
INTELLIGENT CYBERBULLYING DETECTION SYSTEM

Snowba J^a, Southra natchiyar M^b, Mr. G. Rahul Vignesh AP/CSE^{c*}

^aDepartment of Computer Science Engineering, Francis Xavier Engineering College,
Tirunelveli, Tamil Nadu – 627003.

^bDepartment of Computer Science Engineering, Francis Xavier Engineering College,
Tirunelveli, Tamil Nadu – 627003.

^cDepartment of Computer Science Engineering, Francis Xavier Engineering College,
Tirunelveli, Tamil Nadu – 627003.

Article Received: 9 February 2026, Article Revised: 01 March 2026, Published on:

***Corresponding Author: Mr. G. Rahul Vignesh**

Department of Computer Science Engineering, Francis Xavier Engineering College, Tirunelveli, Tamil Nadu – 627003.

DOI: <https://doi-doi.org/101555/ijarp.6460>

ABSTRACT

An Intelligent Cyberbullying Detection System is a piece of software designed to identify and prevent cyberbullying on online forums, social networking sites, and messaging applications. The system analyzes text data, detects abusive language, and instantaneously classifies dangerous content using state-of-the-art technologies such as artificial intelligence (AI), machine learning, and natural language processing (NLP). By examining trends, keywords, mood, and user behavior, the system can accurately distinguish between regular conversations and abusive ones. This reduces online abuse, safeguards users, especially children and teenagers, and fosters a safer online environment. By enabling automatic content filtering, generating moderator notifications, and providing early detection, the proposed approach aims to improve online community safety and mental health. The system uses artificial intelligence, machine learning techniques like Naïve Bayes, Support Vector Machines, and Deep Learning, as well as natural language processing techniques to search for offensive or damaging content in text, emojis, hashtags, and even images. By analyzing user behavior patterns, mood, context, and keywords, the system can accurately determine if interactions are harassing or normal in real time.

KEYWORDS: Natural language processing (NLP), text classification, sentiment analysis, deep learning, content moderation, online safety, social media monitoring, artificial intelligence (AI), machine learning (ML), and cyberbullying detection.

INTRODUCTION

The rapid growth of online communication tools, messaging applications, and social media platforms has made cyberbullying a significant issue in today's digital world. Unlike traditional bullying, cyberbullying occurs through electronic devices and may spread quickly, reaching a broad audience quickly. Anxiety, depression, low self-esteem, and mental pain are common among victims of cyberbullying. Internet platforms create enormous amounts of user content every day, which makes it exceedingly difficult to manually monitor and handle dangerous information. Therefore, an automated and intelligent solution is needed to recognize and prevent such harmful activities.

The Intelligent Cyberbullying Detection System uses machine learning, artificial intelligence, and natural language processing techniques to evaluate text data and identify abusive or threatening information in order to address this problem. By automatically classifying and filtering harmful communications, the technology helps administrators effectively reduce occurrences of cyberbullying and promotes a safer online environment. The vast amount of user-generated content that is posted every day makes it difficult for human moderators to manually identify and control cyberbullying activities.

Cyberbullying may be detected and prevented early because to the system's real-time identification of harmful behaviour patterns, harsh language, and negative attitude. The technology helps lessen online abuse and fosters a courteous, safe, and uplifting digital environment by offering automatic monitoring, notifications, and content screening.

MATERIALS AND METHODS

The Intelligent Cyberbullying Detection System uses data processing techniques, software, and hardware. A computer system with sufficient processing power, internet connectivity, and storage is required for the system to handle large datasets. Python, machine learning libraries (NumPy, Pandas, Scikit-learn), and natural language processing libraries are among the software tools used for data preparation and model building. A label dataset comprising both cyberbullying and non-cyberbullying text is collected from public available sources or social media sites.

Preprocessing is applied to the dataset, which includes stemming, tokenization, and the elimination of punctuation, noise, and stop words. Feature extraction techniques such as TF-IDF or word embeddings are used to convert text into numerical representation. Following

that, the data is categorized using the training and testing of machine learning techniques including Deep Learning models, Support Vector Machines, and Naïve Bayes. The trained model is evaluated using performance metrics such as accuracy, precision, recall, and F1-score.

Another aspect of the materials and methods used in the Intelligent Cyberbullying Detection System is a well-structured system architecture that integrates modules for data collection, processing, and categorization. Data is first gathered from chat applications, online forums, and social media sites to provide a balanced dataset of bullying and non-bullying information. Data labeling is done to clearly classify the text, which is a step that is required in supervised learning. To simplify and improve the quality of data, preprocessing techniques including lemmatization, stop-word removal, special character removal, and lowercasing are employed. The system has a real-time detection module that monitors any new messages or postings that come in. When the technology finds cyberbullying content, it alerts users and saves the incident in a secure database for later use.

Table 1. Performance Comparison of Cyber Detection Techniques.

Method	Technique Used	Accuracy (%)	Remarks
Traditional Method	Keyword-Based Filtering	65–70	Easy to use but gives low accuracy and many false results.
Machine Learning	Naïve Bayes	75–80	Fast and simple, but not good at understanding complex language.
Machine Learning	Support Vector Machine (SVM)	82–88	Provides better accuracy and reliable classification of bullying content.
Deep learning	CNN / LSTM	90–94	Accurate real-time Detection.
Proposed System	Hybrid ML + NLP	94–97	Most effective method with high accuracy and real-time detection.

Table 1 presents a comparative analysis of various cyberbullying detection methods is explained in the performance comparison table. Because they simply employ particular words and don't take into account the context of communications, traditional keyword-based approaches have poor accuracy. Although machine learning techniques like Naïve Bayes enhance detection by identifying patterns in data, they are still unable to handle complicated language. Support Vector Machine successfully distinguishes between bullying and non-bullying information, improving classification accuracy and dependability. Although they

need more data and processing power, deep learning models such as CNN and LSTM perform better by comprehending contextual and sequential information in text.

RESULTS AND DISCUSSION

When compared to conventional techniques, the Intelligent Cyberbullying Detection System's results show that sophisticated machine learning and deep learning algorithms greatly increase the accuracy of cyberbullying detection. According to the experimental investigation, keyword-based methods have a lower accuracy and more false positives since they are unable to comprehend context. Learning patterns from labeled data improves the performance of machine learning models like Naïve Bayes and Support Vector Machine (SVM), which produce more accurate and dependable classification results. By incorporating contextual and semantic data from text, deep learning models enhance performance even more, leading to increased recall and accuracy.

Real-Time Cyberbullying Detection Efficiency: The proposed hybrid system, which integrates machine learning and natural language processing techniques, achieves the greatest overall performance in terms of accuracy, precision, and recall. It effectively identifies cyberbullying content in real time and reduces the misclassification of ordinary texts. The system also exhibits scalability and efficiency, making it suitable for application in a realistic online platform. The discussion focuses on how combining a number of astute tactics results in a more comprehensive and reliable cyberbullying detection system, encouraging safer and more secure online communication environments.

Comparison of Detection Techniques: Support Vector Machine and deep learning models perform better than simple machine learning techniques, according to a study of several detection methods. The best accuracy, precision, and recall are attained by the suggested hybrid approach, which blends machine learning and natural language processing methods. This illustrates how combining many clever techniques might result in accurate cyberbullying detection.

Real-Time Alert System: The alert mechanism of the Intelligent Cyberbullying Detection mechanism is crucial for ensuring timely response against harmful material. When the system detects abusive or bullying communications using machine learning and natural language processing (NLP) analysis, it immediately generates an alarm notification. Administrators, moderators, parents, or concerned users may notice these alerts, depending on how the

application is set up. The alert system helps with quick preventative action by alerting the user, banning the content, or temporarily restricting the account. It also stores the found instances in a secure database for further examination and monitoring. By reducing the spread of dangerous content, accelerating reaction times, and providing real-time alerts, the alert system enhances internet safety in general.

In addition to warning the user, the alert system can also alert moderators, parents, or administrators for further action. This enables the timely and appropriate management of significant events. The system may also categorize notifications based on the severity of the material to assist authorities in prioritizing critical situations. Every instance that is discovered is also securely stored in a database for further examination and tracking. This record-keeping helps to improve the detection model and find repeating violations over time. All things considered, the warning system enhances real-time response, inhibits the spread of dangerous material, and encourages the creation of a safer online environment.

CONCLUSION

The Intelligent Cyberbullying Detection System uses cutting-edge technology like machine learning and natural language processing to effectively detect and stop harmful online conduct. In order to stop cyberbullying from spreading, the technology effectively analyzes text content, finds abusive language, and instantly provides notifications. The suggested hybrid strategy outperforms conventional techniques in terms of accuracy, precision, and dependability. The technology contributes to the creation of a more secure and encouraging online environment by automating content monitoring and providing moderators with immediate notifications. All things considered, this experiment shows how sophisticated and automated detection systems are essential for reducing cyberbullying and encouraging appropriate online conversation.

The system can recognize patterns in cyberbullying, evaluate user-generated language, and classify data appropriately. This helps restrict the rapid spread of the unwanted messages on digital channels and reduces the need for human monitoring.

The findings of the proposed system show that advanced and hybrid detection procedures significantly outperform traditional keyword-based methods. The technology's capacity to understand the tone and context of messages improves dependability and lowers false

positives. Customers are protected from continued abuse by a real-time alarm system that guarantees swift response against reported cases of cyberbullying.

The Intelligent Cyberbullying Detection System is crucial to addressing the escalating issues of online harassment in today's digital age. As social media usage continues to rise, there is a growing need for automated and intelligent monitoring solutions. This experiment demonstrates how combining machine learning and natural language processing approaches may significantly improve the detection of harmful and abusive content in comparison to traditional methods.

Apart from precisely identifying instances of cyberbullying, the technology enables real-time monitoring and alarm generation, allowing for timely preventive measures. The proposed model reduces false positives and improves classification performance, which ensures reliable findings. The safe database archiving of recorded incidents also helps with tracking repeat offenders and analyzing behavioral trends for possible improvements. This system might be expanded with advanced deep learning models, language support, and connection with popular social media platforms to increase its utility. With continued training and growth, the detection accuracy might be increased even more. The study highlights the importance of intelligent systems in promoting polite online discourse, protecting users' mental health, and improving digital safety in general.

REFERENCES

1. Dadvar, F. de Jong, R. Ordelman, and D. Trieschnigg, "Improved cyberbullying detection using gender information," *Proceedings of the 12th Dutch-Belgian Information Retrieval Workshop*, 2012.
2. S. Hinduja and J. W. Patchin, "Bullying, cyberbullying, and suicide," *Archives of Suicide Research*, vol. 14, no. 3, pp. 206–221, 2010.
3. Y. Zhao, J. Mao, and J. Wang, "Learning deep features for cyberbullying detection," *International Conference on Web Intelligence*, IEEE, 2016.
4. P. Badjatiya, S. Gupta, M. Gupta, and V. Varma, "Deep learning for hate speech detection in tweets," *Proceedings of the 26th International World Wide Web Conference*, 2017.
5. T. Davidson, D. Warmsley, M. Macy, and I. Weber, "Automated hate speech detection and the problem of offensive language," *Proceedings of ICWSM*, 2017.
6. N. Kumar and S. Sachdeva, "Cyberbullying detection on social media using machine learning," *International Journal of Computer Applications*, vol. 180, no. 25, 2018.

7. Xu, D. Lin, and J. Mao, "Cyberbullying detection based on sentiment analysis," *Journal of Information Security and Applications*, Elsevier, 2019.
8. K. Dinakar, B. Jones, C. Havasi, H. Lieberman, and R. Picard, "Common sense reasoning for detection, prevention, and mitigation of cyberbullying," *ACM Transactions on Interactive Intelligent Systems*, 2012.
9. S. Salawu, Y. He, and J. Lumsden, "Approaches to automated detection of cyberbullying: A survey," *IEEE Transactions on Affective Computing*, 2020.
10. Z. Waseem and D. Hovy, "Hateful symbols or hateful people? Predictive features for hate speech detection on Twitter," *Proceedings of NAACL-HLT*, 2016.
11. K. Reynolds, A. Kontostathis, and L. Edwards, "Using machine learning to detect cyberbullying," *Proceedings of the 10th International Conference on Machine Learning and Applications*, IEEE, 2011.