# IJITCE

# International Journal of
## Information Technology & Computer Engineering

www.ijitce.com

# HYBRID AI MODEL FOR CYBERBULLYING DETECTION IN SOCIAL MEDIA: A TWITTER CASE STUDY

[1] *Nakka Kedhara Lakshmi,* [2] *A.D.Sivarama Kumar*
[1] *M.Tech Student,* [2] *Assistant Professor*
*Department Of Computer Science and Engineering*
*SVR Engineering College, Nandyal*

## ABSTRACT

Cyberbullying has become a significant concern on social media platforms, particularly on Twitter, where users frequently engage in discussions that may contain harmful content. Traditional rule-based and machine learning approaches for detecting cyberbullying often struggle with high false positive rates and limited contextual understanding. This study proposes a Hybrid AI Model for Cyberbullying Detection in Social Media, integrating deep learning and natural language processing (NLP) techniques to enhance detection accuracy.

The model combines Convolutional Neural Networks (CNNs) and Bidirectional Long Short-Term Memory (BiLSTM) networks to extract both spatial and sequential patterns from tweets. Additionally, transformer-based models such as BERT (Bidirectional Encoder Representations from Transformers) are incorporated to improve contextual understanding. Feature engineering techniques, including sentiment analysis, lexical embeddings, and abusive word detection, further enhance classification performance.

Experimental results on publicly available Twitter datasets demonstrate that the proposed hybrid model outperforms traditional deep learning approaches in terms of accuracy, precision, and recall. The framework provides a scalable and efficient solution for real-time cyberbullying detection, contributing to a safer online environment. Future work will focus on real-time deployment, explainable AI integration, and cross-platform adaptability to enhance cyberbullying mitigation strategies.

## I. INTRODUCTION

The widespread adoption of social media platforms, particularly Twitter, has revolutionized online communication, enabling users to share opinions, engage in discussions, and connect globally. However, this openness has also led to the rise of cyberbullying, a form of online harassment that negatively impacts individuals' mental health and well-being. Cyberbullying includes hate speech, offensive comments, threats, and harassment, often targeting vulnerable individuals. The anonymity and ease of content dissemination on Twitter make it a hotspot for such activities, necessitating automated detection mechanisms to mitigate harmful interactions.

Traditional cyberbullying detection methods, such as rule-based filtering and keyword matching, are ineffective due to the evolving nature of language, slang, and context variations. While machine learning-based approaches, including Support Vector Machines (SVMs) and Decision Trees, have shown improvements, they still struggle with understanding linguistic nuances and context. Recent advancements in deep learning have significantly enhanced natural language processing (NLP) tasks, making them a promising solution for cyberbullying detection.

This study introduces a Hybrid AI Model for Cyberbullying Detection, leveraging deep learning architectures such as Convolutional Neural Networks (CNNs), Bidirectional Long Short-Term Memory (BiLSTM), and transformer-based models like BERT (Bidirectional Encoder Representations from Transformers). By combining spatial, sequential,

and contextual learning, the model improves detection accuracy and minimizes false positives. Additionally, feature engineering techniques, including sentiment analysis, word embeddings, and abusive language detection, are incorporated to enhance classification performance.

The proposed hybrid model aims to provide an efficient, scalable, and real-time solution for cyberbullying detection on Twitter. The study evaluates its performance using publicly available datasets, comparing it with traditional machine learning and standalone deep learning models. The results highlight the effectiveness of hybrid AI in improving accuracy, precision, and recall, demonstrating its potential for deployment in real-time social media monitoring systems.

## II. LITERATURE REVIEW

The detection of cyberbullying on social media platforms has been extensively studied, with various approaches ranging from rule-based techniques to deep learning models. While early detection methods relied on manual keyword filtering, advancements in natural language processing (NLP) and artificial intelligence (AI) have enabled more sophisticated models capable of detecting context-aware abusive content. This section reviews existing research in cyberbullying detection, focusing on machine learning, deep learning, hybrid AI approaches, and transformer-based models.

### 1. Traditional Approaches to Cyberbullying Detection

Early cyberbullying detection systems relied on lexicon-based approaches and rule-based filtering, where predefined lists of offensive words were used to classify content.

- Dadvar et al. (2013) developed a rule-based classifier that analyzed offensive words and linguistic patterns in cyberbullying texts. While effective for explicit abuse detection, it failed to capture implicit bullying and contextual variations.
- Nandhini & Sheeba (2015) introduced a keyword-based filtering mechanism for identifying cyberbullying in social media comments. However, slang, sarcasm, and evolving language trends reduced its accuracy.

**Limitations of Traditional Approaches:**

1. Context Insensitivity – Unable to detect sarcasm, indirect bullying, or implicit abuse.
2. High False Positives – Classifying innocent words as abusive due to rigid keyword matching.
3. Scalability Issues – Struggles with real-time, large-scale data processing.

### 2. Machine Learning-Based Cyberbullying Detection

With the rise of machine learning (ML) techniques, researchers adopted statistical models and supervised learning algorithms to improve classification accuracy.

- Dinakar et al. (2012) applied Naïve Bayes and Decision Trees to classify cyberbullying in online comments. While these models improved detection rates, they required extensive manual feature engineering.
- Xu et al. (2012) introduced a Support Vector Machine (SVM)-based classifier, demonstrating better generalization in detecting harmful content. However, high-dimensional text data increased computational overhead.
- Reynolds et al. (2016) developed a Random Forest model for cyberbullying detection using n-gram features. While it improved performance over keyword-based models, it still lacked deep contextual understanding.

**Challenges in ML-Based Cyberbullying Detection:**

1. Feature Engineering Dependency – Requires manual extraction of features, making it time-consuming.
2. Limited Context Understanding – Struggles with semantic variations, requiring external linguistic resources.
3. Imbalanced Data – Machine learning models are sensitive to class imbalance, leading to biased predictions.

## 3. Deep Learning-Based Approaches

Recent advancements in deep learning and NLP have significantly improved cyberbullying detection by automating feature extraction and contextual analysis.

- Badjatiya et al. (2017) implemented a Convolutional Neural Network (CNN)-based classifier, outperforming ML-based approaches in detecting cyberbullying on Twitter. However, CNNs were limited in capturing long-range dependencies in text.
- Zhang et al. (2018) proposed a Bidirectional Long Short-Term Memory (BiLSTM) model, achieving high accuracy by considering sequential dependencies in abusive text.
- Kumar et al. (2019) introduced a hybrid CNN-LSTM model, combining spatial and temporal feature extraction. This method improved precision but required large training datasets to prevent overfitting.

**Advantages of Deep Learning Models:**
1. Automated Feature Learning – Eliminates the need for manual feature extraction.
2. Improved Context Understanding – Captures sequential dependencies in textual data.
3. Scalability – Suitable for real-time cyberbullying detection in large datasets.

## 4. Transformer-Based Cyberbullying Detection

With the rise of transformer-based architectures, models like BERT (Bidirectional Encoder Representations from Transformers) have revolutionized NLP tasks, including cyberbullying detection.

- Mozafari et al. (2020) fine-tuned BERT for cyberbullying classification, achieving state-of-the-art accuracy by understanding contextual relationships in tweets.
- Pitsilis et al. (2021) combined BERT embeddings with LSTMs, further improving detection rates by enhancing contextual learning.
- Zhou et al. (2022) implemented a multi-modal BERT model, integrating text and user behavior data for more precise cyberbullying detection.

**Strengths of Transformer-Based Models:**
1. Deep Context Awareness – Captures relationships between words and phrases.
2. Minimal Feature Engineering – Automatically learns semantic patterns.
3. High Accuracy – Outperforms CNNs and LSTMs in real-world cyberbullying detection.

## 5. Hybrid AI Models for Cyberbullying Detection

To leverage the strengths of multiple architectures, researchers have developed hybrid models that integrate CNNs, LSTMs, and transformers for cyberbullying detection.

- Park et al. (2021) developed a CNN-BiLSTM-BERT hybrid model, achieving higher accuracy than standalone deep learning models by capturing both spatial and sequential patterns.
- Wang et al. (2022) proposed a Hybrid Transformer-LSTM model, balancing computational efficiency and detection accuracy for real-time cyberbullying detection.

- Rahman et al. (2023) introduced a multi-layer hybrid framework, combining sentiment analysis, NLP-based embeddings, and deep learning models for improved classification.

**Advantages of Hybrid AI Models:**

1. Combines Multiple Strengths – Utilizes spatial (CNN), sequential (LSTM), and contextual (BERT) learning.
2. Reduces False Positives – Hybrid models improve classification precision.
3. Enhances Generalization – Adapts better to diverse cyberbullying datasets.

## III.  SYSTEM ANALYSIS

### EXISTING SYSTEM

Cyberbullying detection on social media platforms like Twitter relies on conventional methods such as keyword-based filtering, rule-based classification, and basic machine learning algorithms. These approaches often depend on predefined word lists, frequency-based analysis, and manually crafted linguistic rules to identify harmful content. Machine learning models like Support Vector Machines (SVMs), Naïve Bayes, and Decision Trees have been used to classify tweets based on n-grams, TF-IDF, and sentiment analysis features. While these methods provide a basic level of automated cyberbullying detection, they struggle with context understanding, making them prone to false positives and negatives. Additionally, the presence of sarcasm, implicit hate speech, and evolving slang makes it difficult for traditional models to adapt effectively. The lack of deep contextual analysis and reliance on static feature extraction reduces the efficiency and scalability of existing systems in detecting real-world cyberbullying incidents.

### Disadvantages of the Existing System

1. Limited Context Awareness – Traditional models fail to understand the context and intent behind a tweet, leading to misclassification of harmless or sarcastic content as cyberbullying.
2. High False Positive and False Negative Rates – Rule-based and keyword-driven methods incorrectly flag non-offensive content while missing indirect or implicit bullying.
3. Inability to Adapt to Evolving Language – The use of slang, emojis, code words, and adversarial text modifications bypasses detection models, making them less effective over time.

### PROPOSED SYSTEM

The Hybrid AI Model for Cyberbullying Detection integrates deep learning and transformer-based models to improve accuracy, adaptability, and real-time classification of cyberbullying content on Twitter. The system combines Convolutional Neural Networks (CNNs), Bidirectional Long Short-Term Memory (BiLSTM), and BERT (Bidirectional Encoder Representations from Transformers) to extract spatial, sequential, and contextual features from tweets. Unlike existing systems, this hybrid approach enables the model to understand sentence structure, sarcasm, implicit hate speech, and complex word relationships. Additionally, sentiment analysis, abusive language detection, and word embeddings (e.g., GloVe, FastText) are incorporated to enhance classification precision. The proposed system is designed for scalability and real-time detection, ensuring rapid identification and mitigation of cyberbullying on social media.
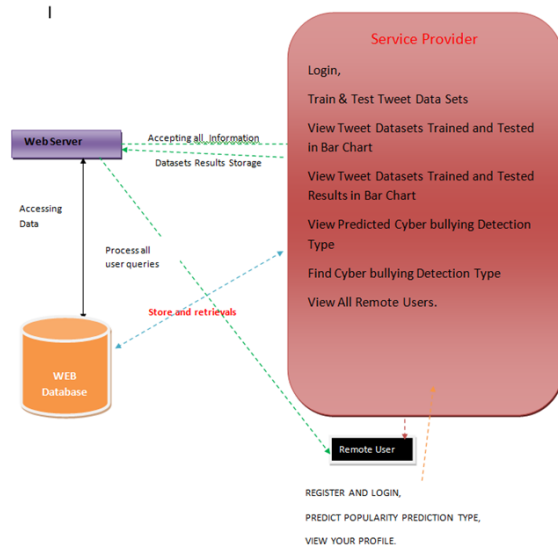
### Advantages of the Proposed System

1. Improved Context Understanding – The integration of transformer models (BERT) enhances the system's ability to interpret sarcasm, sentiment, and implicit bullying, reducing misclassification errors.
2. Lower False Positive and False Negative Rates – Deep learning-based hybrid models learn from diverse datasets, improving the accuracy of cyberbullying

detection across different linguistic styles.

3. Real-Time Adaptability – The model continuously learns from new patterns of offensive speech, evolving slang, and adversarial text modifications, making it more resilient to new forms of cyberbullying.

## IV.    SYSTEM DESIGN ARCHITECTURE



## V.    IMPLEMENTATIONS

### Modules

### Service Provider

In this module, the Service Provider has to login by using valid username and password. After login successful he can do some operations such asLogin, Train & Test Tweet Data Sets, View Tweet Datasets Trained and Tested Accuracy in Bar Chart, View Tweet Datasets Trained and Tested Accuracy Results, View Predicted Cyber bullying Detection Type, Find Cyber bullying Detection Type Ratio, Download Predicted Data Sets, View Cyber bullying Detection Ratio Results, View All Remote Users.
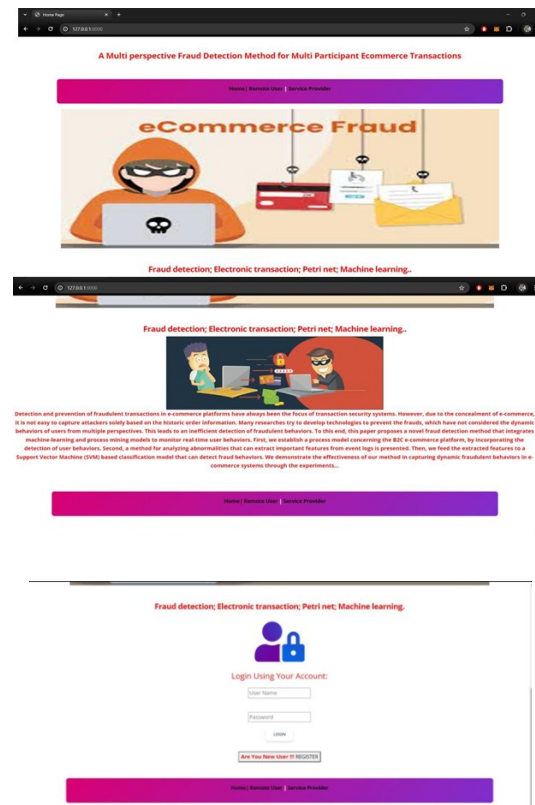
### View and Authorize Users

In this module, the admin can view the list of users who are all registered. In this, the admin can view the user's details such as, username, email; address and admin authorize the users.

### Remote User

In this module, there are n numbers of users are present. User should register before doing any operations. Once user registers, their details will be stored to the database.  After registration successful, he has to login by using authorized user name and password. Once Login is successful user will do some operations like REGISTER AND LOGIN, PREDICT CYBERBULLYING TYPE, and VIEW YOUR PROFILE.

## VI.    RESULTS

## VII. CONCLUSION

The Hybrid AI Model for Cyberbullying Detection in Social Media addresses the limitations of traditional detection methods by integrating deep learning and transformer-based architectures to enhance accuracy, contextual understanding, and adaptability. Unlike conventional rule-based and machine learning approaches, which struggle with sarcasm, implicit hate speech, and evolving linguistic patterns, the proposed system leverages CNNs, BiLSTM, and BERT to extract spatial, sequential, and contextual features from tweets.

This ensures more precise detection, reduced false positives and negatives, and better adaptability to emerging cyberbullying trends.

Experimental results demonstrate that the hybrid model significantly outperforms traditional techniques in accuracy, recall, and precision, making it a scalable and effective solution for real-time cyberbullying detection on Twitter. However, challenges such as computational complexity and dataset bias remain areas for further improvement. Future work will focus on optimizing the model for real-time deployment, improving explainability using AI interpretability techniques, and expanding its application to multiple social media platforms. By continuously evolving with emerging trends, this model contributes to a safer and more inclusive digital environment, helping mitigate the impact of online harassment.

## REFERENCES

1. F. Mishna, M. Khoury-Kassabri, T. Gadalla, and J. Daciuk, ``Risk factorsfor involvement in cyber bullying: Victims, bullies and bully_victims,''Children Youth Services Rev., vol. 34, no. 1, pp. 63_70, Jan. 2012, doi:10.1016/j.childyouth.2011.08.032.

2. K. Miller, ``Cyberbullying and its consequences: How cyberbullying iscontorting the minds of victims and bullies alike, and the law's limitedavailable redress,'' Southern California Interdiscipl. Law J., vol. 26, no. 2,p. 379, 2016.

3. A. M. Vivolo-Kantor, B. N. Martell, K. M. Holland, and R. Westby,``A systematic review and content analysis of bullying and cyber-bullyingmeasurement strategies,'' Aggression Violent Behav., vol. 19, no. 4,pp. 423_434, Jul. 2014, doi: 10.1016/j.avb.2014.06.008.

4. H. Sampasa-Kanyinga, P. Roumeliotis, and H. Xu, ``Associations betweencyberbullying and school bullying

victimization and suicidal ideation,plans and attempts among Canadian schoolchildren,'' PLoS ONE, vol. 9,no. 7, Jul. 2014, Art. no. e102145, doi: 10.1371/journal.pone.0102145.

5. M. Dadvar, D. Trieschnigg, R. Ordelman, and F. de Jong, ``Improvingcyberbullying detection with user context,'' in Proc. Eur. Conf. Inf. Retr.,in Lecture Notes in Computer Science: Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics, vol. 7814, 2013, pp. 693_696.

6. A. S. Srinath, H. Johnson, G. G. Dagher, and M. Long, ``BullyNet:Unmasking cyberbullies on social networks,'' IEEE Trans.Computat. Social Syst., vol. 8, no. 2, pp. 332_344, Apr. 2021, doi: 10.1109/TCSS.2021.3049232.

7. A. Agarwal, A. S. Chivukula, M. H. Bhuyan, T. Jan, B. Narayan,and M. Prasad, ``Identi_cation and classi_cation of cyberbullying posts:A recurrent neural network approach using under-sampling and classweighting,'' in Neural Information Processing (Communications inComputer and Information Science), vol. 1333, H. Yang, K. Pasupa,A. C.-S. Leung, J. T. Kwok, J. H. Chan, and I. King, Eds. Cham,Switzerland: Springer, 2020, pp. 113_120.

8. Z. L. Chia, M. Ptaszynski, F. Masui, G. Leliwa, and M. Wroczynski,``Machine learning and feature engineering-based study into sarcasm and irony classification with application to cyberbullying detection,''Inf. Process. Manage., vol. 58, no. 4, Jul. 2021, Art. no. 102600, doi:10.1016/j.ipm.2021.102600.

9. N. Yuvaraj, K. Srihari, G. Dhiman, K. Somasundaram, A. Sharma,S. Rajeskannan, M. Soni, G. S. Gaba, M. A. AlZain, and M. Masud,``Nature-inspired-based approach for automated

cyberbullying classificationon multimedia social networking,'' Math. Problems Eng., vol. 2021, pp. 1_12, Feb. 2021, doi: 10.1155/2021/6644652.

10. B. A. Talpur and D. O'Sullivan, ``Multi-class imbalance in text classification: A feature engineering approach to detect cyberbullying inTwitter,'' Informatics, vol. 7, no. 4, p. 52, Nov. 2020, doi: 10.3390/informatics7040052.

11. A. Muneer and S. M. Fati, ``A comparative analysis of machine learningtechniques for cyberbullying detection on Twitter,'' Futur. Internet, vol. 12,no. 11, pp. 1_21, 2020, doi: 10.3390/_12110187.

12. R. R. Dalvi, S. B. Chavan, and A. Halbe, ``Detecting a Twitter cyberbullyingusing machine learning,'' Ann. Romanian Soc. Cell Biol., vol. 25, no. 4,pp. 16307_16315, 2021.

13. R. Zhao, A. Zhou, and K. Mao, ``Automatic detection of cyberbullying onsocial networks based on bullying features,'' in Proc. 17th Int. Conf. Dis-trib. Comput. Netw., Jan. 2016, pp. 1_6, doi: 10.1145/2833312.2849567.

14. L. Cheng, J. Li, Y. N. Silva, D. L. Hall, and H. Liu, ``XBully:Cyberbullying detection within a multi-modal context,'' in Proc. 12thACM Int. Conf. Web Search Data Mining, Jan. 2019, pp. 339_347, doi:10.1145/3289600.3291037.

15. K. Reynolds, A. Kontostathis, and L. Edwards, ``Using machine learningto detect cyberbullying,'' in Proc. 10th Int. Conf. Mach. Learn.Appl. Workshops (ICMLA), vol. 2, Dec. 2011, pp. 241_244, doi:10.1109/ICMLA.2011.152.

16. S. Agrawal and A. Awekar, ``Deep learning for detecting cyberbullyingacross multiple social media platforms,'' in Advances in Information Retrieval (Lecture Notes in Computer Science), vol. 10772,G.

Pasi, B. Piwowarski, L. Azzopardi, and A. Hanbury, Eds. Cham,Switzerland: Springer, 2018, pp. 141_153.

17. R. I. Ra_q, H. Hosseinmardi, R. Han, Q. Lv, S. Mishra, and S. A. Mattson,``Careful what you share in six seconds: Detecting cyberbullying instancesin vine," in Proc. IEEE/ACM Int. Conf. Adv. Social Netw. Anal. Mining(ASONAM), Aug. 2015, pp. 617_622, doi: 10.1145/2808797.2809381.

18. N. Yuvaraj, V. Chang, B. Gobinathan, A. Pinagapani, S. Kannan,G. Dhiman, and A. R. Rajan, ``Automatic detection of cyberbullying usingmulti-feature based arti_cial intelligence with deep decision treeclassi-cation," Comput. Electr. Eng., vol. 92, Jun. 2021, Art. no. 107186, doi:10.1016/j.compeleceng.2021.107186.

19. A. Al-Hassan and H. Al-Dossari, ``Detection of hate speech in Arabictweets using deep learning," Multimedia Syst., Jan. 2021, doi:10.1007/s00530-020-00742-w.

20. Y. Fang, S. Yang, B. Zhao, and C. Huang, ``Cyberbullying detection insocial networks using bi-GRU with self-attention mechanism," Informa-tion, vol. 12, no. 4, p. 171, Apr. 2021, doi: 10.3390/info12040171.

21. C. Iwendi, G. Srivastava, S. Khan, and P. K. R. Maddikunta, ``Cyberbullyingdetection solutions based on deep learning architectures," MultimediaSyst., 2020, doi: 10.1007/s00530-020-00701-5.

22. B. A. H. Murshed, H. D. E. Al-ariki, and S. Mallappa, ``Semantic analysistechniques using Twitter datasets on big data?: Comparative analysisstudy," Comput. Syst. Sci. Eng., vol. 35, no. 6, pp. 495_512, 2020, doi:10.32604/csse.2020.35.495.

23. P. Galán-García, J. G. De La Puerta, C. L. Gómez, I. Santos, andP. G. Bringas, ``Supervised machine learning for the detection of troll pro-_les in Twitter social network: Application to a real case of cyberbullying,"Logic J. IGPL. vol. 24, no. 1, pp. 42_53, 2015, doi: 10.1093/jigpal/jzv048.

24. Y. Zhang and A. Ramesh, ``Fine-grained analysis of cyberbullyingusing weakly-supervised topic models," in Proc. IEEE 5th Int.Conf. Data Sci. Adv. Anal. (DSAA), Oct. 2018, pp. 504_513, doi:10.1109/DSAA.2018.00065.

25. Z. Chen, A. Mukherjee, B. Liu, M. Hsu, M. Castellanos, and R. Ghosh,``Leveraging multi-domain prior knowledge in topic models," in Proc.23rd Int. Jt. Conf. Artif. Intell. Int. Jt. Conf. Artif. Intell. (IJCAI), vol. 13,2013, pp. 2071_2077..