# Social Media Data Analysis

**Advisor : Amarnath Gupta**

**Amir Shirkhani, Mohammed Tawashi, Hanu Pathuri, Naveed Mohammed, Vamsi Namuduri**

## Introduction

The project focuses on studying how Twitter images impact the narrative of hashtags.

**Hypothesis:** Media files (Images/videos) are being used to influence the narratives of hashtags.

**Objective:** The project aims at finding specific patterns in tweets where media files (images) are used to change the narrative of the corresponding hashtag and co-occurred hashtags. To study the importance of images in the Tweets by modeling aggregated text of hashtags with and without images.
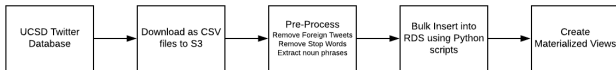
**Targeted Audience:**

❖ Election organizer as part of a campaign can study and detect endorsing and opposing trends and act by countermeasures using similar techniques

❖ Social Media platform and specially Twitter itself can detect patterns and potentially restrict the behavior.

❖ Journalists can report to general public on how a potential small group of influencers can sway a narrative and push various agendas.

## Data Sources

UCSD Twitter database is the main data source for the project. Below are the main entities that are downloaded from this database.

➢ Users
➢ Tweets
➢ Hashtags
➢ Media URLs

The data collected is processed and inserted into RDS Postgres instance. Materialized Views are created on top of the data.
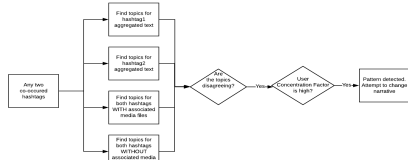


## Methodology

**Hashtag Co-occurrence**: This is pair of hashtags that were used together in a Tweet.

**User Concentration Factor**: Is the number of tweets per User.

**Pattern Recognition:**



**Community recognition based on user mentions:**

NFM Model with NetworkX sub-graph traversal parameter K-Core = [2, 6]
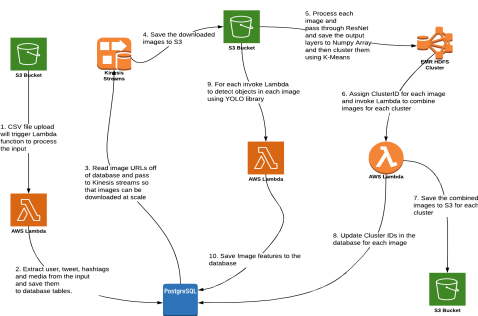(K is the in-degree + out-degree)

No of Topics = 10

No of words per topic = 20.

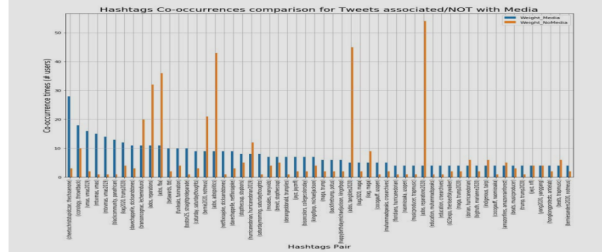At K core = 6, the model converges and provides optimal communities based on the topics.

## Data Pipeline

### Solution Architecture
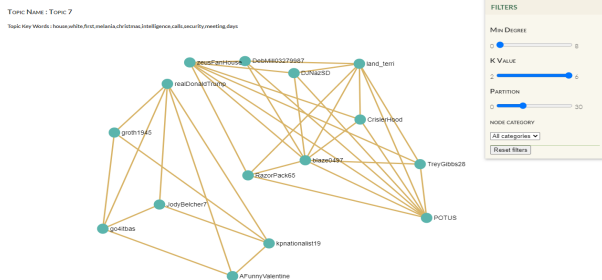


## Results

**Hashtag Co-occurrence with and without media**







## Conclusion

Conclusions that came out of the data analysis from Twitter.

➢ Text extraction from media files associated with tweets are, in most cases, supplementing the narrative of the tweet's text.

➢ Media files are being used as a powerful tool to contaminate the original narrative of single/co-occurred hashtag pairs.

➢ Object detection from images did not result in any additional insights.

## Future Work and Acknowledgements

**Future Work:**

The model we developed can be extended on a few different ideas:

o Tweets related to Trump dominated the dataset. We can do the analysis by filtering those tweets.

o Include videos and GIFs/memes in the analysis.

o Expand the dataset to include wide range of topics instead of just political tweets.

**Acknowledgements:**

We would like to thank:

✓ Our advisor Prof. Amarnath Gupta

✓ All the DSE professors and teaching assistants.

✓ Staff at Supercomputer Center (SDSC).