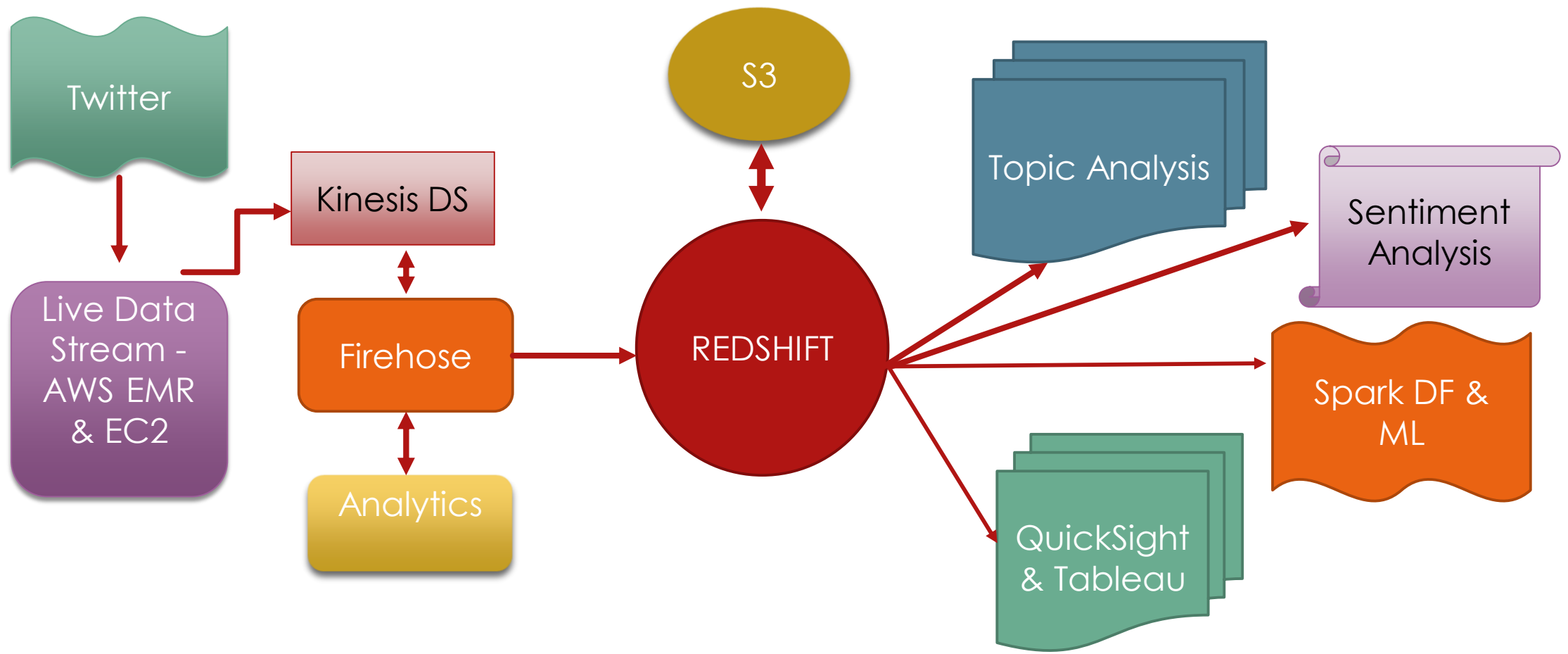# Live Streaming *to* Dashboard *to* Spark ML

– Mohammad Chowdhury

# The Goal

- Build a Big Data Project using data from live twitter stream & find the sentiments & Topics using ML using Big Data technologies like
  - AWS EMR
  - AWS EC2
  - AWS S3
  - AWS Kinesis
  - AWS Redshift
  - AWS QuickSight
  - Spark
  - Tableau
  - Python
  - Linux

# Data Architecture

# Twitter Stream Raw Data sample

# PySpark sample code

```python
from pyspark.sql import SparkSession
from pyspark.sql.types import *

spark = SparkSession.builder.appName('Twitter_Aalysis').getOrCreate()
sc = spark.sparkContext
lines = sc.textFile('trump.txt')
counts = lines.flatMap(lambda line: line.split(" ")) \
        .map(lambda word: (word, 1)) \
        .sortByKey(ascending=True) \
        .reduceByKey(lambda a, b: a + b)
```
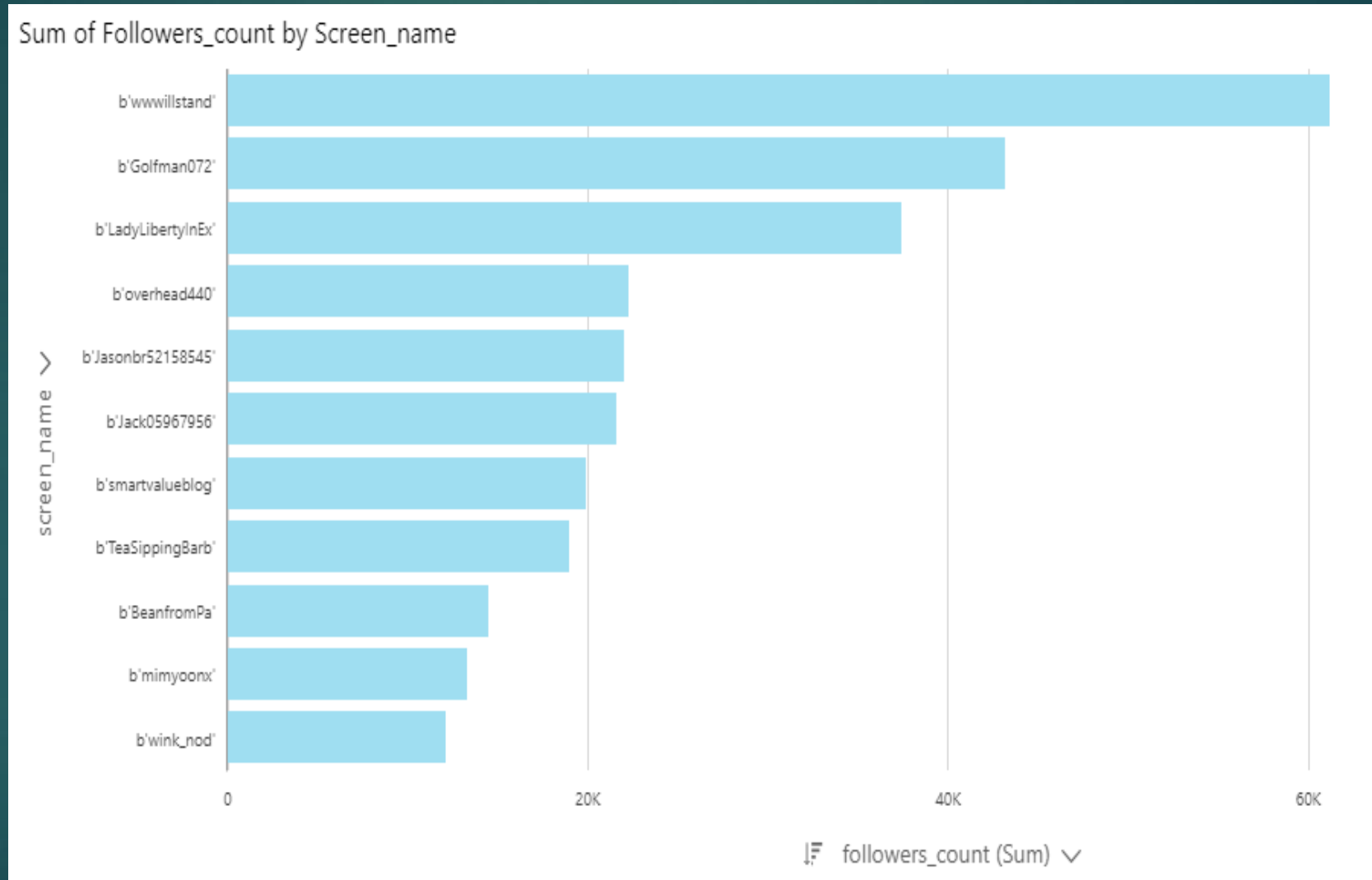
# Data Analysis & Viz

# Top Screen Names by Followers counts from Tweets



Sum of Followers_count by Screen_name

# Location, Timestamp, followers counts from Tweets

# WordClouds

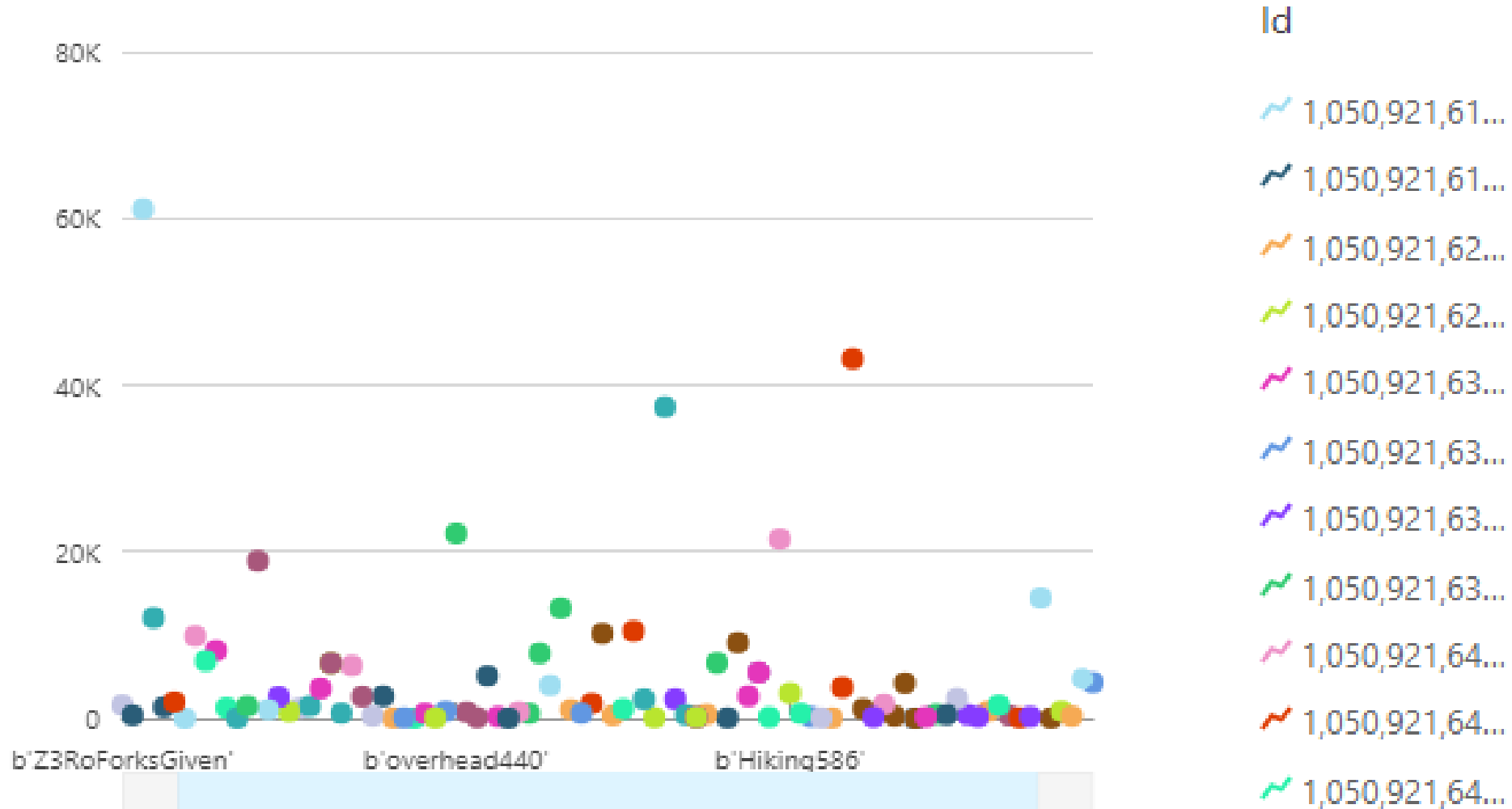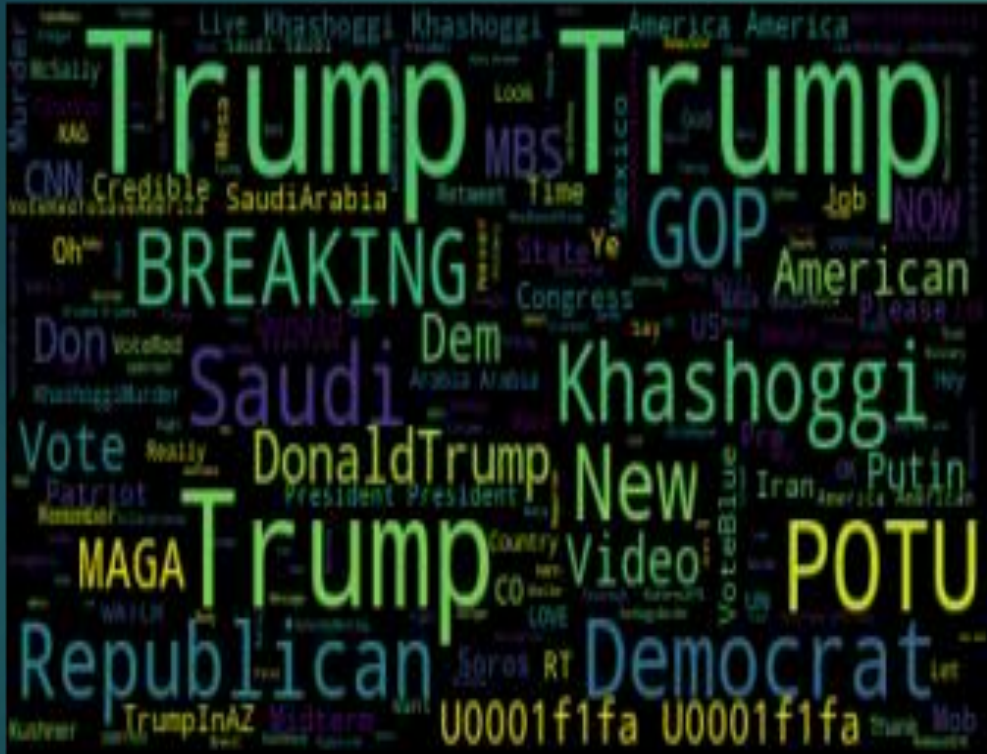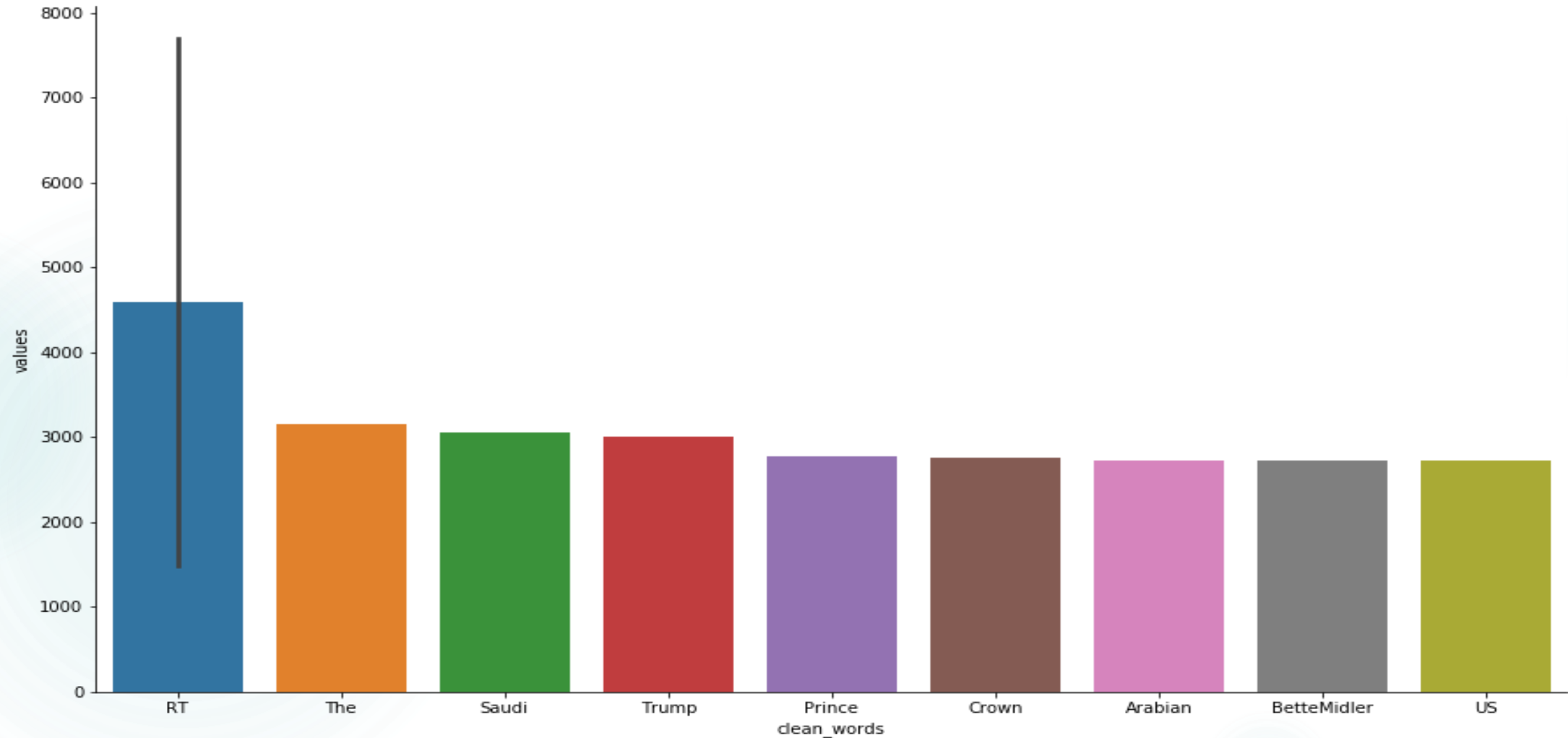# Top 10 Terms – PySpark Data Analysis – No ML

## Intertopic Distance Map (via multidimensional scaling)

PC2

PC1

Marginal topic distribtion

2%

## Top-30 Most Relevant Terms for Topic 1 (5.3% of tokens)

| | 0 | 20 | 40 | 60 | 80 | 100 | 120 |
|---|---|---|---|---|---|---|---|

trump
america
video
mbs
democrats
russia
there
world
well
watch
voteredtosaveamerica
red
russians
message
rally
pres
funny
was
up
water
israel
votered2018
speak
have
veterans
brilliant
complicitgop
bush
everything
trumpatwar

Overall term frequency

Estimated term frequency within the selected topic

# Top trending Topics from Trump tweets

| Topic 1 | Topic 2 | Topic 3 |
|---------|---------|---------|
| Trump | DonaldTrump | mega |
| Amercia | GOP | what |
| Video | American | saudiarabia |
| Mbs | Don | no |
| Democrats | who | arizona |
| Russia | tp | do |

# Sentiment Analysis results

- Positive Sentiments 66.5%
- Negative Sentiments 7.0%
- Neutral Sentiments 26.5

# Thank You