# EECS 4415 - Big Data Systems Project Presentation

# Computing TF-IDF Vectors for Subreddits
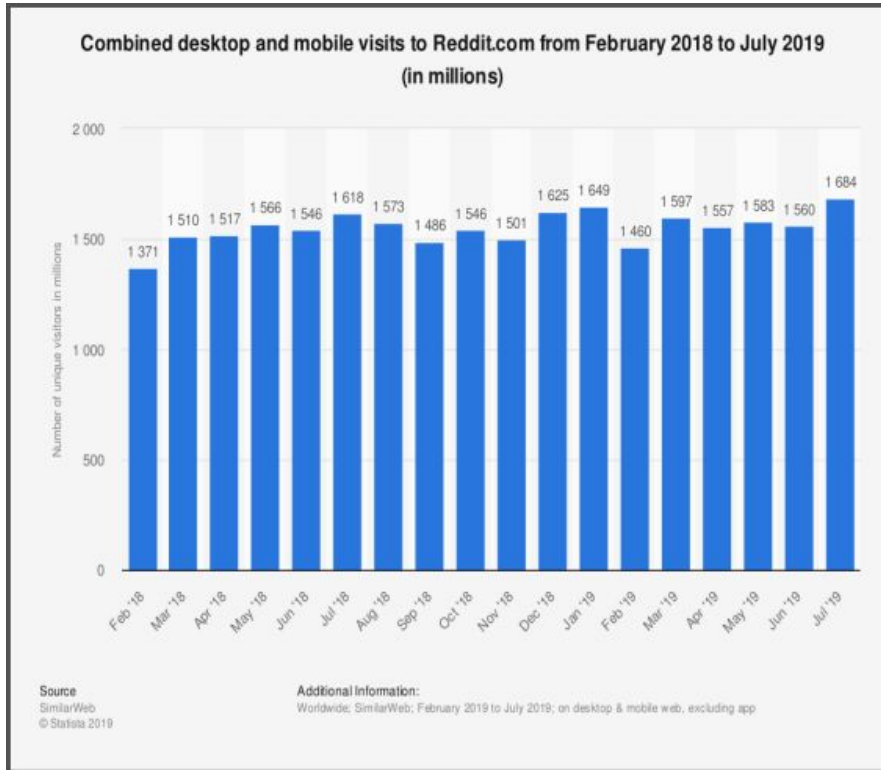
November 27, 2019.

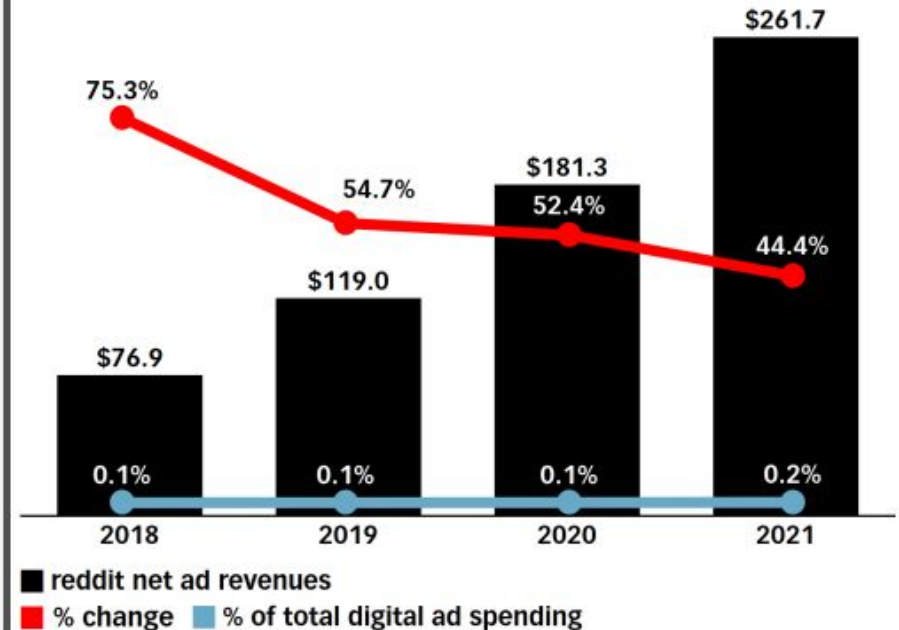Ken Tjhia <hexken@my.yorku.ca>
Qijin Xu <jackxu@my.yorku.ca>
Ibrahim Suedan <isuedan@hotmail.com>

# Motivation



Combined desktop and mobile visits to Reddit.com from February 2018 to July 2019 (in millions)



**Reddit Net Ad Revenues in the US, 2018-2021**
*millions, % change and % of total digital ad spending*

- Almost 1.7 Billion unique visitors and over 145 Million comments in July 2019!
- $261 Million projected spending on digital ads in 2021.

Note: includes advertising that appears on desktop and laptop computers as well as mobile phones, tablets and other internet-connected devices, and includes all the various formats of advertising on those platforms; net ad revenues after companies pay traffic acquisition costs (TAC) to partner sites; excludes nonadvertising revenues (e.g., Reddit Premium, Reddit Coins)
Source: eMarketer, February 2019

T10128                                                          www.eMarketer.com

YORK U
UNIVERSITÉ
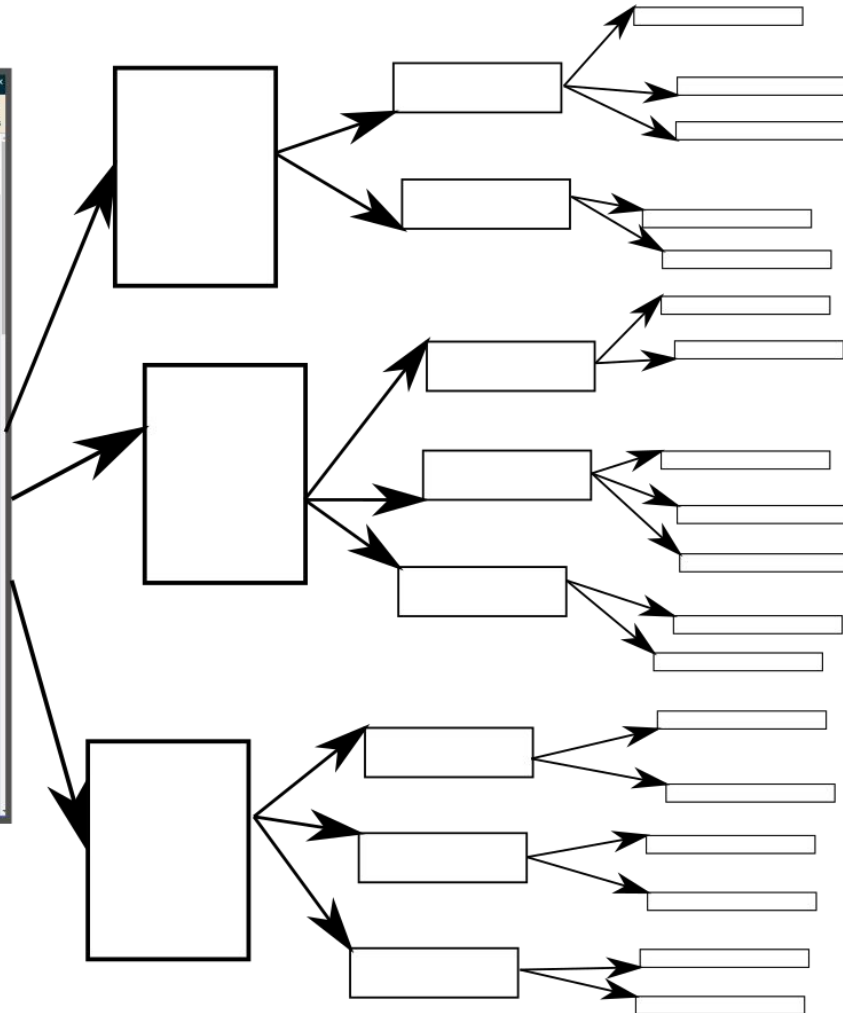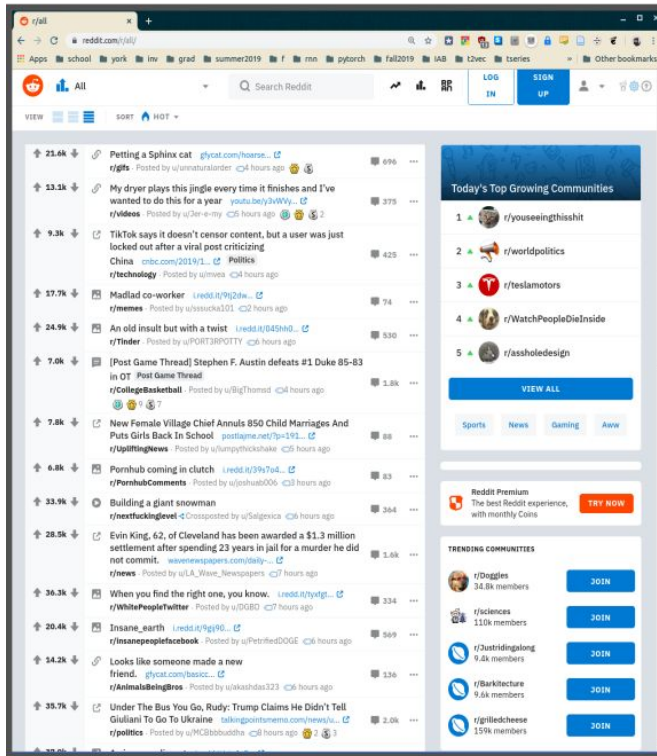UNIVERSITY

# Organization of Reddit
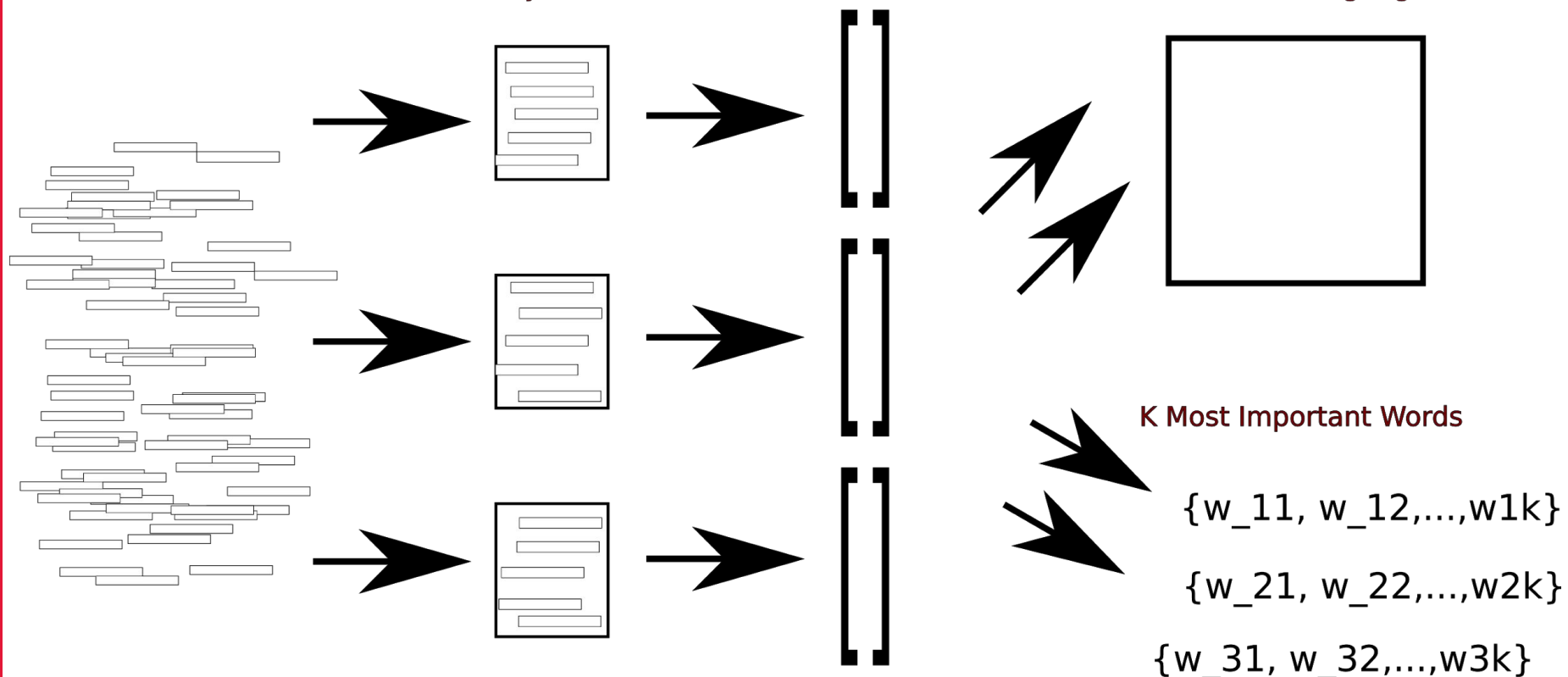
Reddit         Subreddits         Posts         Comments

# Data and Processing Dimensions



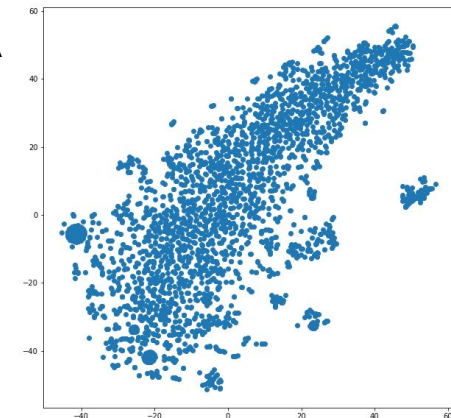Mixed Comments → Comments by Subreddit → Subreddit TF-IDF Vectors → Machine Learning Algorithms

K Most Important Words

$\{w\_11, w\_12, ..., w1k\}$

$\{w\_21, w\_22, ..., w2k\}$

$\{w\_31, w\_32, ..., w3k\}$

# Architecture Overview



Processing

**SPARK**

Exctract Fields
Aggregate subreddit comments
compute TF-IDF
compute top k words
ML algorithms on TF-IDF vectors

**Data Source**

JSON objects

**HDFS**

Store batch data until processing
Store TF-IDF vectors,
top k words

**Tableau**

Visualize
Clusterings

Display top k words

Source          Injestion          Storage          Serving and Visualization

# Results

No PCA



## Example of Words with Top 7 TF-IDF Scores

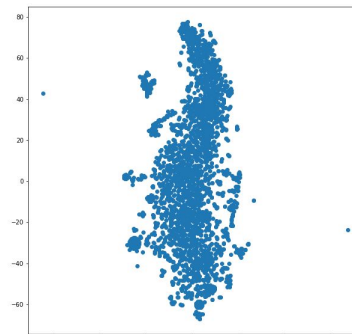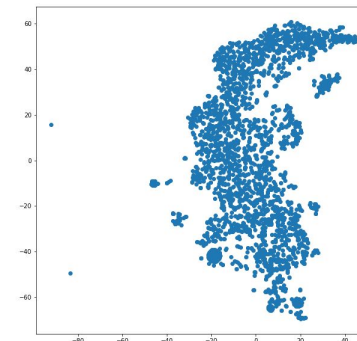| | subreddit | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|---|
| 13 | Botchedsurgeries | fillers | helium | stats | doctor | ass | literally | haters |
| 22 | scifi | stargate | universe | dark | series | tv | matter | extension |
| 2823 | Makeup | freckles | sheer | hide | natural | gorgeous | slightly | skin |
| 2829 | ethereum | matching | multiple | exchange | engine | cpu | services | spread |
| 2837 | spacex | rocket | grounded | fh | satellite | reasonably | risky | alternative |

PCA(10)

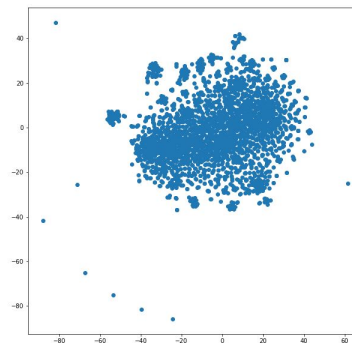

PCA(20)



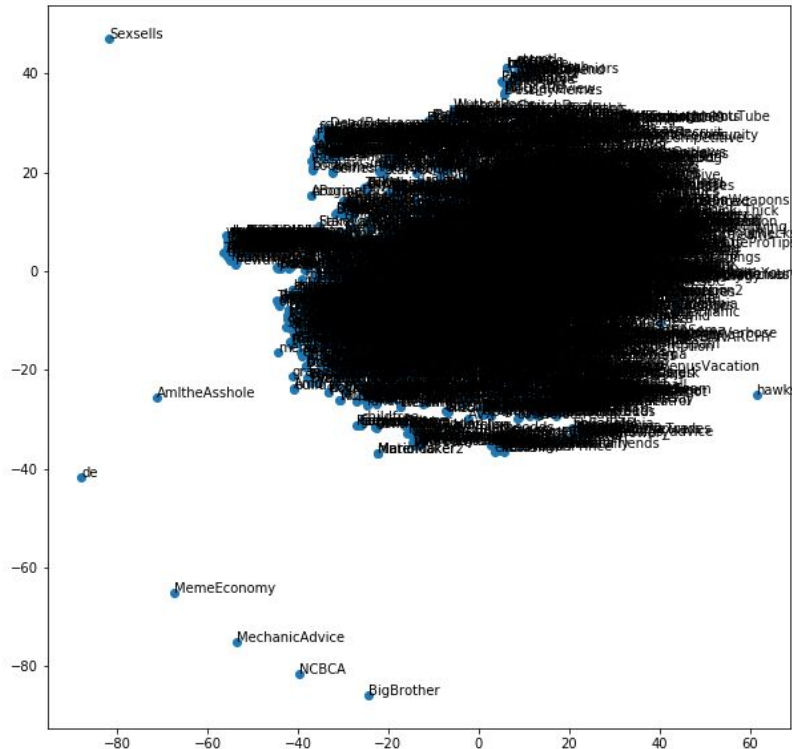PCA(30)



PCA(40)



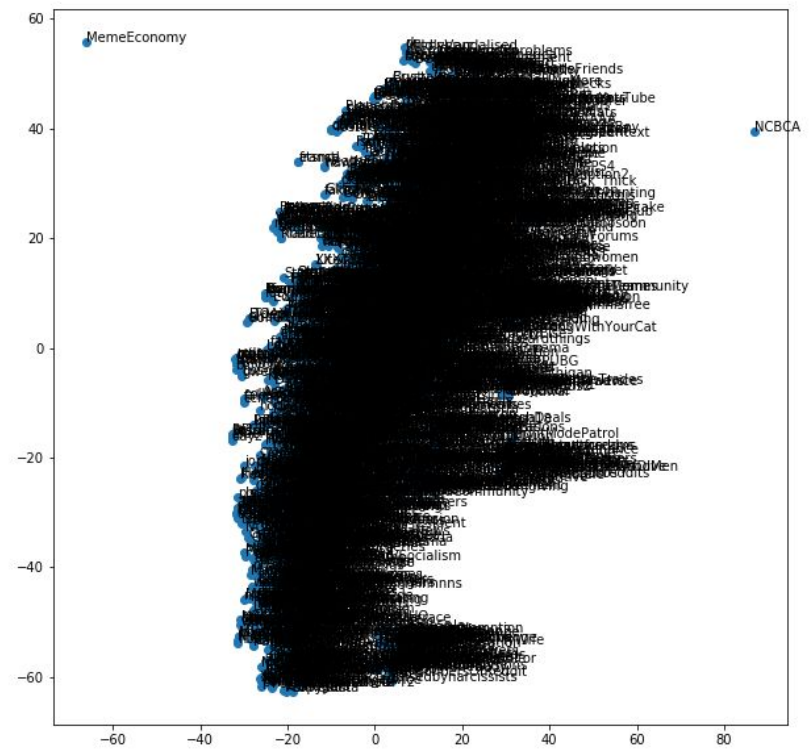PCA(100)



YORK U
UNIVERSITÉ
UNIVERSITY

# Results

PCA(100)

PCA(20)



- Manually choose subreddits we expect to be 'similar' and 'dissimilar', then compare their TF-IDF vectors.
- Cluster the TF-IDF vectors

YORK U
UNIVERSITÉ
UNIVERSITY

# Limitations

1. Currently cannot update TF-IDF vectors incrementally
2. Don't have an actual cluster..
3. Current design requires two passes over the data

# Scalability

1. Horizontally scalable: Uses Spark and HDFS.
2. Fault tolerant, if running on a cluster.
3. Can handle more data.

# Variations and Extensions

1. Add streaming Component, update TF-IDF vectors in real time
2. Consider other document sets, for example

   a document = all of a users text concatenated
1. Use the TF-IDF vectors as features for downstream tasks
2. Compute different kinds of TF-IDF scores, or learn different features in general

# Conclusion

**Mixed Comments**   **Comments by Subreddit**   **Subreddit TF-IDF Vectors**   **Machine Learning Algorithms**

**K Most Important Words**

$\{w\_11, w\_12,...,w1k\}$

$\{w\_21, w\_22,...,w2k\}$

$\{w\_31, w\_32,...,w3k\}$

Processing

**SPARK**

Exctract Fields
Aggregate subreddit comments
compute TF-IDF
compute top k words
ML algorithms on TF-IDF vectors

**HDFS**

Store batch data until processing
Store TF-IDF vectors,
top k words

**Tableau**

Visualize
Clusterings

Display top k words

**Data Source**

JSON objects

Source          Injestion          Storage          Serving and Visualization

YORK U
UNIVERSITÉ
UNIVERSITY

# QUESTIONS?