	message.word	message.count
Apr 17, 2023 @ 18:08:42.554	than	1
Apr 17, 2023 @ 18:08:41.622	michigan	1
Apr 17, 2023 @ 18:08:41.622	million	1
Apr 17, 2023 @ 18:08:40.662	australia	1
Apr 17, 2023 @ 18:08:40.614	fox	1
Apr 17, 2023 @ 18:08:40.455	more	1
Apr 17, 2023 @ 18:08:40.225	news	1
Apr 17, 2023 @ 18:08:22.364	weeks	1
Apr 17, 2023 @ 18:08:21.788	hockey	1
Apr 17, 2023 @ 18:08:21.718	canadas	1