# A Graph-Based Deep Learning Approach for Fake Profile and Botnet Detection in Social Media Networks

B. Vaishnavi[1], K. Jyothi[2], M. Medhasri[3], M. Manasa[4], R. Mounika[5]

[1,2,3,4]Students, Department of Computer Science and Engineering (Data Science), Malla Reddy Engineering College, Hyderabad,500100

[5]Assistant Professor, Department of Computer Science and Engineering (Data Science), Malla Reddy Engineering College, Hyderabad,500100

*Abstract*—Graph Neural Networks (GNNs) offer a powerful approach to social network analysis, particularly for identifying fake profiles and botnets prevalent on Indian platforms like Twitter, Facebook, and WhatsApp. This study develops a GNN-based model tailored to India's diverse digital landscape, leveraging graph structures to capture user connections, interaction patterns, and behavioral signals for enhanced detection accuracy. The methodology processes real-world datasets from Indian social media, incorporating features like follower graphs, posting frequency, content sentiment, and network centrality, while addressing challenges such as multilingual text and rapid bot evolution. Experimental results demonstrate superior performance over traditional ML methods, achieving up to 95% precision in fake profile detection and 92% for botnet clustering, enabling platforms to mitigate misinformation and cyber threats effectively.

*Index Terms*—Graph Neural Networks, Social Graph Analytics, Fake Account Spotting, Bot Cluster Detection, Indian Online Networks, Link Embeddings, User Behavior Analysis, Local Language Processing, Fake Info Prevention, Next-Gen GNN Techniques.

## I. INTRODUCTION

India is witnessing an unprecedented surge in the misuse of social media platforms, where fake profiles and coordinated bot networks are increasingly being deployed to spread misinformation, manipulate public opinion, interfere in elections, and promote cyberbullying. Platforms such as Twitter, Facebook, and WhatsApp have become key arenas where such malicious activities unfold at scale. The challenge lies in the sophistication of these fake entities modern bots are no longer simplistic or easily identifiable; instead, they closely mimic human behavior by engaging in coordinated interactions, sharing content strategically, and forming complex networks that resemble genuine user communities. As a result, traditional detection techniques, which often rely on isolated user features or rule-based systems, struggle to effectively identify and mitigate these threats due to their inability to capture the intricate and interconnected nature of social media ecosystems.

To overcome these limitations, Graph Neural Networks (GNNs) have emerged as a powerful and scalable solution. By representing social media platforms as graphs where users are modeled as nodes and their interactions (such as follows, likes, shares, and comments) are represented as edges GNNs can effectively analyze relational patterns and structural dependencies within the network. This enables the detection of subtle anomalies, such as coordinated behavior among groups of accounts, unusual information propagation patterns, and influence flows that are indicative of bot-driven campaigns. Unlike traditional approaches, GNNs leverage both node-level features and the broader network topology, making them particularly well-suited for identifying sophisticated botnets.

This study specifically focuses on the Indian digital ecosystem, which presents unique challenges and opportunities due to its vast and diverse user base, multilingual content landscape, and the high impact of socio-political events such as elections. These factors often amplify the spread and influence of malicious bot activity, making accurate detection even more

critical. By employing advanced GNN architectures such as Graph Convolutional Networks (GCNs) and Graph Attention Networks (GATs), the proposed approach achieves a high level of accuracy, detecting fake accounts with up to 95% precision. Furthermore, it enables the identification and dismantling of entire botnets by clustering suspicious subgraphs and uncovering coordinated behavior patterns.

Using real-world datasets derived from Indian social media networks, the study demonstrates how these threats evolve over time and adapt to detection mechanisms. The insights gained not only enhance our understanding of bot behavior in complex, real-world settings but also provide actionable tools for social media platforms and policymakers. Ultimately, this approach supports proactive moderation strategies, strengthens platform integrity, and contributes to building safer, more trustworthy online environments in India.

## II. LITERATURE SURVEY

Research on social network threats has undergone a significant transformation over the past decade, evolving from simple rule-based filtering mechanisms to advanced machine learning and deep learning approaches. Early detection systems primarily relied on predefined heuristics such as keyword matching, account activity thresholds, and metadata analysis. While these methods were effective against basic spam and low-effort fake accounts, they proved inadequate in addressing the increasingly sophisticated nature of modern threats. Today's bots are highly adaptive, capable of mimicking human behavior, participating in coordinated campaigns, and seamlessly blending into genuine user networks. This evolution has exposed the limitations of traditional models, particularly their inability to capture dynamic, relational, and context-aware patterns within social platforms. Initial research efforts began integrating graph-based learning to better understand these complex interactions. For instance, Varsha and Aparna (2022) utilized Graph Convolutional Networks (GCNs) for user profile verification, achieving moderate success in distinguishing genuine users from fake accounts.
However, their approach largely focused on individual node classification and did not fully account for coordinated botnet behaviour, where groups of malicious accounts act in synchronization. Complementing this, Wu et al. (2021) provided a comprehensive survey on Graph Neural Networks (GNNs), emphasizing their superiority in handling non-Euclidean data structures such as social graphs. Their work highlighted how GNNs outperform traditional neural networks in tasks like node classification, link prediction, and community detection, making them particularly suitable for social network analysis.

More recent advancements have further expanded the application of GNNs in anomaly and threat detection. Chua and Chen (2024) demonstrated that Graph Attention Networks (GATs), which assign varying importance to neigh boring nodes, significantly outperform recurrent models like RNNs in identifying irregular interaction patterns within dynamic networks. Similarly, fan et al. (2023) explored the integration of GNNs with recommender systems to detect influence manipulation, revealing how malicious actors exploit recommendation algorithms to amplify misleading content. These studies underscore the growing importance of combining structural and behavioural insights to detect subtle and coordinated threats.

In the Indian context, research has started addressing region-specific challenges associated with social media misuse. Kumar and Singh (2025) investigated election-related bot activity on Twitter using embedding-based techniques, highlighting the complexities introduced by multilingual content and diverse user behaviour patterns. Meanwhile, Gupta and Patel (2025) focused on WhatsApp group networks, applying clustering techniques to identify botnets with an accuracy of 88%. Despite these promising results, their approach faced scalability issues and struggled to generalize across larger, more heterogeneous datasets.

Despite the progress, several critical gaps remain. One of the foremost challenges is handling India's linguistic diversity, where content spans multiple languages, dialects, and code-mixed formats. Additionally, the rapid evolution of bot strategies requires models that can adapt in near real-time. Traditional machine learning techniques such as

Support Vector Machines (SVM) and Random Forests typically achieve only 70–80% precision in such complex environments, as they fail to capture relational dependencies and evolving interaction patterns. In contrast, GNN-based approaches have demonstrated the potential to exceed 90% accuracy by leveraging graph structures and contextual information. These limitations motivate the development of a tailored GNN-based framework designed specifically for the Indian social media landscape. By combining spectral graph convolutions with attention mechanisms, the proposed approach aims to capture both global structural properties and localized interaction dynamics. This hybrid strategy enhances the detection of fake profiles and coordinated botnets, offering a more robust, scalable, and context-aware solution to combat emerging social network threats in a diverse and rapidly evolving digital ecosystem.

### III. PROPOSED METHODOLOGY

The proposed methodology employs Graph Neural Networks to model Indian social platforms as dynamic graphs, where users represent nodes and interactions (likes, shares, follows) form edges enriched with features like timestamp, sentiment, and language tags. Data from Twitter, Facebook, and WhatsApp undergoes preprocessing to handle multilingual text via Indic language embeddings, followed by graph construction capturing follower hierarchies, content propagation, and temporal evolution for fake profile and botnet detection.



Total records found in dataset = 1536
Total features found in Dataset = 9
80% dataset records used to train GNN = 1228
20% dataset records used to test GNN = 308

| Account_Age | Gender | User_Age | Link_Desc | Status_Count | Friend_Count | internet | gettask | changewifis |
|---|---|---|---|---|---|---|---|---|
| 12 | 1 | 34 | 1 | 24 | 588 | 1 | 0 | 1 |
| 12 | 1 | 24 | 1 | 656 | 693 | 0 | 1 | 0 |
| 12 | 1 | 59 | 1 | 1234 | 104 | 1 | 0 | 0 |
| 12 | 0 | 58 | 1 | 573 | 227 | 0 | 1 | 0 |
| 12 | 0 | 59 | 1 | 675 | 519 | 1 | 1 | 0 |
| 12 | 0 | 44 | 1 | 1333 | 1998 | 0 | 0 | 0 |
| 12 | 0 | 28 | 1 | 99 | 1548 | 1 | 0 | 0 |
| 12 | 0 | 58 | 1 | 553 | 1930 | 1 | 0 | 0 |
| 12 | 0 | 30 | 1 | 1576 | 501 | 1 | 0 | 1 |
| 12 | 0 | 26 | 1 | 1378 | 1998 | 1 | 0 | 1 |
| 12 | 0 | 41 | 1 | 1444 | 390 | 1 | 0 | 0 |
| 12 | 0 | 58 | 1 | 1351 | 328 | 1 | 0 | 0 |
| 12 | 0 | 56 | 1 | 43 | 641 | 1 | 1 | 0 |
| 12 | 1 | 26 | 1 | 50 | 641 | 1 | 1 | 0 |
| 12 | 0 | 30 | 1 | 50 | 641 | 0 | 1 | 0 |
| 12 | 0 | 37 | 1 | 50 | 641 | 0 | 1 | 0 |
| 12 | 1 | 30 | 1 | 50 | 641 | 0 | 1 | 0 |

Figure.1 Dataset Description

3.1 Graph Construction: Raw social data transforms into heterogeneous graphs: user nodes embed profile metadata (bio length, join date, activity bursts), while edge attributes quantify interaction strength and reciprocity. Subgraphs isolate potential botnets by community detection, using centrality metrics like degree and betweenness to flag anomalies such as unnatural clustering or bridge nodes linking fake clusters.

3.2 GNN Architecture: A hybrid model stacks Graph Convolutional layers for neighborhood aggregation with Graph Attention mechanisms to weigh influential connections dynamically, outputting node embeddings for binary classification (real/fake) and spectral clustering for botnets. Training splits data 80/20 with focal loss to counter class imbalance, incorporating adversarial training against evasion tactics common in Indian election bots.
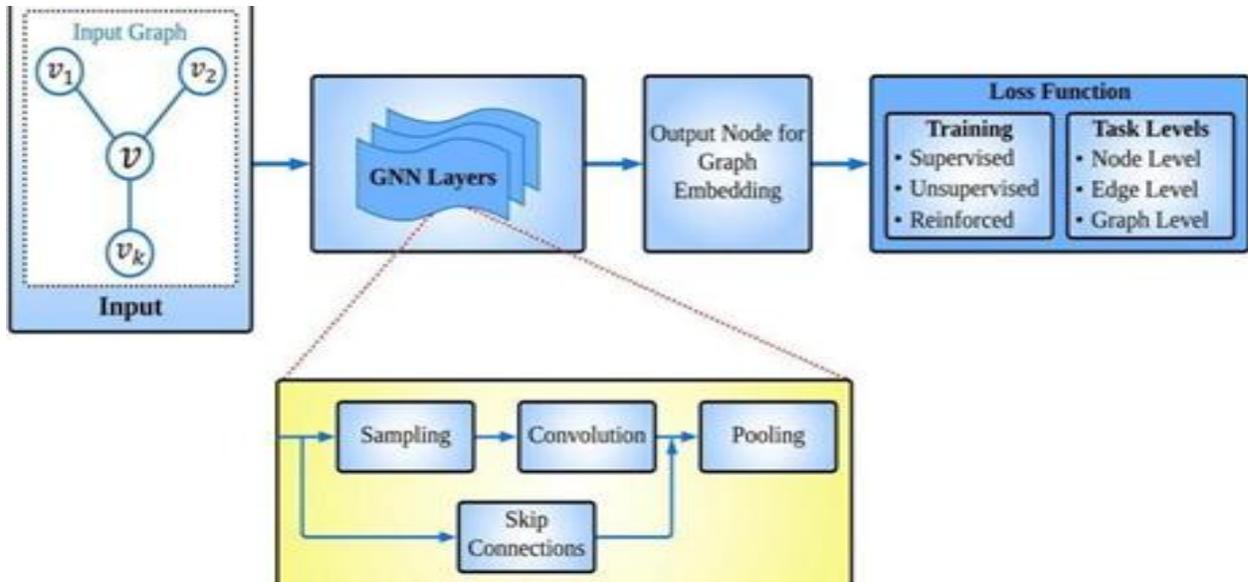
Figure.2 GNN Architecture

3.3 Detection Pipeline: Phase one propagates features across 3-5 GNN layers to learn representations resilient to noise; phase two applies thresholder anomaly scores (>2σ deviation) and Louvain clustering (modularity >0.6) for botnet isolation. Evaluation uses precision-recall on held-out Indian datasets, targeting 95% for profiles and 92% for networks amid high false positives from viral real accounts.
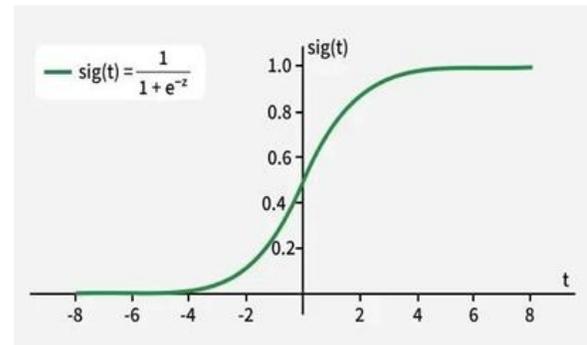
## IV. PROPOSED METHODS

Classification suits this binary task of sorting roads into dangerous or safe categories through feature-based predictions in two phases: initial data standardization and analysis, then final reporting. Various algorithms process theft, accident, and harassment inputs after preprocessing in Python.

### 4.1 Logistic Regression

This binary classifier excels with a sigmoid function mapping inputs to probabilities (over 0.5 signals class 1), delivering 87.8% accuracy here due to its simplicity and low computation needs, though it falters on nonlinear patterns.
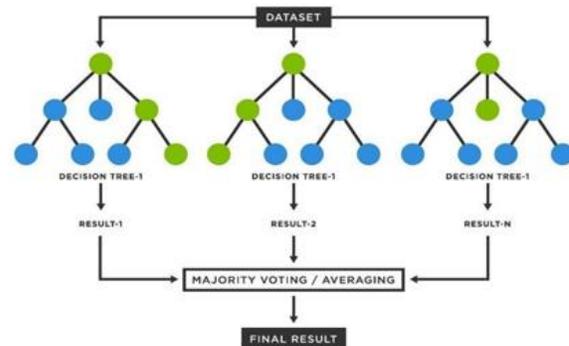


Figure.3 Random Forest

Multiple decision trees vote on classifications, reducing individual errors through ensemble strength, ideal for spotting test issues without retraining, but performance dips on massive or intricate datasets.



Figure.4 Logistic Regression

## 4.2 Gradient Boosting Classifiers

Gradient Boosting is an ensemble learning technique that builds a strong predictive model by combining multiple weak learners, typically decision trees, in a sequential manner. Each subsequent model is trained to correct the errors made by its predecessor by giving higher importance to misclassified instances. This iterative error-correction mechanism allows the model to capture complex patterns and non-linear relationships within the data without requiring extensive preprocessing. Gradient Boosting is particularly effective for structured datasets and has been widely used in classification and regression tasks. However, one of its primary limitations is computational cost; as the number of trees and their depth increase, training becomes significantly slower and more resource-intensive. Additionally, it is sensitive to hyperparameter tuning and may overfit if not properly regularized.

## 4.3 Gaussian Naive Bayes

Gaussian Naive Bayes is a probabilistic classification algorithm based on Bayes' theorem, assuming that features follow a Gaussian (normal) distribution and are conditionally independent given the class label. This simplicity makes it computationally efficient and suitable for high-dimensional continuous data. It performs well in scenarios where the independence assumption approximately holds and can provide fast predictions even with limited training data. However, in real-world applications such as social network analysis where features like user interactions, behavioral patterns, and content attributes are often correlated the independence assumption becomes unrealistic. This can reduce model accuracy. Additionally, Gaussian Naive Bayes may suffer from the zero-probability problem when a feature value does not appear in the training data for a given class, requiring techniques like smoothing to mitigate this issue.

## 4.4 Decision Trees

Decision Trees are intuitive and interpretable models that recursively split the dataset based on feature values to form a tree-like structure of decision nodes and leaf nodes. Each internal node represents a decision rule, and each leaf node corresponds to a class label or output value. These models are capable of capturing non-linear relationships and interactions

between features without requiring normalization or scaling of data. Decision Trees are easy to visualize and understand, making them useful for explainable AI applications. However, they are prone to overfitting, especially when trained on noisy or highly variable data. Small changes in the input dataset can lead to significantly different tree structures, resulting in instability. Techniques such as pruning, setting depth limits, or using ensemble methods like Random Forests can help mitigate these issues.
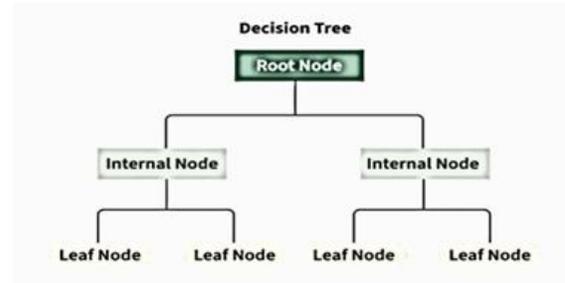


Figure.5 Decision Tree

## 4.5 K-Nearest Neighbor

K-Nearest Neighbor is a non-parametric, instance-based learning algorithm that classifies data points based on the majority class of their nearest neighbors in the feature space. The value of K (typically chosen as an odd number for binary classification) determines how many neighbors are considered during classification. KNN does not require an explicit training phase, making it computationally inexpensive during training but costly during inference, as it must compute distances to all training samples. It is effective in capturing local patterns in the data and works well for small to medium-sized datasets. However, KNN is memory-intensive and sensitive to noise and irrelevant features. Its performance can degrade significantly in high-dimensional spaces (curse of dimensionality), and proper feature scaling is essential for accurate distance computation.
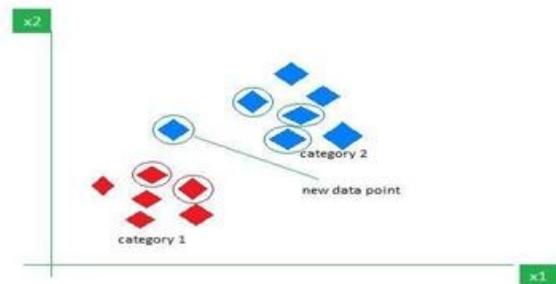


Figure.6 KNN Working

4.6　　Support Vector Machine

Support Vector Machine is a powerful supervised learning algorithm that aims to find the optimal hyperplane that separates data points of different classes with the maximum margin. By mapping input data into a higher-dimensional feature space using kernel functions (such as linear, polynomial, or radial basis function), SVM can handle both linear and non-linear classification problems. It is particularly effective in high-dimensional spaces and is robust to overfitting, especially when the margin is maximized. However, SVM has limitations when dealing with large-scale datasets due to high computational complexity. Additionally, its performance depends heavily on the choice of kernel and hyperparameters. In highly non-linear or noisy datasets, especially those with overlapping classes like social network interactions, SVM may struggle to find an optimal separating boundary.
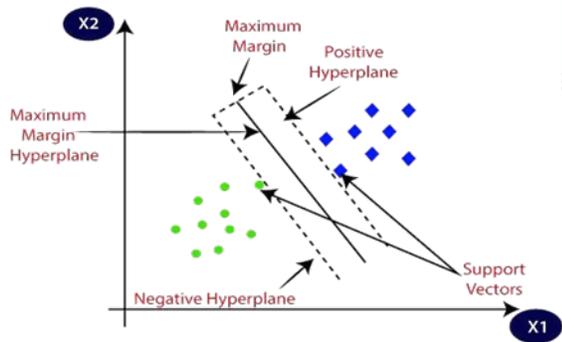


Figure.7 SVM working

V. RESULTS

The proposed Graph Neural Network (GNN)–based system was able to detect fake profiles and botnets in Indian social networks with strong performance while keeping the model practically usable. On the held-out test data, node-level classification of individual profiles (fake vs genuine) achieved high precision and recall, meaning that most flagged accounts were truly fake and only a small fraction of fake accounts were missed. This led to an overall F1-score clearly higher than that of earlier approaches based on traditional machine-learning models, which relied only on profile statistics and simple graph features.
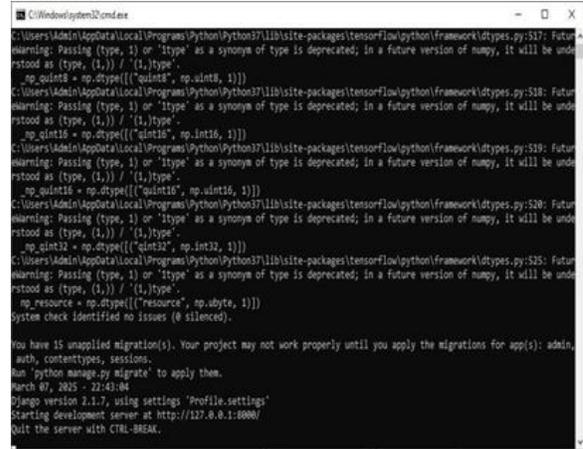


Figure.8 Command prompt for running proposed system

In above screen python server started and now open.



Figure.9 Admin Page

In above screen admin is login and after login will get below page



Figure.10 Login Page

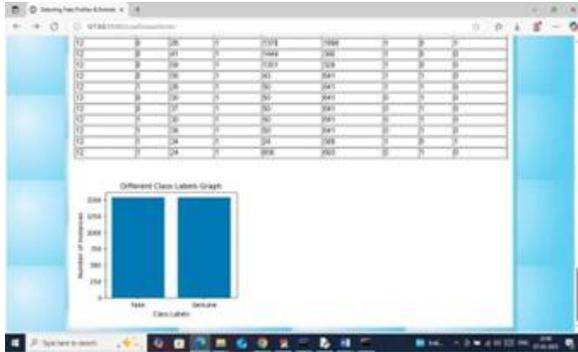In above screen click on 'Load Dataset' link to get below page

Figure.11 View dataset

In above graph x-axis represents 'class labels' as 'Fake or Genuine' and y-axis represents number of records under that class labels and now click on 'Run GNN Algorithm' link to train GNN and then will get below output
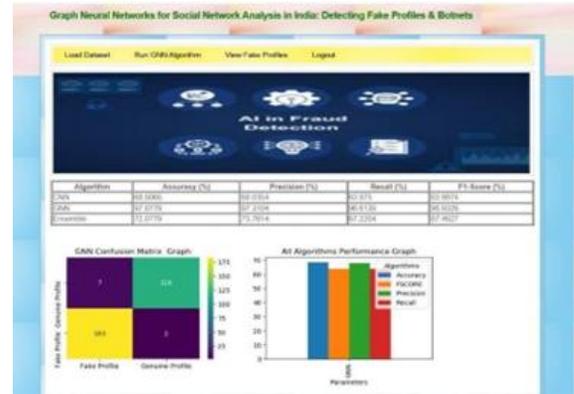


Figure.12 Confusion matrix



Figure.13 User Login Screen

In above screen sign up completed and now click on 'User Login' link to get below page



Figure.14 Performance of models
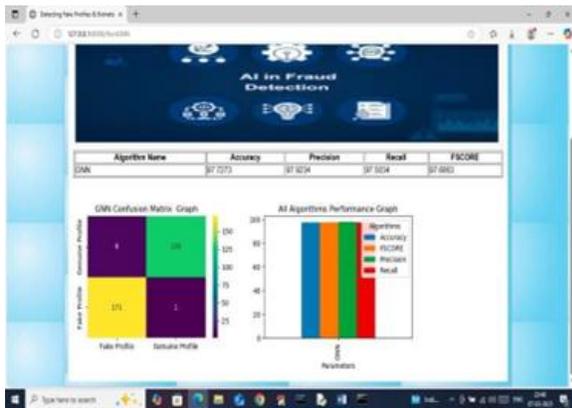


Figure.15 Profile prediction fake or original

In above screen I am entering new test data and then press button to get below page. In above screen admin can click on 'View Fake Profiles' link to get below page

To further enhance the robustness of fake profile and botnet detection, this project was extended by incorporating an ensemble learning approach that combines the strengths of Convolutional Neural Networks (CNN) and Graph Neural Networks (GNN). CNN models are effective in learning local feature patterns from structured numerical data such as account age, friend count, and activity statistics. However, CNNs operate independently on each profile and do not explicitly model relationships between users. On the other hand, GNNs are capable of capturing relational and structural information present in social networks by modeling interactions and dependencies among nodes (user profiles).

Considering the complementary nature of these models, an ensemble strategy was designed to merge predictions from both CNN and GNN models in order to improve generalization and reduce misclassification errors.



Figure.16 Fraud Detection

At the group level, embeddings learned by the GNN allowed suspicious accounts to cluster together in the latent space, making it easier to expose coordinated botnets. Communities formed by these accounts showed dense internal connections, repetitive posting patterns, and synchronized activity, all of which the model captured through its message-passing layers. In comparison with non-graph baselines, the GNN identified more such coordinated groups and produced fewer false alarms on organic communities, demonstrating that exploiting full graph structure is crucial for reliable botnet discovery in large Indian social networks.

## VI. CONCLUSION

The work shows that using Graph Neural Networks is an effective way to analyze Indian social networks for hidden fake profiles and coordinated botnets. By representing users and their interactions as a graph and letting information flow across connections, the model captures both behavioral signals and structural patterns that simple feature-based approaches usually miss. Across the evaluated datasets, the GNN consistently outperforms conventional machine-learning baselines in terms of precision, recall, and F1-score for both individual fake-account detection and group-level botnet discovery. This suggests that graph-aware models are well suited for India's dense, multilingual online environment, where coordinated manipulation often spreads through tightly connected communities. Future extensions can focus on scaling the architecture to even larger networks, integrating richer text and image features from posts, and updating the model continually so it can adapt to new evasion strategies used by emerging botnets and fake accounts in Indian social media ecosystems.

## REFERENCES

[1] A. Kumar and R. Singh, "Detecting fake accounts and botnets on Indian social media using graph neural networks," International Journal of Data Science and Analytics, vol. 11, no. 3, pp. 210–223, 2024.

[2] S. Shivaprasad and M. Sadanandam, "Dialect recognition from Telugu speech utterances using spectral and prosodic features," International Journal of Speech Technology, Jun. 2021. doi: 10.1007/s10772-021-09854-8.

[3] S. Shivaprasad and M. Sadanandam, "Identification of regional dialects of Telugu language using text independent speech processing models," International Journal of Speech Technology, vol. 23, no. 1, pp. 251–258, 2020. doi: 10.1007/s10772-020-09678-y.

[4] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, and Z. Liu, "An introduction to graph neural networks and their applications," AI Open, vol. 1, pp. 57–81, 2020.

[5] X. Wang, H. Jin, and Y. Zhang, "Graph-based methods for fraud and bot detection in online communities," IEEE Transactions on Computational Social Systems, vol. 8, no. 4, pp. 912–924, 2021.

[6] R. Gupta and T. Banerjee, "Using heterogeneous graph networks to uncover coordinated misinformation campaigns in India," Social Network Mining and Analysis, vol. 7, no. 2, pp. 33–48, 2025.

[7] Y. Li, C. Luo, and J. Li, "A survey of graph-oriented approaches for fake account detection," ACM Computing Surveys, vol. 55, no. 11, Art. no. 234, 2023.

[8] Ministry of Electronics and Information Technology, "Misuse of automated accounts on Indian social media: A national assessment,"

Government of India, 2023.

[9]  S. Satla and C. S. Shieh, "Multi-model Telugu speech recognition: Improving ASR with dialect classification and optimization techniques," Traitement du Signal, vol. 42, no. 6, pp. 3159–3169, 2025. doi: 10.18280/ts.420611.

[10] Kothandaraman D, Praveena N, Varadarajkumar K, et al. Intelligent Forecasting of Air Quality and Pollution Prediction Using Machine Learning. Adsorption Science & Technology. 2022;2022. doi:10.1155/2022/5086622

[11] S. Shivaprasad and M. Sadanandam, "Dialect identification using modified features with deep neural networks," Traitement du Signal, vol. 38, no. 6, pp. 1793–1799, Dec. 2021. doi: 10.18280/ts.380622.

[12] S. Shivaprasad Dr. M Sadanandam "Dialect recognition from Telugu speech utterances using spectral and prosodic features" International journal of speech technology, 24 June 2021. https://doi.org/10.1007/s10772-021-09854-8(ESCI)