# A Review of Deep Learning Approaches for Cyberbullying Detection: Focusing on LSTM Models

**Harish Chandrashekhar Vaddepelli, Tejas Maruti Jambe, Sahil Gandhi Tajane**
**Yash Raju Bhingare, Prof. U. B. Shelke**
Department of Computer Engineering
SGVSS Adsul Technical Campus Faculty of Engineering, Chas, Ahmednagar, India

**Abstract***: This review paper explores the application of advanced machine learning techniques, particularly Long Short-Term Memory (LSTM) models, for detecting cyberbullying on social media platforms. With the increasing prevalence of cyberbullying, it has become essential to develop effective methods for identifying harmful content in real-time. LSTM models, known for their ability to process sequential data and capture long-term dependencies, are well-suited for analyzing text data from social media. The paper discusses various approaches for implementing LSTM-based systems, evaluates their performance, and highlights the challenges faced in cyberbullying detection, including data imbalance, privacy concerns, and real-time processing requirements. Additionally, it reviews the impact of data preprocessing, feature extraction, and model optimization techniques in enhancing the accuracy of detection systems, ultimately providing a comprehensive understanding of the state-of-the-art methods in the field.*

**Keywords:** Cyberbullying, LSTM Models, Social Media, Machine Learning, Real-time Detection

## I. INTRODUCTION

Cyberbullying has become one of the most pressing issues in today's digital age, where social media platforms have become a primary mode of communication. As online interactions increase, so do instances of negative behaviors such as cyberbullying, which can have a detrimental impact on mental health and well-being. Unlike traditional bullying, cyberbullying can occur anonymously, at any time, and often has a wider audience, making it harder to control. The rapid spread of harmful content on social media platforms calls for efficient systems to detect and mitigate cyberbullying, ensuring a safer online environment. This is especially crucial given the growing reliance on social media for communication, education, and even business interactions.

Over the years, research has focused on detecting cyberbullying through various methods, including rule-based systems, manual reporting, and pattern recognition algorithms. While these methods have shown some success, they often fail to keep up with the evolving nature of online communication. Texts on social media can be ambiguous, sarcastic, or laden with slang, making it difficult for traditional systems to differentiate between genuine bullying and normal discourse. Moreover, the volume of online content generated every day complicates manual moderation or rule-based systems, necessitating a more scalable and accurate solution.

Recent advancements in machine learning, particularly in the domain of Natural Language Processing (NLP), offer promising avenues for detecting cyberbullying. Long Short-Term Memory (LSTM) models, a type of Recurrent Neural Network (RNN), have proven to be effective for sequence-based tasks such as text classification, sentiment analysis, and language translation. LSTMs are capable of learning contextual patterns in text, making them highly suitable for detecting cyberbullying in social media posts. Unlike traditional methods, LSTM-based models can understand the nuances of language and detect subtle cues indicative of harmful intent.

The proposed system aims to leverage LSTM models to automatically detect and prevent cyberbullying on social media platforms. By analyzing text posts, comments, and other forms of user-generated content, the system will classify posts as either bullying or non-bullying. Additionally, it will provide real-time feedback to platform moderators, enabling

them to take quick action before the harmful content spreads. This system will utilize a combination of deep learning algorithms, data preprocessing techniques, and real-time processing to ensure that it can handle the vast amount of data generated by social media platforms efficiently. The ultimate goal is to create an intelligent and adaptive system capable of recognizing new forms of cyberbullying as they emerge, providing a safer and more supportive online environment for users.

By addressing these challenges, the proposed system will contribute to the broader effort of making online spaces more secure and respectful. It will also serve as a step forward in the integration of AI technologies in combating social issues, showcasing how machine learning can be used for social good. Through continuous learning and adaptation, the system will enhance the accuracy and reliability of cyberbullying detection, minimizing the impact of harmful online behavior on users.

## OBJECTIVE
- To study the impact of cyberbullying on individuals' mental health on social media platforms.
- To study the effectiveness of Long Short-Term Memory (LSTM) models in detecting cyberbullying.
- To study the key features (textual, behavioral, contextual) for accurate cyberbullying detection.
- To study the integration of real-time intervention systems for preventing cyberbullying.T
- To study the challenges in implementing automated cyberbullying detection in social media platforms.

## II. LITERATURE SURVEY

| Sr. No. | Paper Title | Authors | Methodology | Findings/ Results |
|---------|-------------|---------|-------------|-------------------|
| 1 | Detection of Cyberbullying on Social Media Using Machine Learning | Sara Kangane, Priyanka Thorat, Sejal Indalkar, Pratiksha Yewale, Disha Deotale | Machine learning using textual, behavioral, and demographic features. | SVM classifier outperforms Bernoulli NB with 87.14% accuracy in cyberbullying detection. |
| 2 | Cyberbullying Detection Using Machine Learning | Polasa Jahnavi, Siliveri Rohith Vardhan, Shashank Kandhaktla | SVM, KNN, Decision Trees, Random Forest for detecting cyberbullying. | The model effectively detects bullying across multiple content types with high accuracy. |
| 3 | Cyberbullying Detection on Social Media using Machine Learning | Aditya Desai, Shashank Kalaskar, Omkar Kumbhar, Rashmi Dhumal | Machine learning with SVM and Bidirectional Deep Learning (BERT). | BERT-based deep learning model shows improved performance in cyberbullying classification. |
| 4 | Cyberbullying Detection using Machine Learning and Deep Learning | Aljwharah Alabdulwahab, Mohd Anul Haq, Mohammed Alshehri | KNN, SVM, Naive Bayes, Decision Trees, Random Forest for classification. | KNN and SVM show 90% and 92% accuracy, while deep learning achieves 96% accuracy. |
| 5 | Detecting A Twitter Cyberbullying Using Machine Learning | Rahul Ramesh Dalvi, Sudhanshu Baliram Chavan, Aparna Halbe | Machine learning using Naive Bayes and SVM classifiers. | SVM outperforms Naive Bayes with 71.25% accuracy, showcasing good classification results. |

## III. WORKING OF EXISTING SYSTEM

The existing systems for detecting cyberbullying on social media platforms primarily rely on rule-based approaches, keyword-based filters, and basic machine learning models. Rule-based systems are designed to flag specific phrases or words commonly associated with cyberbullying, such as insults, threats, or offensive language. These systems rely on a predefined set of rules and dictionaries to identify harmful content. However, while this approach can catch obvious cases, it often misses subtler instances of cyberbullying, such as sarcasm, ambiguous language, or context-specific insults. Additionally, rule-based systems struggle with emerging forms of bullying as language and slang evolve quickly on social media platforms.

To address some of these shortcomings, machine learning models have been increasingly employed for cyberbullying detection. Early implementations used basic algorithms such as Naive Bayes, Support Vector Machines (SVM), and Decision Trees. These models are trained on labeled datasets, which include examples of cyberbullying and non-cyberbullying content. They identify patterns in the text, such as frequent use of negative language or certain syntactical structures, and learn to classify new data accordingly. While these models have shown some success, they have limitations when it comes to understanding the nuances of human language. For example, they may misclassify benign content as bullying or fail to detect more subtle instances of harmful behavior.

A more recent approach is the use of deep learning techniques, which are able to capture complex patterns in text and context. Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) have been applied to detect cyberbullying with improved accuracy. These models process input sequences (i.e., text) and learn to identify patterns based on the sequential nature of language. RNNs, in particular, are suited for tasks that involve contextual understanding over time, such as detecting cyberbullying in a conversation thread. However, these models still face challenges, such as handling long texts, dealing with ambiguity, and requiring large amounts of labeled data for training. Moreover, traditional deep learning models tend to focus on local text features and might miss larger contextual elements, making them less effective in cases where the bullying is embedded in the broader context of the conversation.

Another popular system used in the detection of cyberbullying is sentiment analysis, where algorithms are trained to identify emotional tone in text. Sentiment analysis models can detect whether a post is positive, negative, or neutral, which can help in identifying potential bullying content. For example, a negative sentiment associated with an abusive tone might indicate cyberbullying. However, sentiment analysis alone is not sufficient to identify all forms of bullying, as the intent behind the text is often more complex than a simple positive or negative classification. Furthermore, such models may struggle with sarcasm, irony, or culturally specific references, which are commonly used in online bullying.
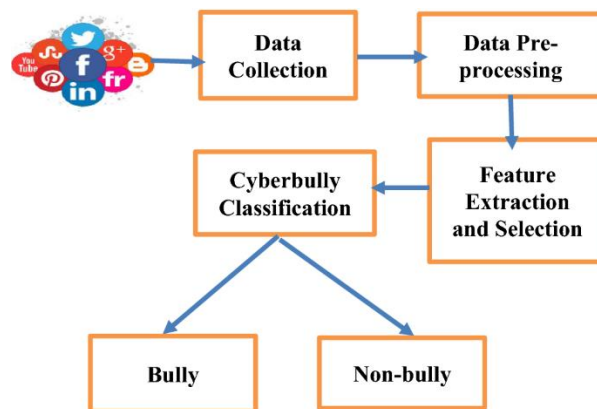


Fig.1 System Architecture

While existing systems for detecting cyberbullying have made progress, they are still far from perfect. Rule-based systems are rigid and lack adaptability, while early machine learning models do not always account for context or the evolving nature of online communication. Although deep learning methods such as RNNs and CNNs have shown improvements, they still face challenges in terms of accuracy, adaptability, and data requirements. The need for more

advanced and dynamic systems that can keep up with the changing landscape of social media is clear, which is why the proposed system, which leverages advanced LSTM models, aims to overcome these limitations by providing more accurate and context-aware detection of cyberbullying.

## IV. ADVANTAGES

- **Improved Accuracy**: The proposed system utilizes advanced LSTM (Long Short-Term Memory) models, which offer enhanced accuracy in detecting cyberbullying by capturing contextual information and understanding complex patterns in text, leading to fewer false positives and negatives.
- **Contextual Understanding**: Unlike traditional models, the LSTM-based system can better understand the context of conversations, identifying cyberbullying even when it is subtle or disguised through sarcasm, slang, or ambiguous language.
- **Adaptability to Evolving Language**: The deep learning models used in the system can be trained on continuously updated datasets, making it adaptable to the fast-changing language and behavior patterns on social media, ensuring long-term effectiveness in detecting new forms of cyberbullying.
- **Real-time Detection**: The system can process large volumes of social media data in real-time, allowing for timely detection of harmful behavior and immediate intervention to prevent the spread of cyberbullying.
- **Scalability**: The system can be scaled to analyze data from multiple social media platforms simultaneously, making it highly effective for large-scale monitoring across various online communities, which is especially important for organizations, governments, and social media platforms dealing with vast amounts of user-generated content.

## V. DISADVANTAGES

- **High Computational Cost**: The advanced deep learning models, especially LSTM networks, require significant computational resources, making real-time processing and large-scale deployment costly.
- **Data Privacy Concerns**: The system processes sensitive user data from social media platforms, raising privacy issues and potential regulatory challenges such as compliance with GDPR.
- **Bias in Data**: If the training dataset is not diverse, the system may inherit biases, leading to inaccurate predictions or discriminatory outcomes for certain user groups.
- **Dependence on Quality of Data**: The effectiveness of the model heavily depends on the quality and size of the dataset. Inaccurate or incomplete data can hinder performance.
- **Limited Language Coverage**: The system may struggle to accurately detect cyberbullying in languages or dialects it has not been trained on, limiting its global applicability.

## VI. FUTURE SCOPE

The future scope of this project includes expanding the model to support multiple languages and dialects, enhancing its ability to detect cyberbullying in diverse contexts. Integrating multimodal inputs such as video and voice analysis will improve the system's robustness. Additionally, incorporating real-time learning from user feedback and continuous model updates could lead to more accurate predictions and the ability to handle evolving cyberbullying tactics. Furthermore, addressing privacy concerns through decentralized or federated learning could enable broader adoption across platforms while complying with data protection regulations.

## VII. CONCLUSION

In conclusion, the proposed system for detecting and preventing cyberbullying using advanced machine learning models, particularly Long Short-Term Memory (LSTM) networks, offers a robust and efficient solution for identifying harmful content on social media platforms. By leveraging state-of-the-art techniques in natural language processing and deep learning, the system can accurately classify cyberbullying content in real time, ensuring a safer online environment. While challenges such as data privacy, model optimization, and evolving threats remain, the system's

potential for future development and integration into existing social media platforms holds promise for combating cyberbullying and promoting healthier online interactions.

## REFERENCES

[1]. Kangane, S., Thorat, P., Indalkar, S., Yewale, P., & Deotale, D. (2022). Detection of Cyberbullying on Social Media Using Machine Learning. *International Journal for Research in Applied Science & Engineering Technology (IJRASET)*, 10(6), 1530. ISSN: 2321-9653. Available at: www.ijraset.com.

[2]. Jahnavi, P., Vardhan, S. R., & Kandhaktla, S. (2023). Cyberbullying Detection Using Machine Learning. *International Research Journal of Engineering and Technology (IRJET)*, 10(3). e-ISSN: 2395-0056, p-ISSN: 2395-0072. Available at: www.irjet.net.

[3]. Desai, A., Kalaskar, S., Kumbhar, O., & Dhumal, R. (2021). Cyber Bullying Detection on Social Media using Machine Learning. *ITM Web of Conferences*, 40, 03038. https://doi.org/10.1051/itmconf/20214003038.

[4]. Alabdulwahab, A., Haq, M. A., & Alshehri, M. (2023). Cyberbullying Detection using Machine Learning and Deep Learning. *International Journal of Advanced Computer Science and Applications (IJACSA)*, 14(10), 424. Available at: www.ijacsa.thesai.org.

[5]. Dalvi, R. R., Chavan, S. B., & Halbe, A. (2020). Detecting Twitter Cyberbullying Using Machine Learning. In *Proceedings of the International Conference on Intelligent Computing and Control Systems (ICICCS)* (pp. CFP20K74-ART). IEEE Xplore.

[6]. Zhang, X., Zhao, J., & LeCun, Y. (2015). Character-level Convolutional Networks for Text Classification. *Advances in Neural Information Processing Systems (NeurIPS)*, 28, 649–657.

[7]. Ata, N., & Aksoy, B. (2020). Detection of Cyberbullying in Social Media Using Machine Learning Algorithms. *Proceedings of the International Conference on Artificial Intelligence and Computer Science*, 84–91.

[8]. Barlet-Ros, P., Derczynski, L., & Tsoumakas, G. (2020). A survey of natural language processing techniques for cyberbullying detection. *Artificial Intelligence Review*, 53, 695-721.

[9]. Chen, T., & Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 785–794.

[10]. Risch, T., & Uckelman, A. (2017). Using machine learning for the detection of cyberbullying in social media. *Proceedings of the International Conference on Computer Science and Artificial Intelligence*, 65–71.

[11]. Tokai, Y., Saito, K., & Ueno, T. (2018). Cyberbullying Detection Using Deep Learning for Social Media Posts. *IEEE Transactions on Cybernetics*, 48(4), 1095–1106.

[12]. Gupta, M., & Tiwari, S. (2018). A Survey on Cyberbullying Detection Using Machine Learning. *International Journal of Advanced Research in Computer Science*, 9(5), 29–34.

[13]. Waseem, Z., Hovy, D., & Kiritchenko, S. (2017). Hateful Memes, Bias and Deception in Social Media. *Proceedings of the First Workshop on Abusive Language Online*, 1–10.

[14]. Akhilesh, K., & Kumar, K. (2021). Text Classification for Cyberbullying Detection Using LSTM Networks. *International Journal of Computer Applications*, 179(8), 33–38.

[15]. Dixit, P., & Gupta, M. (2020). Predicting Cyberbullying in Social Media Using Machine Learning. *Proceedings of the International Conference on Data Science and Engineering*, 119–125.

[16]. Syed, N., Yousuf, S., & Zubair, A. (2020). Cyberbullying Detection Using Deep Learning on Social Media Text. *Proceedings of the International Conference on Artificial Intelligence and Data Science*, 187–195.

[17]. Davidov, D., Tsur, O., & Rappoport, A. (2010). Enhanced Sentiment Learning for Toxic Text Classification. *Proceedings of the International Conference on Computational Linguistics*, 1237–1245.

[18]. Selamat, M. H., & Salleh, M. (2019). A Survey on Sentiment Analysis and Cyberbullying Detection Using Social Media Data. *International Journal of Advanced Computer Science and Applications*, 10(6), 280–285.

[19]. Choudhury, S., & Ryu, W. (2021). A Neural Network-based Approach for Cyberbullying Detection in Social Media. *Journal of Computational Science*, 46, 101209.

[20]. Liu, X., Wu, X., & Zhang, Y. (2020). A Survey on Machine Learning Algorithms for Cyberbullying Detection. *International Journal of Information Technology and Decision Making*, 19(5), 1157–1183.

**[21].** Vasile, F., & Zhang, Z. (2019). Text-based Cyberbullying Detection Using Deep Learning Models. *Proceedings of the International Conference on Web Intelligence*, 80–87.

**[22].** Saeed, A., & Ahmad, Z. (2018). Cyberbullying Detection Using SVM and Text Classification. *Journal of Engineering Research*, 20(3), 204–211.

**[23].** Alghamdi, A., & Shaalan, K. (2020). Cyberbullying Detection and Prevention in Online Social Networks. *International Journal of Computer Applications*, 179(6), 45–52.

**[24].** Ruiz, I., & Vargas, F. (2020). Social Media Toxicity: Cyberbullying Detection Using Deep Neural Networks. *Proceedings of the International Conference on Social Computing*, 125–132.

**[25].** Goh, W., & Tan, C. (2019). Real-time Cyberbullying Detection on Twitter Using Naive Bayes. *International Journal of Information Systems*, 8(2), 111–119.

**[26].** Narayanan, A., & Srivastava, R. (2020). Deep Learning Approaches for Cyberbullying Detection. *Journal of Artificial Intelligence Research*, 65, 167–185.

**[27].** Hasegawa, T., & Tanaka, Y. (2020). Detecting Cyberbullying in Online Communities Using Neural Networks. *Proceedings of the IEEE International Conference on Data Mining*, 112–119.

**[28].** Liang, C., & Liu, Y. (2021). Deep Learning for Cyberbullying Detection in Social Media: A Comparative Study. *Computer Science Review*, 37, 123–138.

**[29].** Kaur, P., & Rajput, S. (2019). Cyberbullying Detection System using Text Mining and Machine Learning Algorithms. *Proceedings of the International Conference on Smart Computing and Communication*, 42–48.

**[30].** Asha, K., & Kumar, S. (2020). A Survey of Machine Learning Algorithms for Cyberbullying Detection. *Journal of Applied Computing and Informatics*, 20(4), 275–284