



Twitter Community detection

Unveiling the Twittersphere

IFQIR MOHAMED
SEDJARI YASSINE

Supervised by:
Pr. Mohamed LAZAAR

Outline

- Introduction
- Communities in social media
- Dataset description
- Feature extraction
- Clustering phase
- Results & Experimentation
- Conclusion

Introduction





Communities in social media

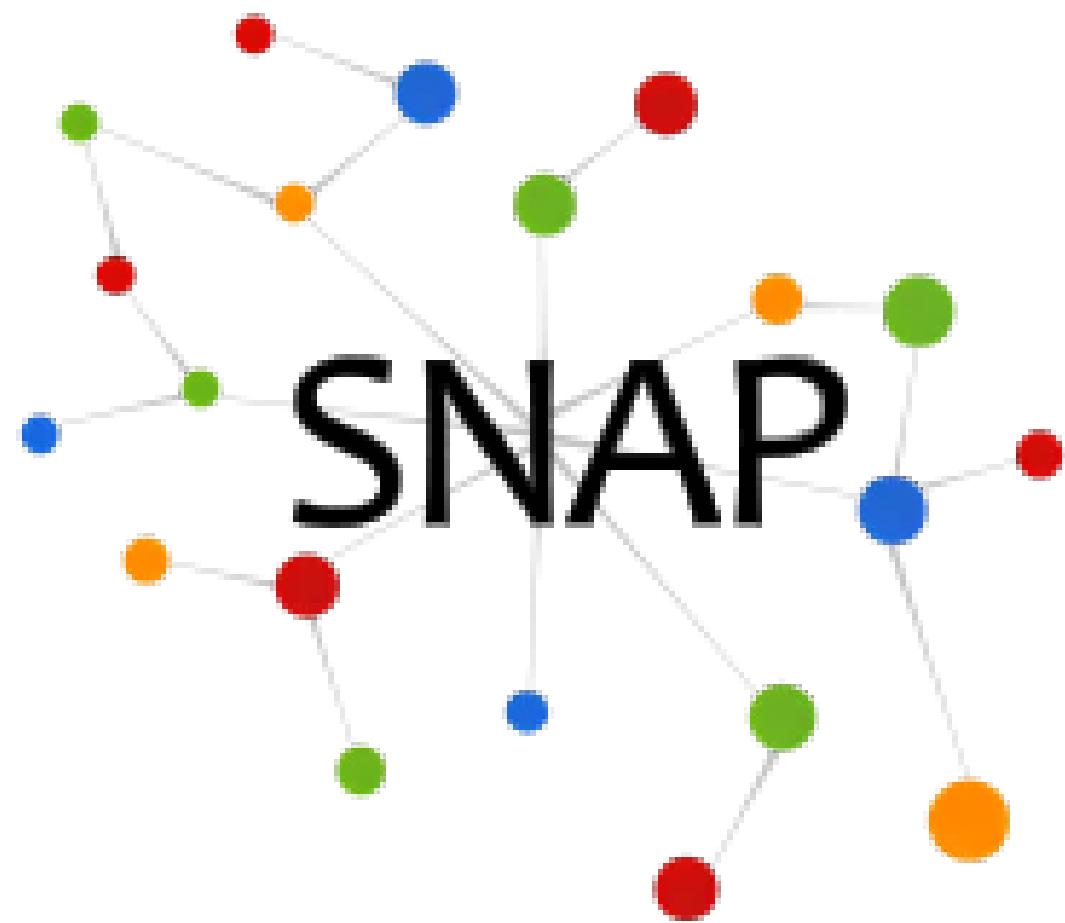
- The social hierarchy of social media users









Dataset Description





Dataset statistics	
Nodes	81306
Edges	1768149
Nodes in largest WCC	81306 (1.000)
Edges in largest WCC	1768149 (1.000)
Nodes in largest SCC	68413 (0.841)
Edges in largest SCC	1685163 (0.953)

Average clustering coefficient	0.5653
Number of triangles	13082506
Fraction of closed triangles	0.06415
Diameter (longest shortest path)	7
90-percentile effective diameter	4.5



Files hierarchy

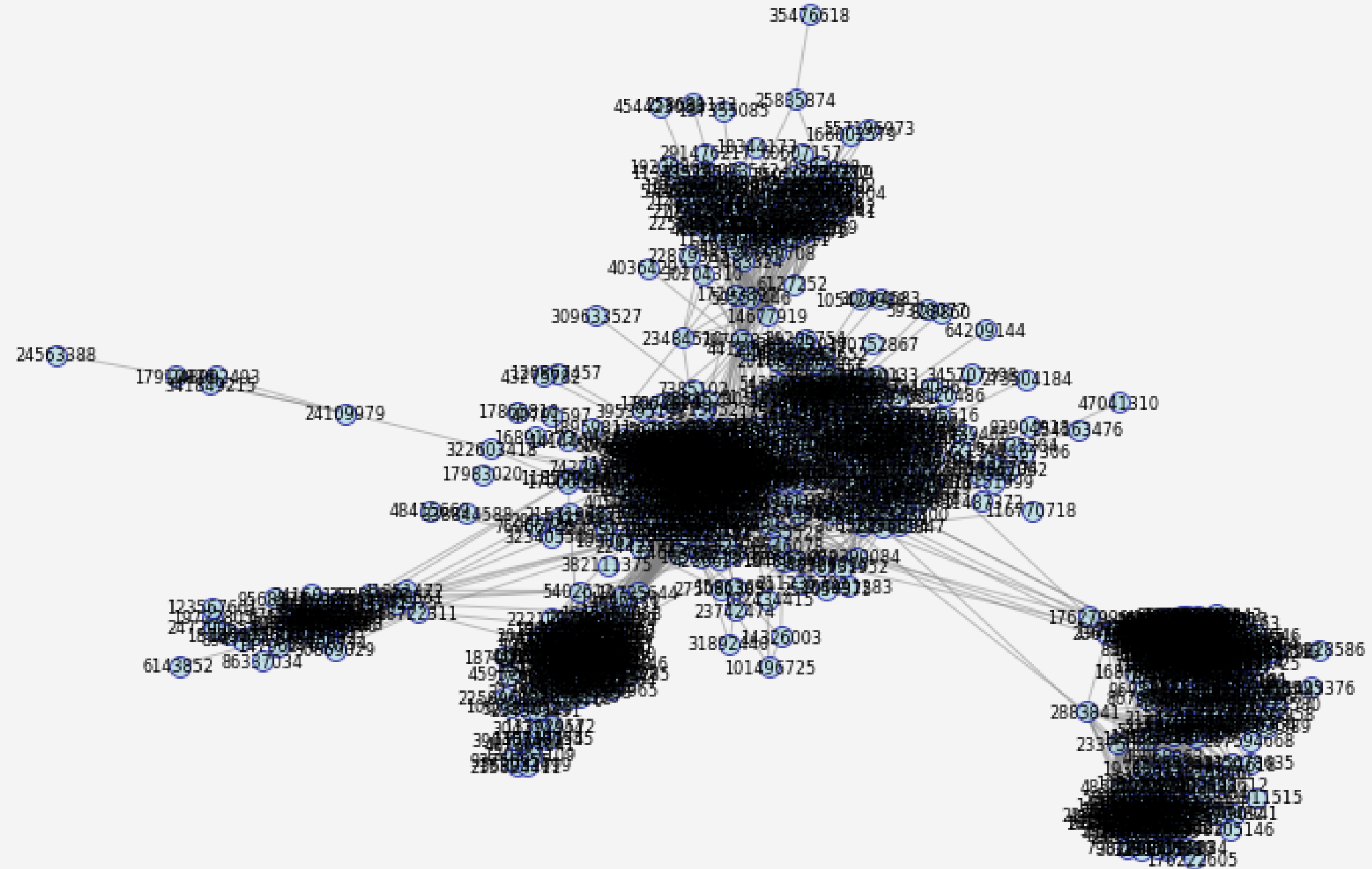
Project directory

- **twitter_combined.txt**
- **Twitter directory**

NodeID.egofeat

NodeID.featnames

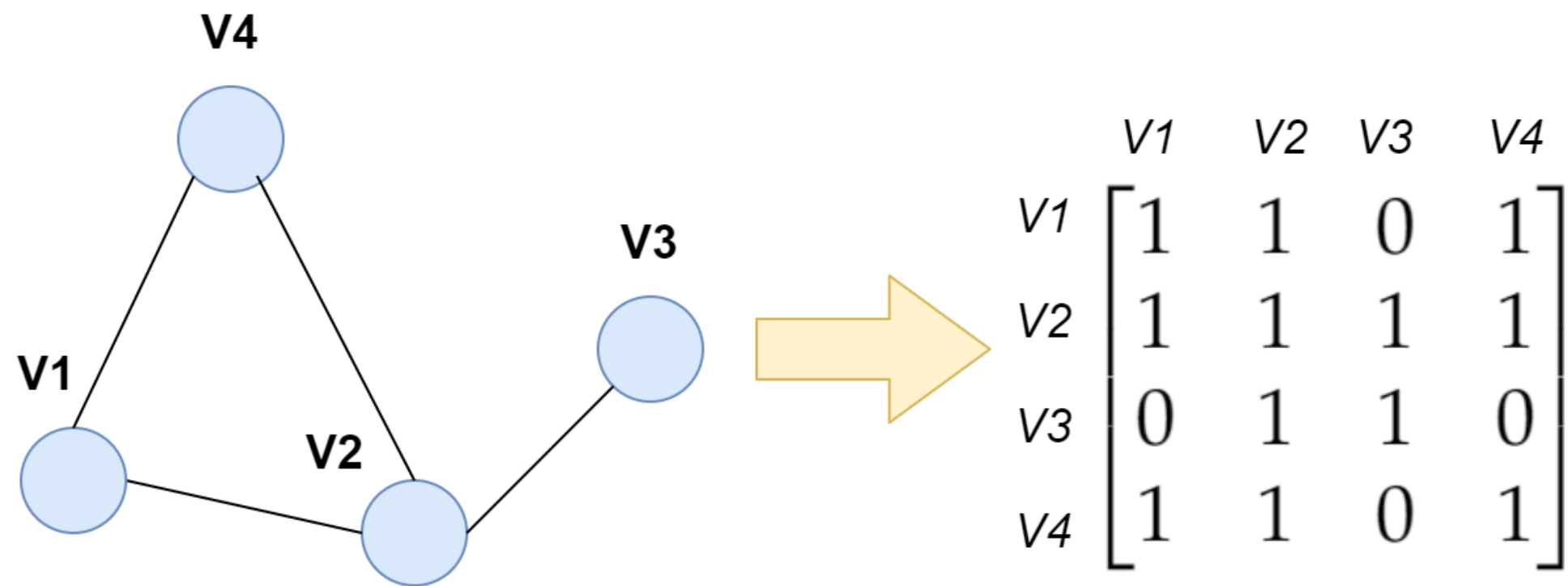
Graph network



Directed graph $A \rightarrow B$

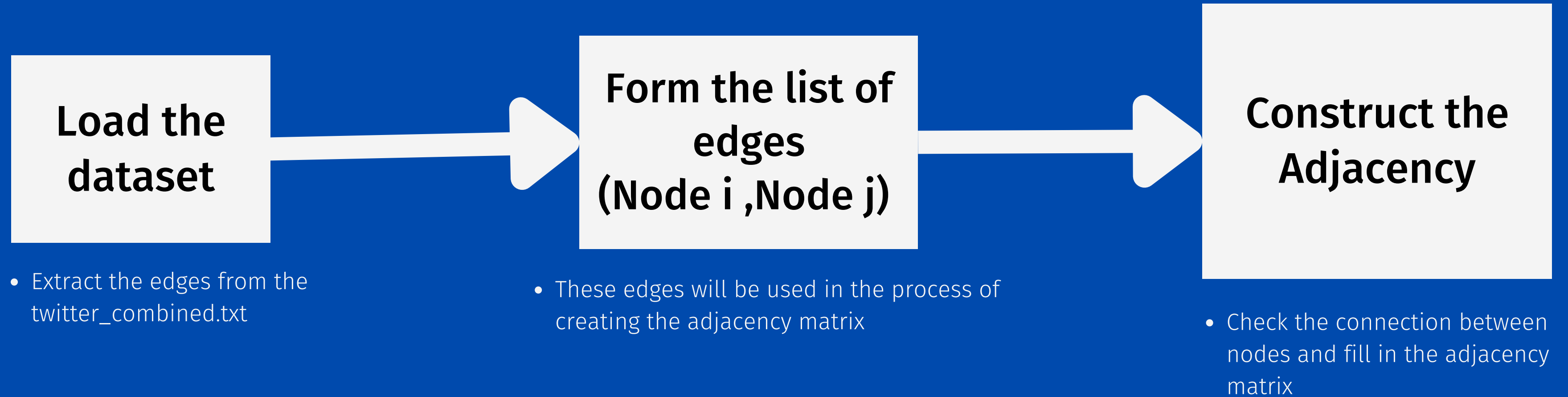
Node A follows Node B

Feature Extraction

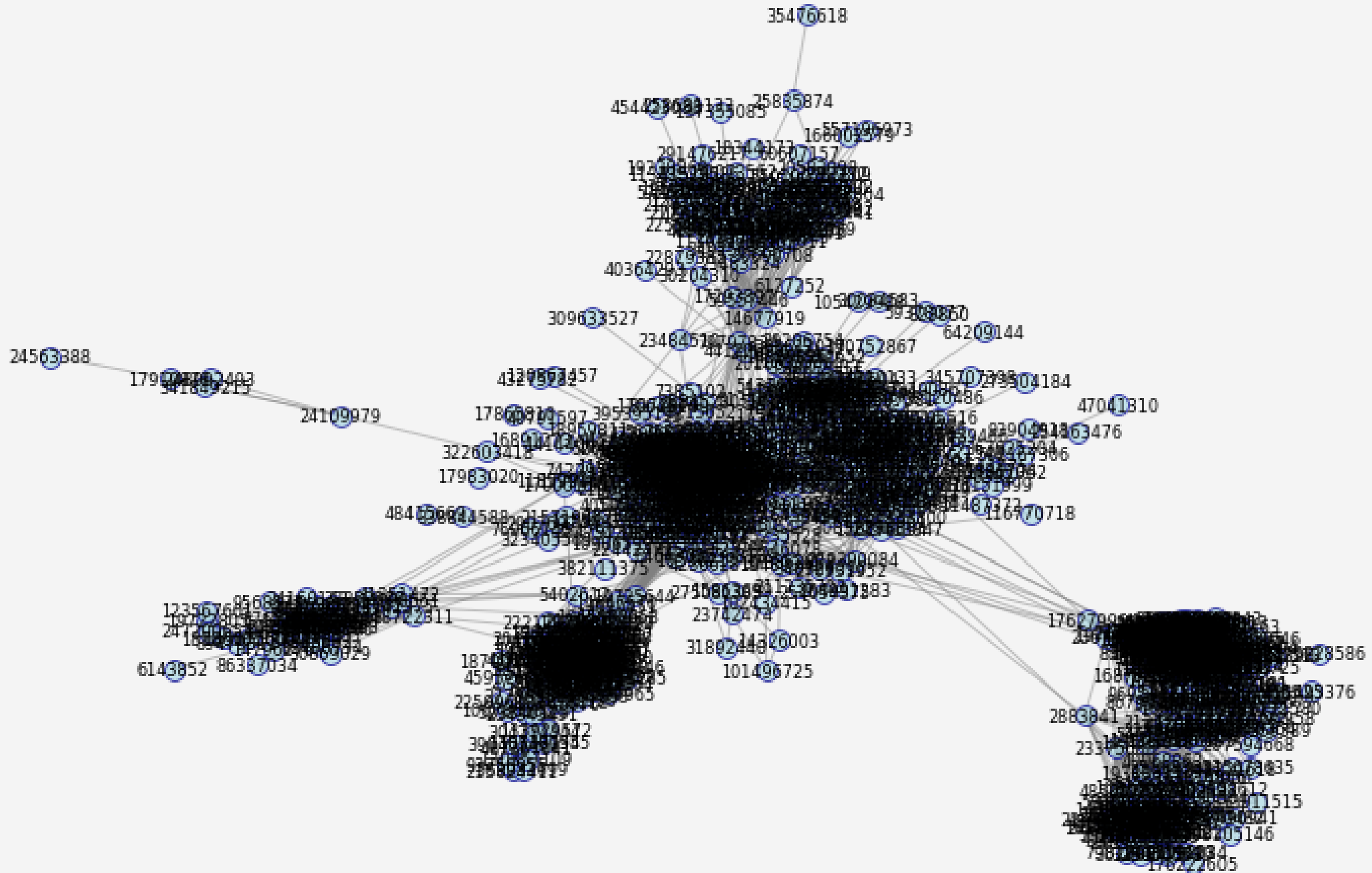


Edge-based Approach

Extract Adjacency Matrix



Graph network



Feature-based approach

Extract Adjacency Matrix

Load the dataset

- Extract the edges from the twitter_combined.txt

**Form the list of edges
(Node i ,Node j)**

- These edges will be used in the process of creating the adjacency matrix

**Extract account features from
nodeId.feature names**

- Check the connection between nodes and fill in the adjacency matrix

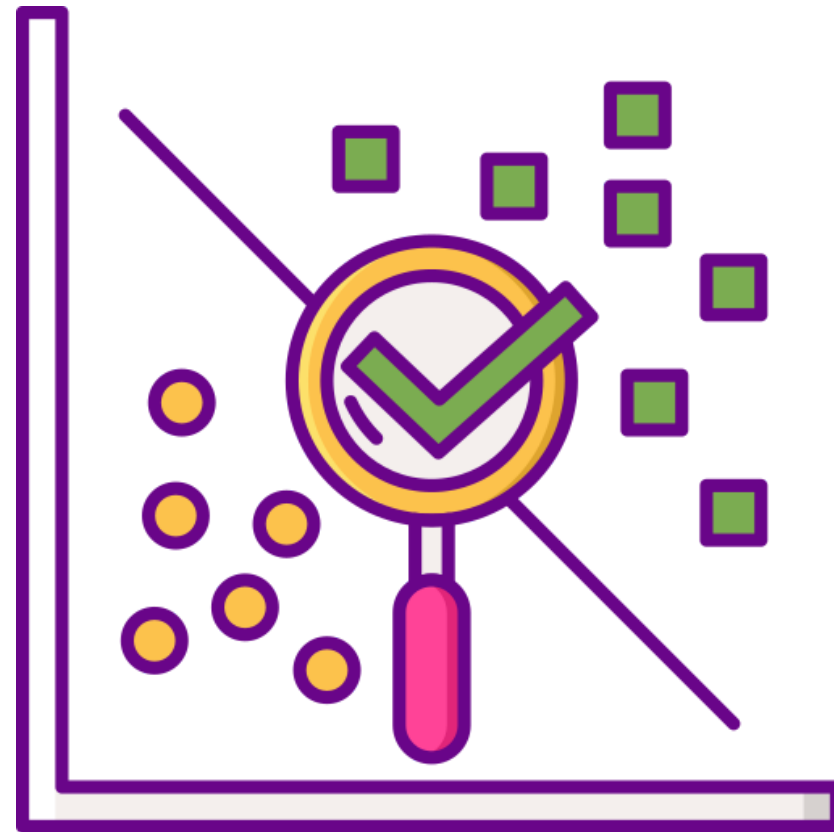
Compute degree, closeness and betweenness centrality

**Store extracted data in a Dataframe
(972,24250)**

Final Dataset

	nodeId	#OCTAVIA	#THEHELP	#ff	@BAFTA	@FuckYesEmma	@JUDAOcombr	@astowellcom	@emmastonebr	@helpmovie	...	@chococat
968	14528221	0	0	1	0	0	0	0	0	0	...	0
969	14840869	0	0	0	0	0	0	0	0	0	...	0
970	82726142	0	0	0	0	0	0	0	0	0	...	0
971	255790981	0	0	0	0	0	0	0	0	0	...	0
972	36618690	0	0	0	0	0	0	0	0	0	...	1
@helpmovie	...	@chococat	@dannyBstyle,	@jon_blaze55	@kylepulver:	@terrycavanagh	@twitchtv.	Degree Centrality	Closeness Centrality	Betweenness Centrality		
0	...	0	0	0	0	0	0	0.001030	0.238165	0.000000		
0	...	0	0	0	0	0	0	0.003090	0.239990	0.001701		
0	...	0	0	0	0	0	0	0.002060	0.206245	0.000958		
0	...	0	0	0	0	0	0	0.009269	0.273984	0.002988		
0	...	1	1	1	1	1	1	0.010299	0.294064	0.010228		

Clustering Phase



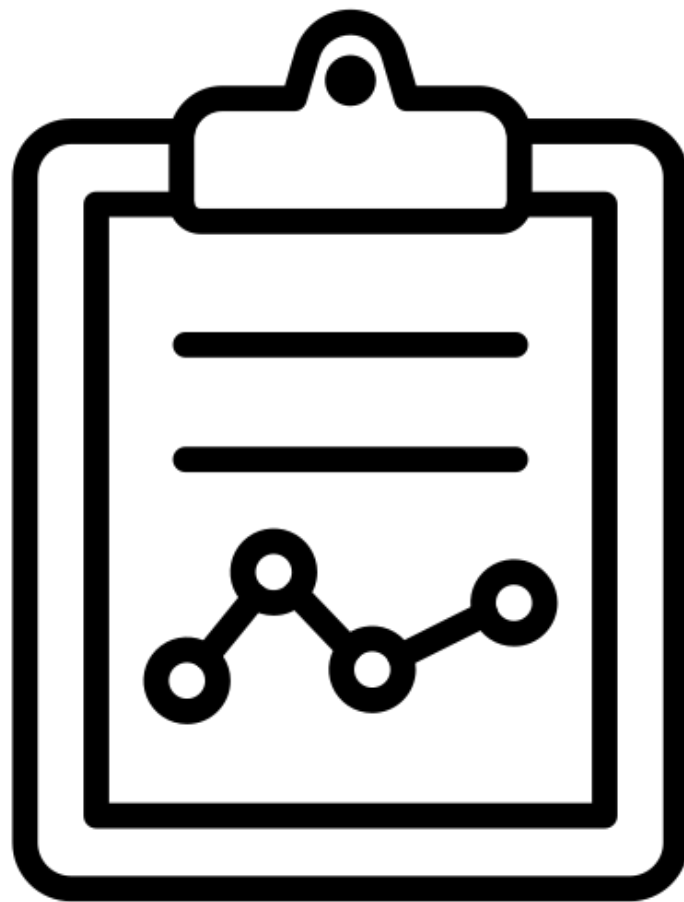
Edge-based Approach

Parameter	Value
Number of clusters	3
Affinity	Precomputed
Algorithm	Spectral Clustering

Feature-based approach

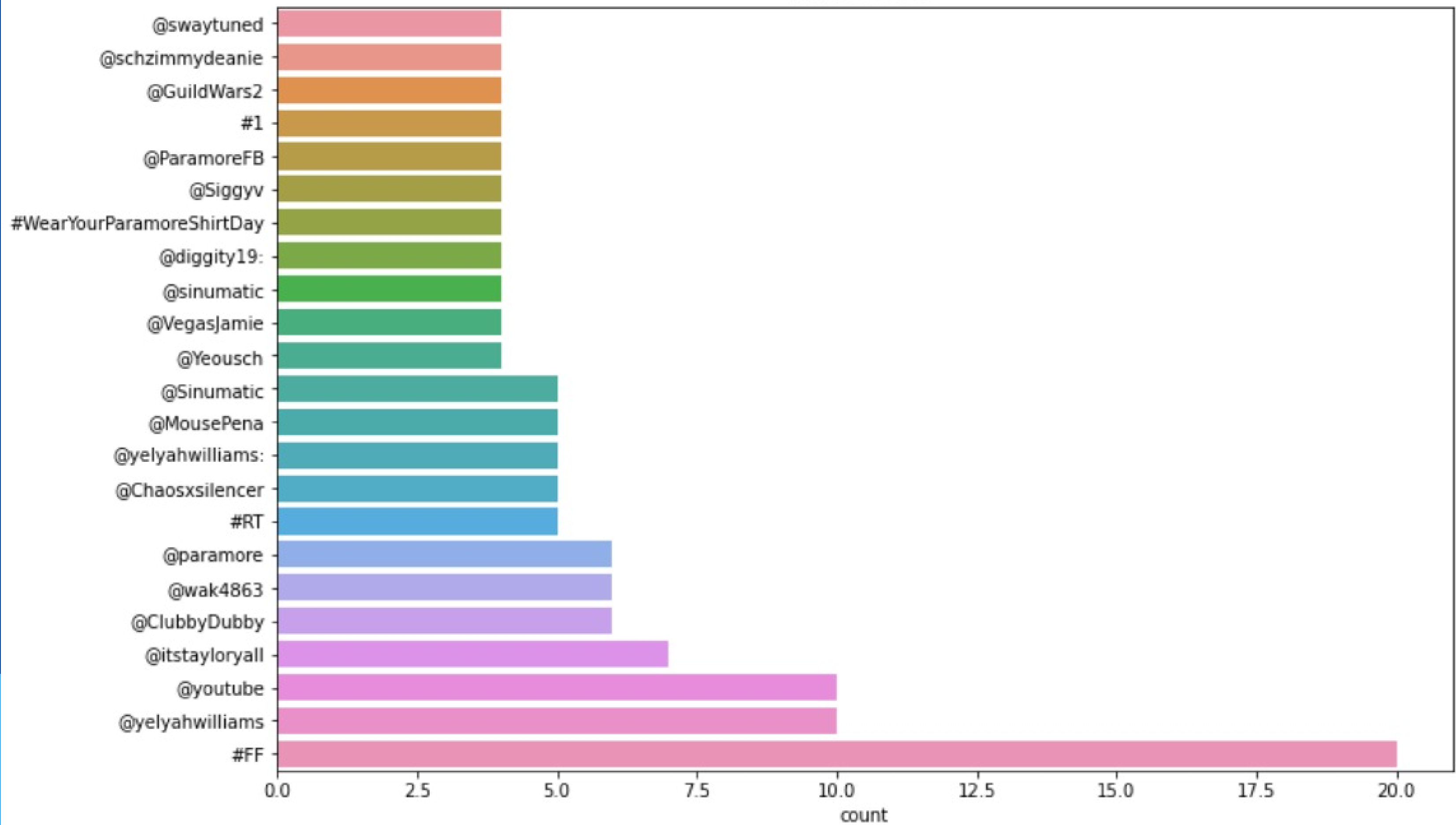
Algorithm	Number of Clusters	Data	Random Seed
k-means	3	df	42
hierarchical	3	df	-
spectral	3	df	42

Results & Experimentation

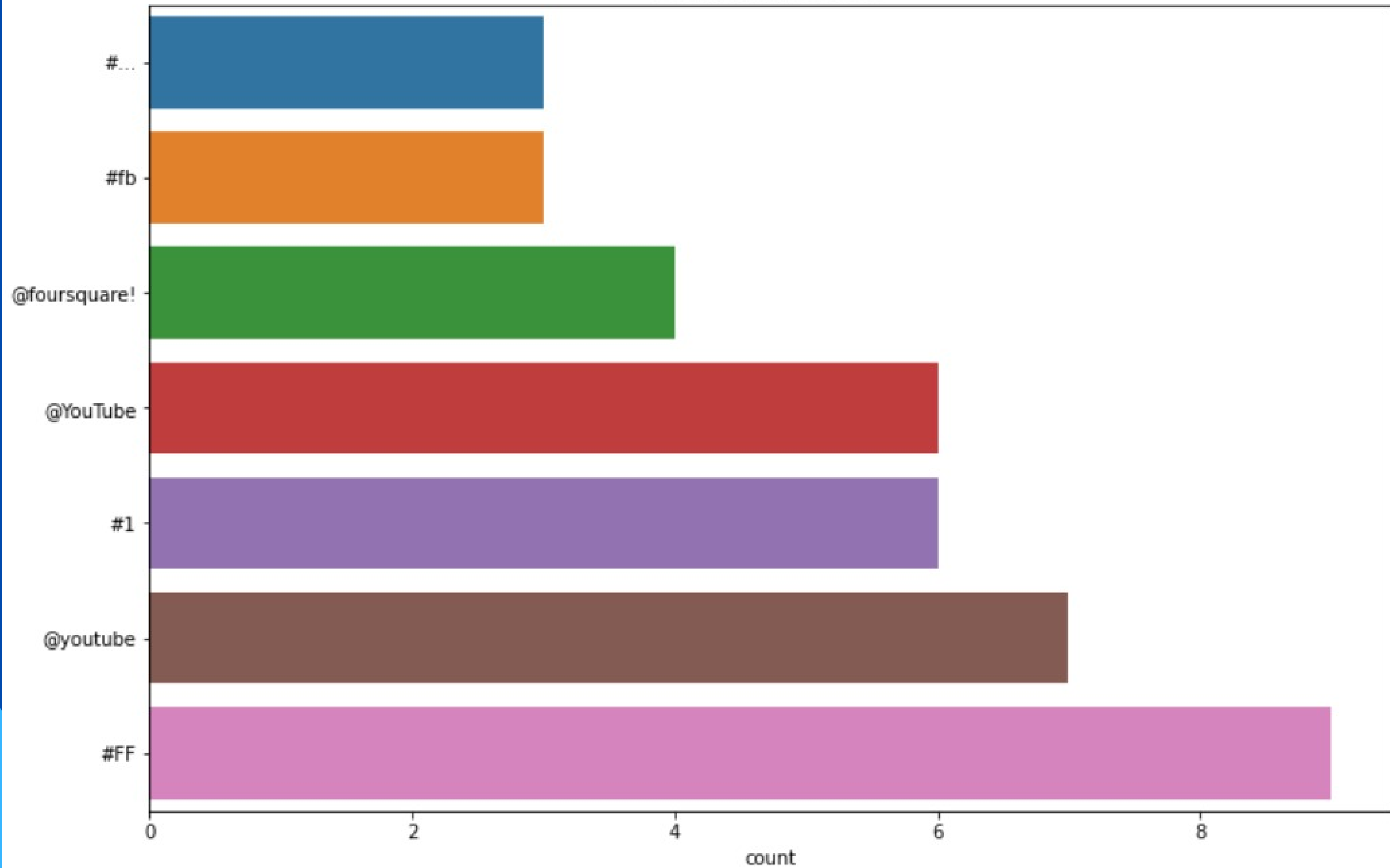


Method	Clustering Algorithm	Silhouette Score
Method 1: Edge-based Approach	Spectral Clustering	-0.1887
Method 2: Feature-based Approach	Spectral Clustering	-0.0343
	Hierarchical Clustering	0.6756
	K-MEANS	0.6917

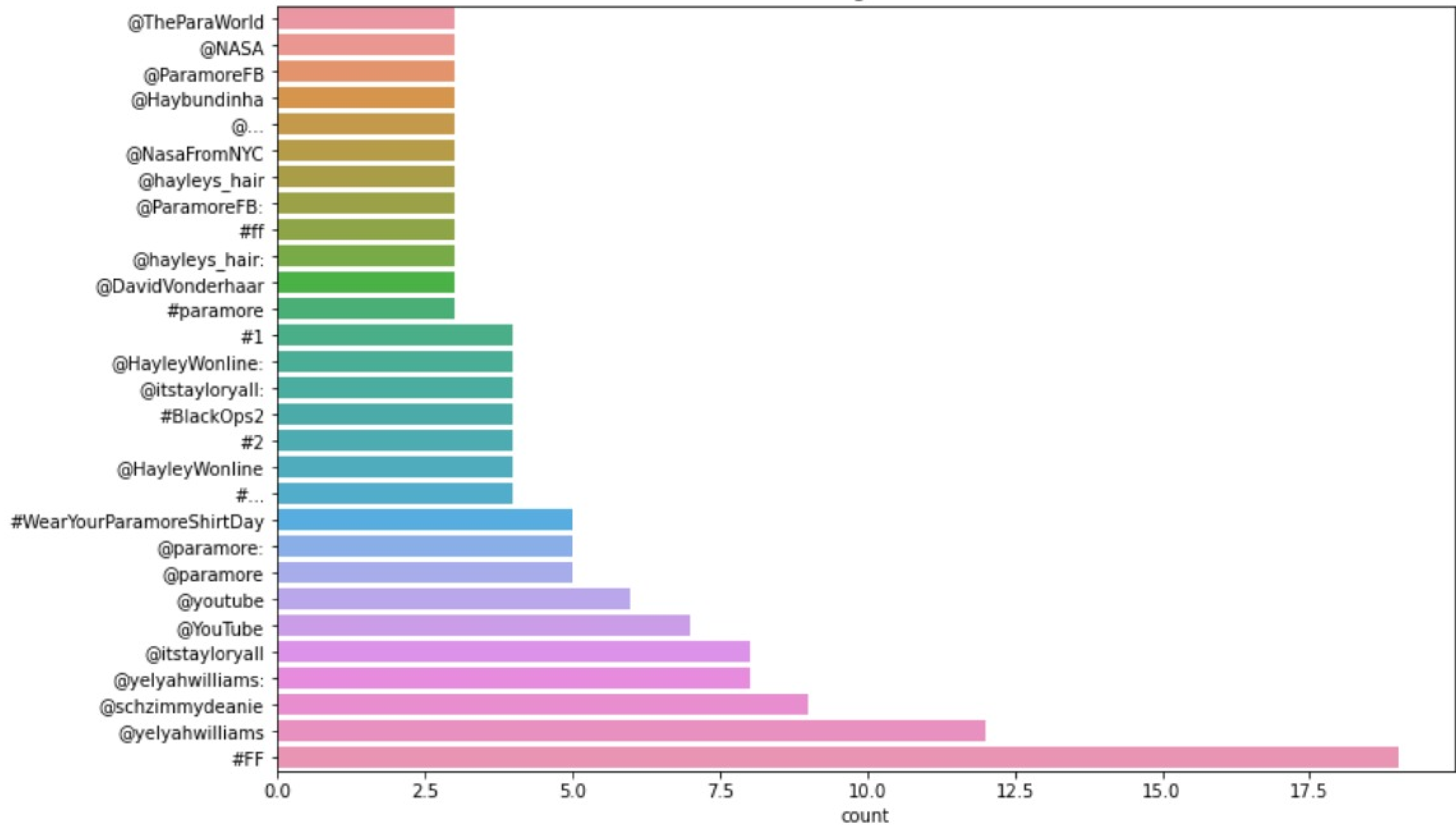
Music Cluster



Social Media Cluster



Gaming Cluster



Conclusion



Thank you for your attention

Feel free to ask questions