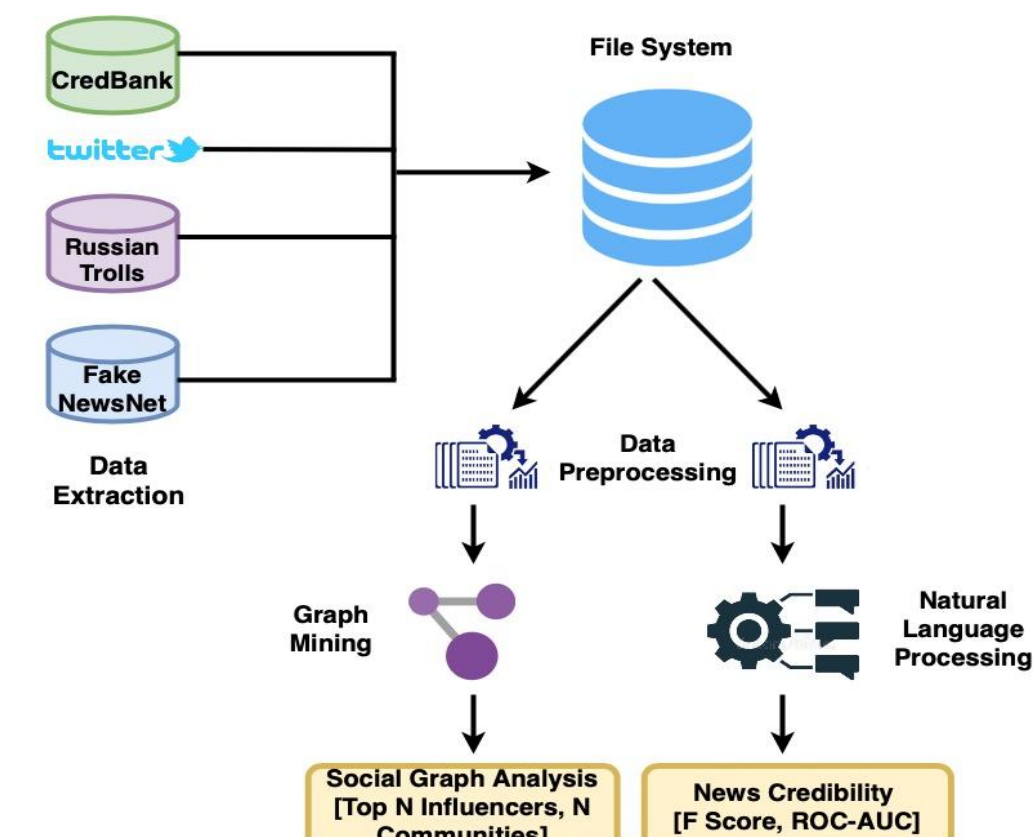




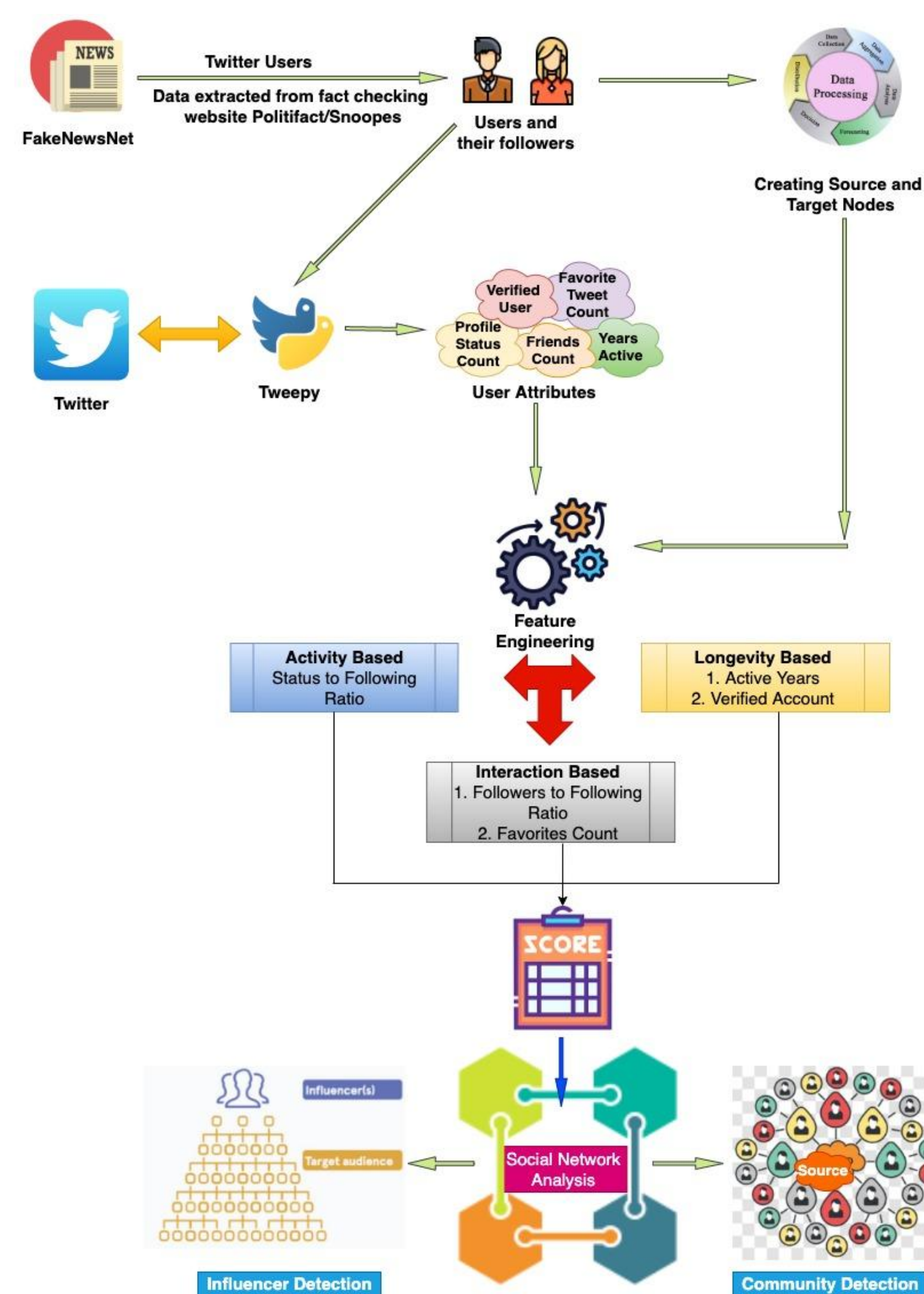
## Introduction

The advancement in technology and social media platforms has eased the creation and publication of news and so the power to spread misinformation has come in the hands of normal people. The news outlet and social media platforms may publish fake news to build a readership, or even more than that as a part of psychological warfare. Although fake or misleading information has been around us since long, after the 2016 US presidential elections, this topic grabbed more attention. There are several fact-checking online resources, like wikipedia, snopes.com, fact-check.org and few others, available and accessible to every individual. Despite this convenience being provided, it is a monotonous and tiresome activity to manually verify the sources for the originality of the news forecasted. Our research aims to design a system that can detect tweets with fake news and prevent fake news from spreading further.



## Methodology

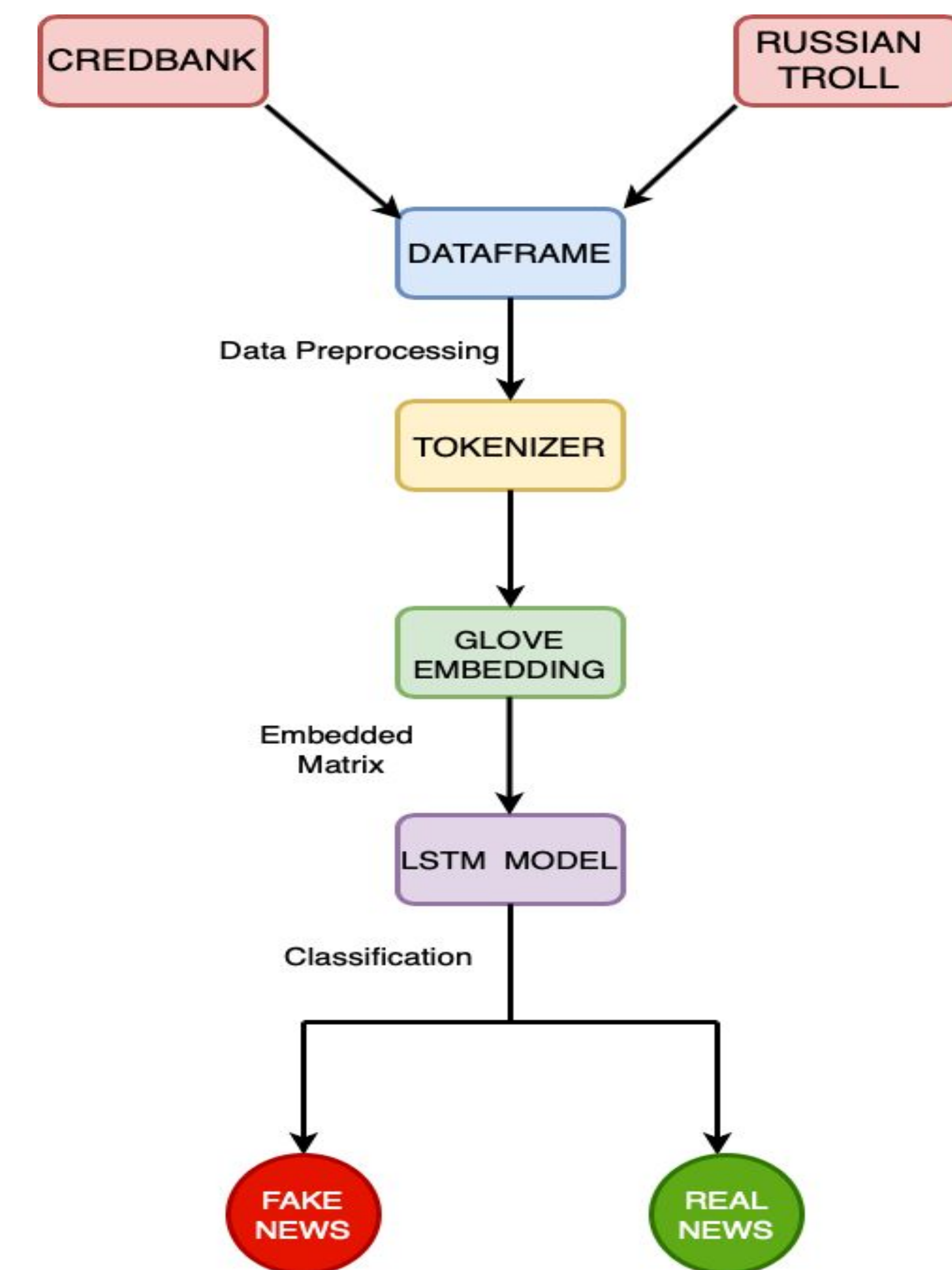
## Fake News Propagation



## Methodology

The user and their followers are collected from FakeNewsNet Dataset and their corresponding user attributes (number of followers and following, active years, verified status, tweet count, favorites count) are extracted from twitter with the help of Tweepy API. The edge weights are computed based on the score assigned to these user attributes on the basis of their quantile range. Due to the computational constraints in Neo4j, the influencer and community detection is done on a reduced dataset where the top 25% of the followers are considered based on their user attribute scoring. The influencer detection is performed with the help of Closeness centrality, Betweenness Centrality, Degree Centrality and Page Rank due to their ability to measure the influence a node has in a social graph. The Community detection is performed using Louvain, Label Propagation and Strongly connected component algorithms due to their ability to quickly and efficiently identify subcommunities in a graph.

## Content based



A balanced dataset is created by using labelled data from Credbank and Russian trolls data sets. These data sets are combined into a dataframe and is further preprocessed. After pre-processing, tokenization is done to create vector sequences. We have used GloVe embedding for creating the contextually related vectors. Once we have the embedding matrix, we feed this to our LSTM model that has five layers (Embedding, LSTM, GlobalMaxpool and 2 dense layers with 50 and 2 neurons respectively). We trained the LSTM model with various combinations of batch sizes and number of epochs and finally found that a batch size of 128 trained on 10 epochs works the best.

## Analysis and Results

Fake News Propagation  
Influencer Detection

*Fake News:*

Node	Degree Centrality	Page Rank	Betweenness Centrality	Closeness Centrality
19672966	412	167.8201207	167.8201207	0.2251199232
21237884	378	165.8195036	165.8195036	0.1836982333
20974554	386	158.559881	158.559881	0.2232423176

*Real News:*

Node	Degree Centrality	Page Rank	Betweenness Centrality	Closeness Centrality
57419364	343	142.0523125	7718611.921	0.2079099437
16255515	235	99.98408311	7562899.863	0.2035352719
1426071613	224	99.45528214	24976	1

## Community Detection

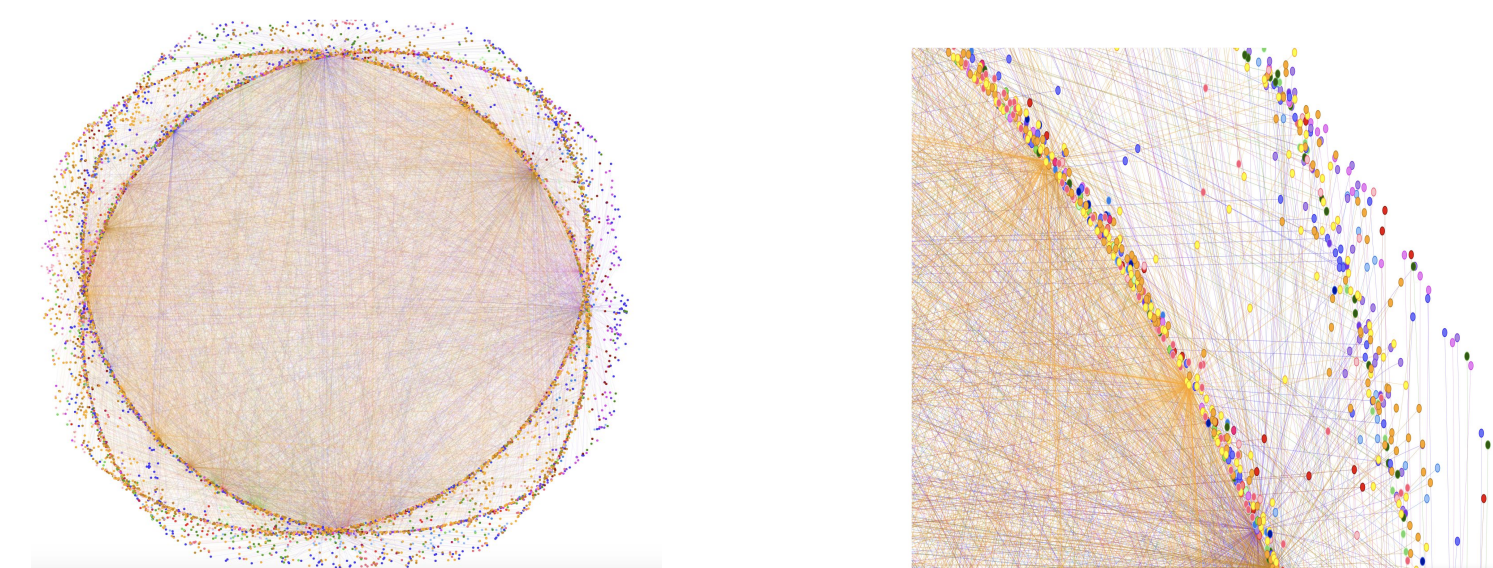
*Comparative analysis of community detection algorithms:*

Community Algorithm	Type of Users	Total Number of Communities	Size of the largest Community	Size of the smallest community	Number of overlapping communities formed
Label Propagation	Fake	47	4031	7	N/A
	Real	41	8805	2	N/A
Louvain	Fake	122	373	7	40
Modularity	Real	127	303	2	37
Connected Components	Fake	24	14080	7	N/A
	Real	26	13853	2	N/A

*Significance of Influencers in Community Detection:*

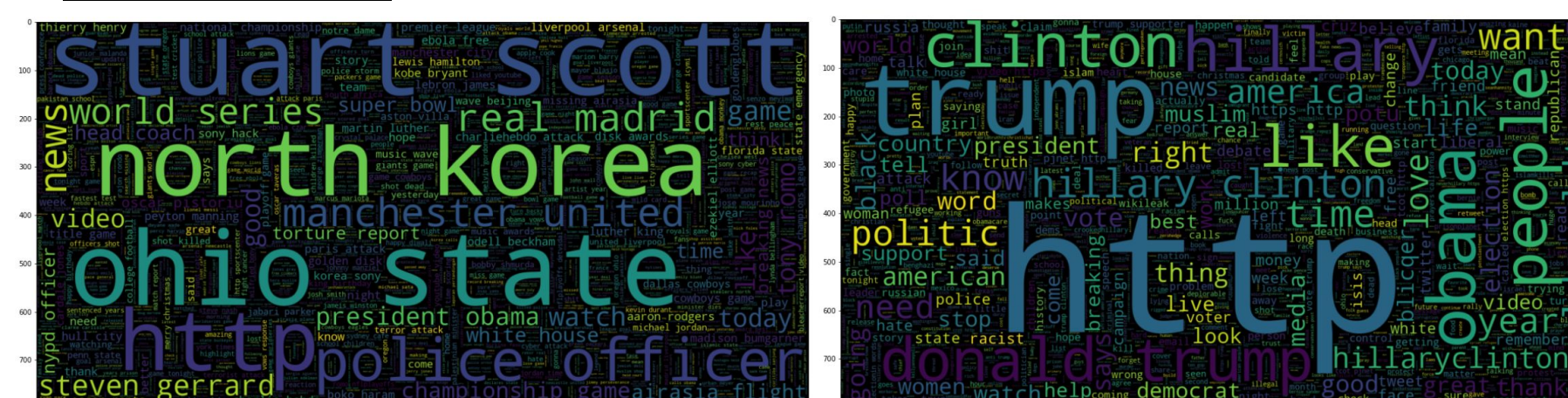
Community Algorithm	Type of Users	Top 3 Influencers (user id)	Community label in which Influencer node is present	Size of the community	Community Rank based on size (Nth largest community)
Label Propagation	Fake	19672966	21104	4031	1
	Real	21237884	18748	2177	3
	Real	20974554	21104	4031	1
	Real	57419364	1	8805	1
Louvain	Fake	1426071613	150	1610	2
	Real	16255515	14889	225	8
	Real	19672966	6509	350	4
	Real	21237884	9440	373	1
	Real	20974554	3865	338	6
	Real	57419364	8420	303	1
Connected Components	Fake	1426071613	14889	225	2
	Real	16255515	2592	218	13
	Real	19672966	0	14080	1
	Real	21237884	0	14080	1
	Real	20974554	0	14080	1
	Real	57419364	0	13853	1
	Real	1426071613	14888	255	3
	Real	16255515	0	13853	1

*Visualization for Community Detection:*



## Twitter content analysis

*Word clouds:*

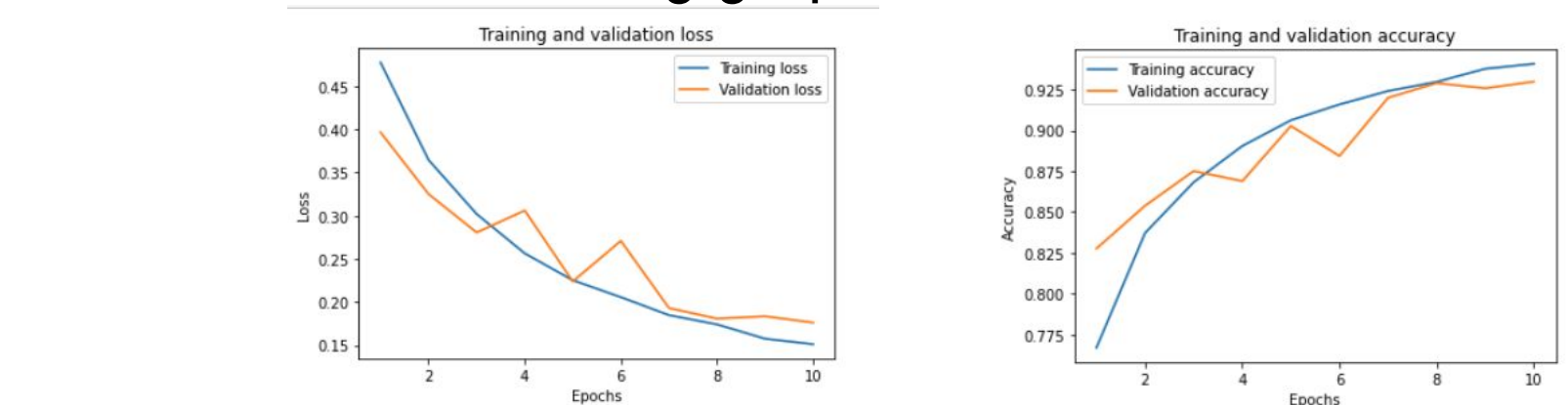


As can be seen from the above two word clouds, usage of words are different in the two types - Real and Fake. While the real tweets contain very less names, the fake news contain a lot of names of people. In addition fake news contain a big usage of catchy and emotion generating words such as "racist" and "breaking".

Our LSTM model was able to capture the distinguishing features between real and fake news tweets with an F1 score of ~0.91. The confusion matrix of our model looks as below. Since it is more important to capture fake news and tagging real news as fake is not that harmful (high false positives), our model works well.

ACTUAL \ PREDICTED	Real	Fake
Real	TP 7033	FP 1068
Fake	FN 292	TN 7908

To optimize the batch size and number of epochs, we experimented with various combinations and with a batch size of 128 and 10 epochs obtained the following validation and training graphs:



## Conclusion and Future Work

Thus, the Social Network Analysis was performed and the significance of influencer and community detection were understood. The future work for fake news propagation analysis would be to use weighted centrality measures for influencer detection and ensemble community detection algorithm which uses the result of all community detection algorithms analyzed.

The features of a fake news tweet were analysed and with the limited amount and variety of data available, we could build a model that could predict if a tweet was fake or real approximately 91% of the time. In future more balanced and larger dataset could be employed to avoid the bias that our dataset has. (Positively labelled data is from the 2016 elections).

Hybrid fake news detection engines that combines the content and user characteristics would be the future scope of this research.

## Key References

- [1] P. E. N. Lutu, "Using Twitter Mentions and a Graph Database to Analyse Social Network Centrality," 2019 6th International Conference on Soft Computing & Machine Intelligence (ISCMI), 2019.
- [2] A. M. Pudjajana, D. Manongga, A. Iriani, and H. D. Purnomo, "Identification of Influencers in Social Media using Social Network Analysis (SNA)," 2018 International Seminar on Research of Information Technology and Intelligent Systems (ISRITI), 2018}.
- [3] T. Li, M. Hua and X. Wu, "A Hybrid CNN-LSTM Model for Forecasting Particulate Matter (PM2.5)," in IEEE Access, vol. 8, pp. 26933-26940, 2020, doi: 10.1109/ACCESS.2020.2971348.

## Acknowledgements

We would like to thank Professor Magdalini Eirinaki and the Department of Computer Engineering for all the support, guidance and encouragement in completing the Master's Project.