

# Real Time Analysis of Twitter hashtags using Apache Spark

- Hritwik Singhal 18UCS055
- Raunak Goyal 18UCC162
- Utkarsh Gupta 18UCC140
- Aarchi Gupta 18UCS156
- Abhay Singhal 18UCS011
- Sameer Gupta 18UCS008



Project Link: [github.com/HritwikSinghal/Spark-tweet](https://github.com/HritwikSinghal/Spark-tweet)

# Keywords

- Twitter API v2
- Apache Spark
- Spark Streaming
- Live tweets streaming
- Tweet Analysis
- Flask
- ApexCharts



# Table of contents

1. Introduction
2. Flow chart Diagram
3. Twitter API v2
4. Limitations of Twitter API
5. Apache Spark
6. Spark Structured Streaming
7. Flask
8. ApexCharts
9. Final Output
10. Use cases of our project
11. Some key points to navigate the code
12. References

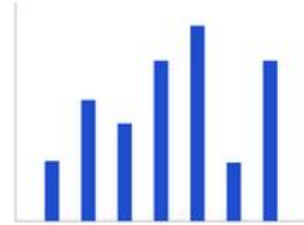
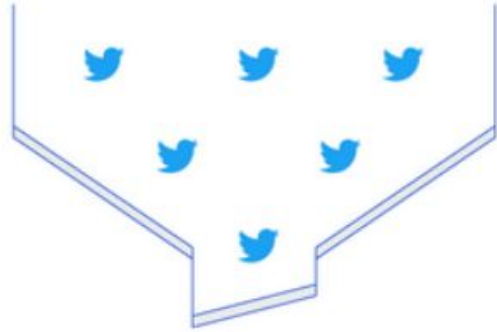


# Introduction

We are Using Apache Spark streaming, Real-Time Analytics engine, to process tweets retrieved from Twitter API and identify the trending hashtags from them based on a certain keywords and, finally, represent the data in a real-time dashboard using flask web framework.



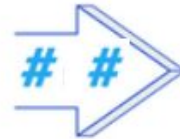
# Flow Diagram



Sockets.io



Spark



Twitter HTTP Client App

# Twitter API v2

- Searching for tweets.
- Using endpoints for finding relevant data more easily.
- There are 2 types of endpoints.
  - Recent Tweets Search (accessible to all)
  - Full Tweets Archive Search (for enterprise only)
- [api.twitter.com/2/tweets/search/recent](https://api.twitter.com/2/tweets/search/recent)
- The data that we are getting from Twitter is send to Spark through Socket.io

# Limitations

- 450 queries per 15 minutes.
- 500K queries per month.
- We cannot get general tweets from Twitter. We have to get tweets based on some keywords.

Explained by: Abhay Singhal

## Request

```
query_params = {'query': keyword,  
               'end_time': end_date,  
               'max_results': max_results,  
               'tweet.fields': 'id,text,author_id,geo,conversation_id,created_at,lang,entities',  
               'next_token': next_token  
}
```

Twitter API endpoint used:

[api.twitter.com/2/tweets/search/recent](https://api.twitter.com/2/tweets/search/recent)

## Response

```
{  
  "data": [  
    {  
      "text": "Looking to get started with the Twitter API? @jessicagarson She'll use  
examples using our v2 endpoints. #TwitterDev #TwitterAPI",  
      "author_id": "2244994945",  
      "id": "1373001119480344583",  
      "lang": "en"  
    }  
  ],  
  "includes": {  
    "users": [  
      {  
        "id": "2244994945",  
        "entities": {  
          "description": {  
            "hashtags": [  
              {  
                "start": 17,  
                "end": 28,  
                "tag": "TwitterDev"  
              },  
              {  
                "start": 105,  
                "end": 116,  
                "tag": "TwitterAPI"  
              }  
            ]  
          }  
        }  
      }  
    ],  
    "created_at": "2013-12-14T04:35:55.000Z",  
    "username": "TwitterDev",  
    "name": "Twitter Dev"  
  }  
],  
  "meta": {  
    "next_token": "1373001119480344583"  
  }  
}
```



# Apache Spark

- Open source unified analytics engine for large scale data processing.
- Utilises in-memory computation and optimize query execution for fast analytic queries for data of any size.
- Used to process the tweets using Apache Spark Streaming to identify hashtags counts.

Explained by: Aarchi Gupta

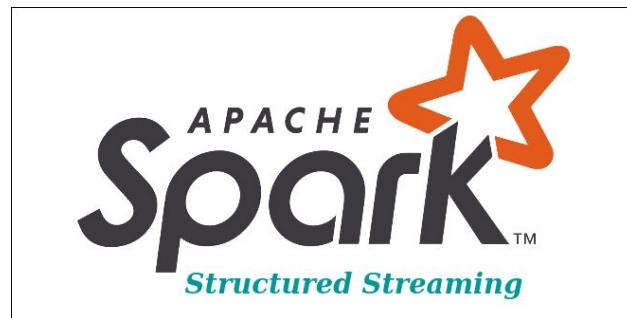
# Spark Streaming



- Extension of core Spark API.
- Used to process real time data .
- Fast in recovery from failures .
- Better load balancing and resource usage.
- Native integration with advances processing libraries.

Explained by: Aarchi Gupta

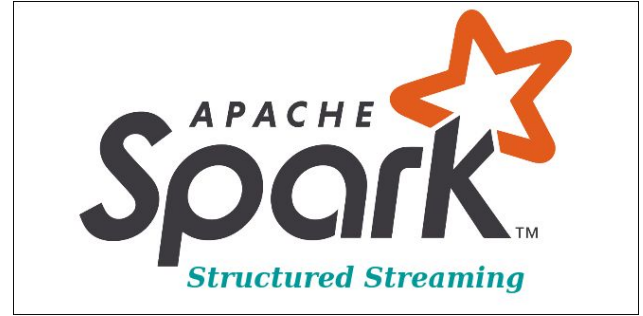
# Structured Streaming



- From the Spark 2.x release onwards
- More high end than streaming and RDD
- Built on the Spark SQL library
- provides more declarative, high-level functional API
- uses DataFrame and Dataset APIs

Explained by: Hritwik Singhal

# Structured Streaming



- We first create a local SparkSession, the starting point of all functionalities related to Spark
- Then we create a streaming DataFrame that represents text data received from a server listening on Localhost using sockets.
- Use spark built-in Split and explode, to split each line into multiple rows.
- Then we Generate running word count using 'groupBy' function of spark.
- At last we send the results from spark to flask Using REST API

Refer to the code in github repo for details

Explained by: Hritwik Singhal



- Flask is a micro web application framework written in Python which reduces development time and allows programmers to build app faster and smarter.
- It consists only of what developers put in it, with no unnecessary code responsible for features we don't use.
- In our scenario, Flask server has collected hashtag data from spark and directed it to the frontend application.

Explained by: Raunak Goyal

- Spark application will POST processed data to server through “updateData” endpoint.
- Frontend application will call “refreshData” endpoint to GET hashtag data-points from server in an interval of every 2 seconds.
- Hashtag data will repeatedly be updated at frontend through Jinja2 template engine.

```
@app.route("/")
def home():
    return render_template('index.html', dataValues=dataValues, categoryValues=categoryValues)

@app.route('/refreshData')
def refresh_graph_data():
    global dataValues, categoryValues
    print("labels now: " + str(dataValues))
    print("data now: " + str(categoryValues))
    # dataValues[0]=dataValues[0]+1
    return jsonify(dataValues=dataValues, categoryValues=categoryValues)

@app.route('/updateData', methods=['POST'])
def update_data_post():
    global dataValues, categoryValues
    if not request.form or 'data' not in request.form:
        return "error", 400
    categoryValues = ast.literal_eval(request.form['label'])

    for i, ele in enumerate(categoryValues):
        try:
            new_ele = re.findall(r'bytearray\(b\'(.*?)\'\)\'', ele)[0]
            print(new_ele)
            categoryValues[i] = new_ele
        except:
            continue

    dataValues = ast.literal_eval(request.form['data'])
    print(f"labels received: {str(categoryValues)}")
    print(f"data received: {str(dataValues)}")
    return "success", 201
```



- ApexCharts is a modern charting library that empowers developers to create dynamic, responsive and interactive visualizations for web pages.
- ApexCharts has a NPM support and works better with bigger datasets.
- In our scenario, Javascript calls the server API to GET hashtag data and then it passes to ApexCharts API to get a responsive and dynamic bar graph.

Explained by: Utkarsh Gupta

templates > <> index.html > ...

```
31 var options = {
32   title: {
33     text: "Trending Hashtags",
34     align: "left",
35     style: {
36       fontSize: "20px"
37     }
38   },
39   chart: {
40     width: "100%",
41     height: 600,
42     type: "bar",
43     foreColor: "#fff",
44   },
45   fill: {
46     type: "gradient",
47     gradient: {
48       gradientToColors: ["#F55555", "#6078ea", "#6094ea"],
49       shade: "dark",
50       type: "horizontal",
51       shadeIntensity: 0.5,
52       inverseColors: true,
53       opacityFrom: 1,
54       stops: [0, 100]
55     }
56   },
57   tooltip: {
58     theme: "dark",
59   },
60   grid: {
61     borderColor: "#40475D"
62   },
63   colors: ["#FCCF31", "#17ead9", "#f02fc2"],
64   plotOptions: {
65     bar: {
66       horizontal: true,
67     }
68   },
69   dataLabels: {
70     enabled: true,
71   },
72   stroke: {
73     width: 0,
74   },
75   series: [
76     {
77       name: 'Count',
78       data: dataValues
79     },
80   ],
81   xaxis: {
82     title: {
83       text: 'Count',
84       style: {
85         fontSize: "17px"
```

templates > <> index.html > ...

```
81 xaxis: {
82   title: {
83     text: 'Count',
84     style: {
85       fontSize: "17px"
86     }
87   },
88   type: 'category',
89   categories: categoryValues,
90   axisBorder: {
91     color: "#333"
92   },
93 },
94 yaxis: {
95   title: {
96     text: 'Hashtag',
97     style: {
98       fontSize: "17px"
99     }
100 },
101 labels: {
102   offsetX: 10
103 },
104 },
105 legend: {
106   position: "right",
107   verticalAlign: "top",
108   containerMargin: {
109     left: 35,
110     right: 60
111   }
112 },
113 responsive: [
114   {
115     breakpoint: 1000,
116     options: {
117       fill: {
118         type: "gradient",
119         gradient: {
120           shade: "dark",
121           type: "vertical",
122           shadeIntensity: 0.5,
123           inverseColors: false,
124           opacityFrom: 1,
125           stops: [0, 100]
126         }
127       },
128       plotOptions: {
129         bar: {
130           horizontal: false
131         }
132       },
133       legend: {
134         position: "bottom"
135       },
136       xaxis: {
137         title: {
138           text: 'Hashtag'
139         },
140         axisBorder: {
141           color: "#333"
142         }
143       },
144       yaxis: {
145         title: {
146           text: 'Count',
147           style: {
148             fontSize: "17px"
149           }
150         }
151       },
152     }
153   },
154 ],
155 },
156 plotOptions: {
157   bar: {
158     horizontal: false
159   }
160 },
161 legend: {
162   position: "bottom"
163 },
164 xaxis: {
```

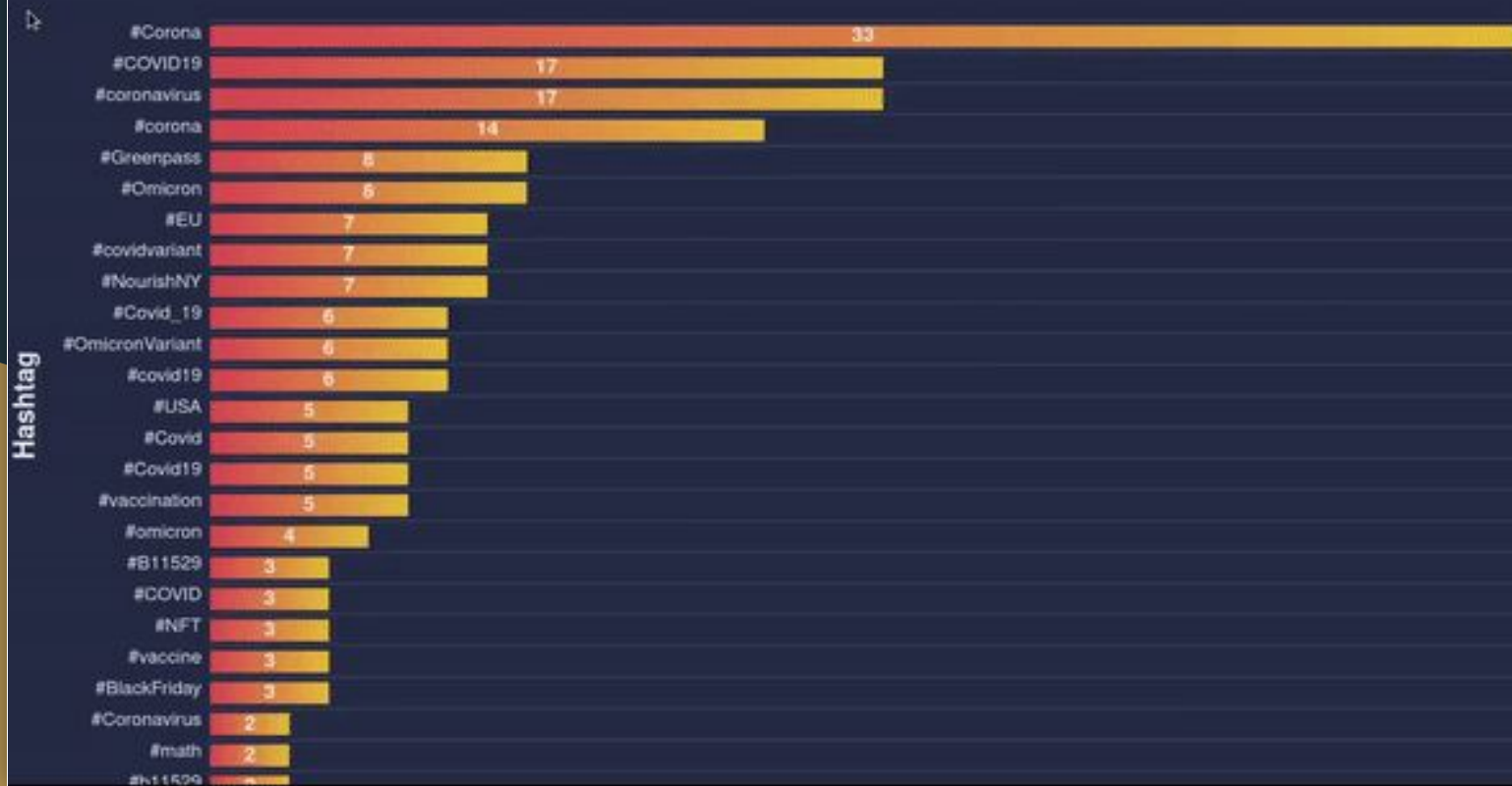
templates > <> index.html > html > body > script

```
109 containerMargin: {
110   left: 35,
111   right: 60
112 },
113 responsive: [
114   {
115     breakpoint: 1000,
116     options: {
117       fill: {
118         type: "gradient",
119         gradient: {
120           shade: "dark",
121           type: "vertical",
122           shadeIntensity: 0.5,
123           inverseColors: false,
124           opacityFrom: 1,
125           stops: [0, 100]
126         }
127       },
128       plotOptions: {
129         bar: {
130           horizontal: false
131         }
132       },
133       legend: {
134         position: "bottom"
135       },
136       xaxis: {
137         title: {
138           text: 'Hashtag'
139         },
140         axisBorder: {
141           color: "#333"
142         }
143       },
144       yaxis: {
145         title: {
146           text: 'Count',
147           style: {
148             fontSize: "17px"
149           }
150         }
151       },
152     }
153   },
154 ],
155 };
156
157 var chart = new ApexCharts(
158   document.querySelector("#responsive-chart"),
159   options
160 );
161
162 chart.render();
163
```

Explained by: Utkarsh Gupta



## Trending Hashtags



Explained by: Utkarsh Gupta

# Use cases of the project

- Marketing a particular product
- Promote contests and giveaways
- Raise awareness about a particular topic
- Build a community around a hashtag
- Discover trending topics



Explained by: Utkarsh Gupta

# Some key points to navigate the code

- Code location : <https://github.com/HritwikSinghal/Spark-tweet/>
- Running the code and other things are explained in README.md
- There are 3 major files in the project name 'twitter\_app.py', 'spark\_app.py' and 'app.py', all of them are at root of the project.
- 'Twitter\_app.py' retrieves data from twitter and sends it to 'spark\_app.py' through sockets.
- 'Spark\_app.py' will run basic spark server using pyspark and process those tweets and finally send the data to 'app.py' using REST API
- 'App.py' is the flask server with multiple endpoints.

Explained by: Hritwik Singhal

# References

- [Twitter API v2 tools & libraries | Docs | Twitter Developer Platform](#)
- [GET /2/tweets/search/recent | Docs | Twitter Developer Platform](#)
- [Overview - Spark 3.2.0 Documentation](#)
- [Spark Streaming - Spark 3.2.0 Documentation](#)
- [Structured Streaming Programming Guide - Spark 3.2.0 Documentation](#)
- [PySpark Documentation — PySpark 3.2.0 documentation](#)
- [Welcome to Flask — Flask Documentation \(2.0.x\)](#)
- [Installation & Getting Started – ApexCharts.js](#)



**Thank You**