# Motivation

We were wondering, "<u>What in the world do people tweet about?!?</u>" 🤔

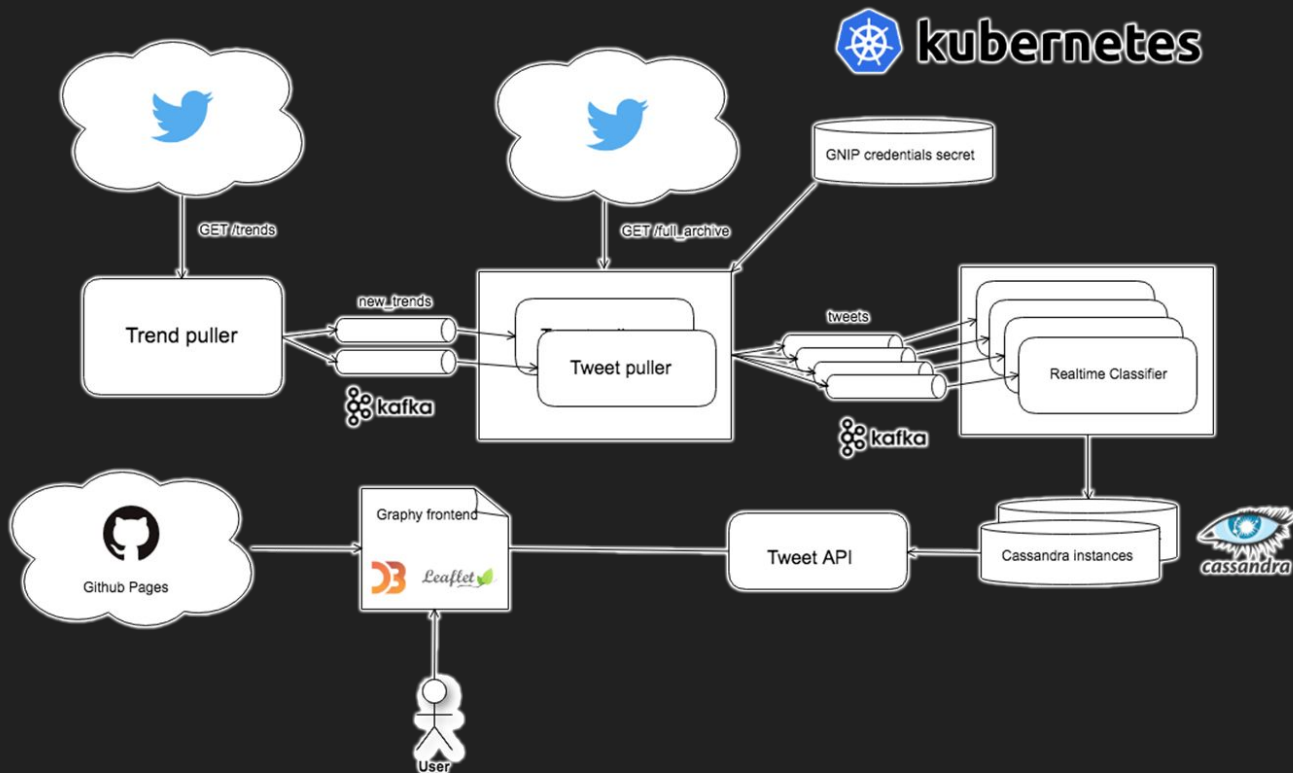Spoiler: Mostly Entertainment and Sports... 😛

Understanding what's happening

We wanted to analyze this further and <u>find the correlation among the popular content</u>.

# Backend + DevOps

# System Architecture

# System Architecture

# Focus on scalability

# System Architecture details

- System details:
    - Kubernetes v1.8.7
    - Deployed with Kops for AWS
- Kubernetes Nodes
    - **Master**: 1 EC2 t2.medium instance
    - **Nodes**: 4 EC2 t2.small instances
- AWS:
    - 2x 100GB EBS (for Cassandra Storage)
    - Elastic Load Balancer to balance if needed to add more APIs servers
- Classifies ~6K tweets/min
- GitHub pages CDN auto-deployment

# Machine Learning

# Classification

Topic Classification

Classifying tweets into the following categories:
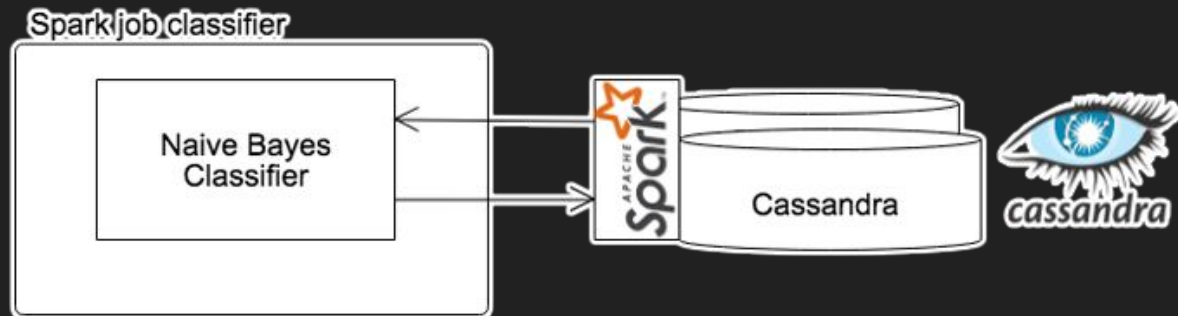
- Sports
- Politics
- Technology
- Mood
- Entertainment

Sentiment Classification - Joy, Sadness, Anger, Neutral.

# Workflow
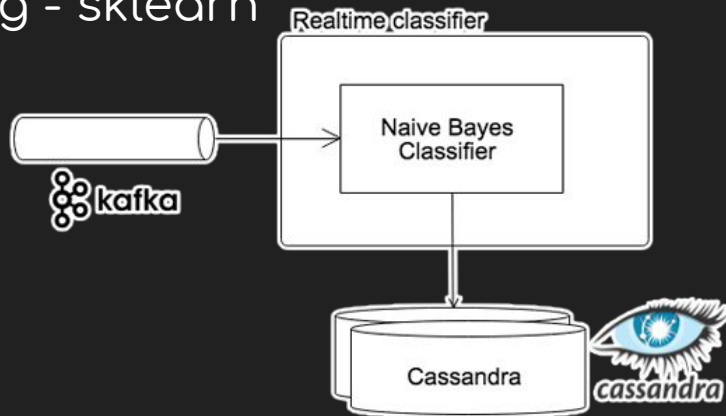
1. Pre-processing
   a. Cleaning Data with regex
   b. Removing twitter specific stop words
   c. Monkey trend labeling
2. Model
   a. Naive Bayes - bag of words model
   b. Suitable for real time pipeline
3. Trend Classification

   - Based on majority classification of it's tweets

4. Results/Accuracy
   a. Precision and recall ~ 93%
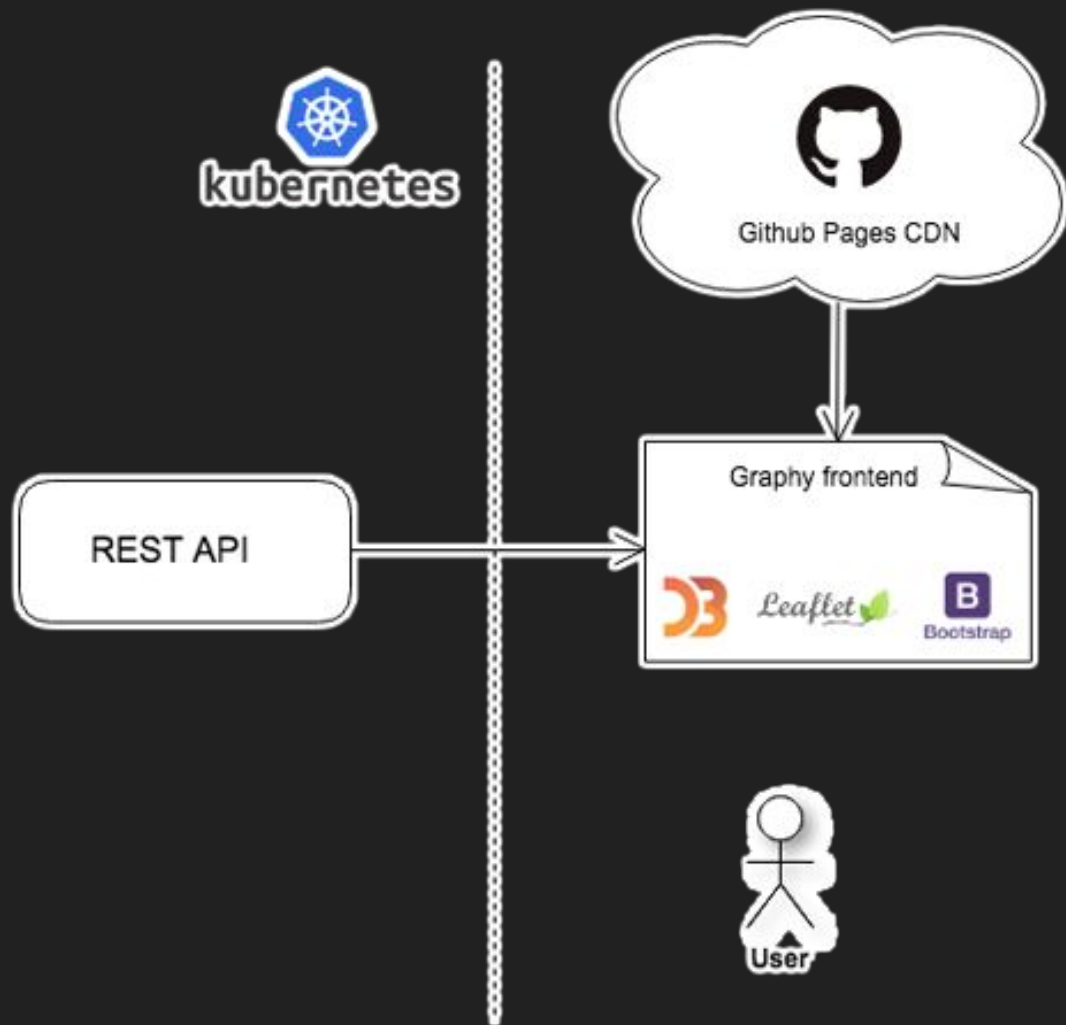   b. 6k classifications per minute

# Infrastructure

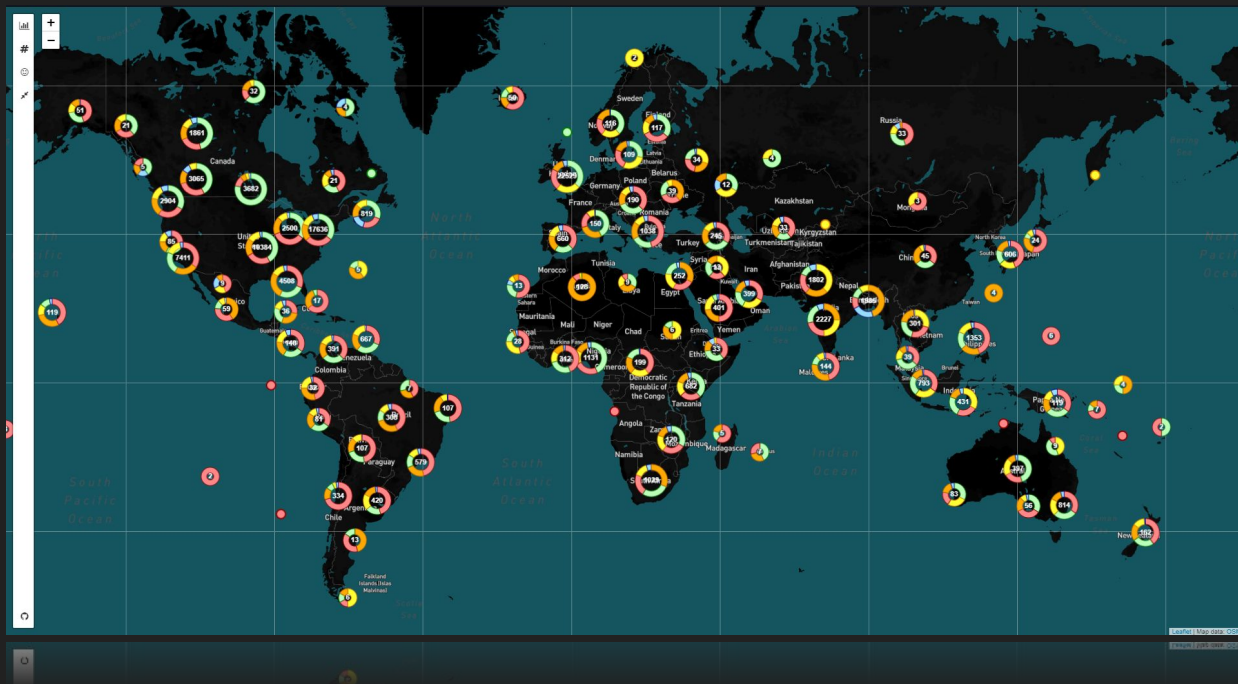- Batch Processing - Spark-mllib



- Real Time Processing - sklearn

# Frontend

kubernetes

Github Pages CDN

Graphy frontend

REST API

User

# Background Map

Map - created with Mapbox (customized URL) and Leaflet.
ClusterPies, Marker Cluster, D3.

# Sidebar Features

- **Tweets by Topic: C3 Chart**
- A D3-based reusable chart library
- Present bar chart of most popular tweets
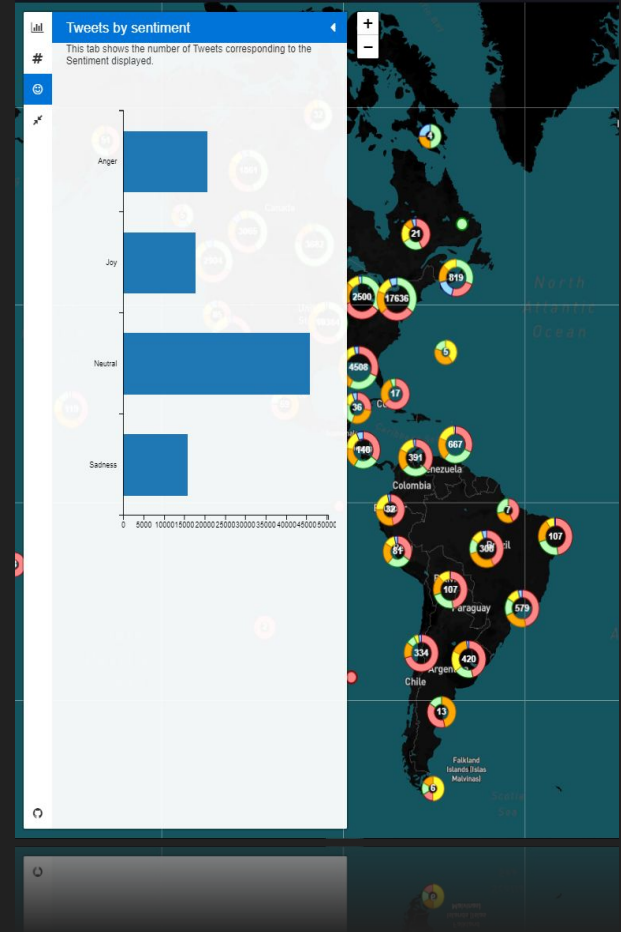
# Sidebar Features

- Trends by Topic: Bootstrap Badges
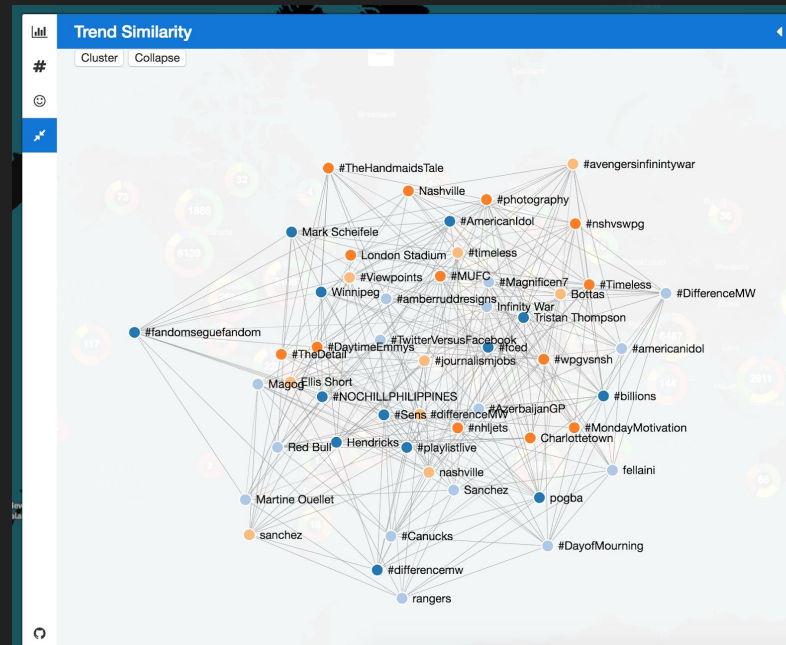- Displays hottest tweets by each trend

# Sidebar Features

- **Sentiment Analysis**
- Emojis for each tweet on map
- Description…if needed…will be added.

# Sidebar Features

- Trend Similarity: D3 Network Graph

- Clustering Trends with
  **Doc2Vec** and
  **Clauset Algorithm**
  - Professor Aaron Clauset's Algorithm

# Issues we struggled upon

- Spark job deployment

- Deploying Kubernetes in AWS is not easy

- D3 learning curve

- Difficulties to label data

  - Non-existent labeled data

  - Tweets are sometimes ambiguous

# Future Work

# Thank You!