# CYBER BULLYING DETECTION BY USING MACHINE LEARNING

T Anusha<sup>1</sup> Asst. Professor Department of CSE(DS) TKR College of Engineering & Technology anushathirumani@tkrcet.com V Sai Teja<sup>2</sup> B.Tech(Scholar) Department of CSE(DS) TKR College of Engineering & Technology vadlasaiteja2403@gmail.com

MD Mafaiz Uddin<sup>4</sup> B.Tech(Scholar) Department of CSE(DS) TKR College of Engineering & Technology mafaizuddinm@gmail.com M Suraj Kumar <sup>5</sup> B.Tech(Scholar) Department of CSE(DS) TKR College of Engineering & Technology <u>mnsurajkumar@gmail.com</u> G Manasa<sup>3</sup> B.Tech(Scholar) Department of CSE(DS) TKR College of Engineering & Technology ganganimanasa120gmail.com

### ABSTRACT

Cyberbullying poses a significant threat in the digital age, especially among youth on social media platforms. Traditional approaches to detecting and addressing cyberbullying rely heavily on manual reporting, which is often subjective and inefficient. This study introduces a machine learning-based system to automate the detection cyberbullying, of utilizing advanced natural language processing (NLP) techniques and supervised learning models. The system is trained on diverse datasets featuring text with varying levels of abusive language, sarcasm, and context, ensuring adaptability across different platforms. Preprocessing methods such as tokenization, stemming, and stopword removal are employed to clean and structure the data for analysis. Feature extraction techniques, including TF-IDF and word embeddings, enable the models to capture semantic nuances effectively. Experimental results reveal high detection accuracy in identifying

bullying patterns, with performance validated benchmark against multiple datasets. Future directions include integrating the system with real-time social media monitoring tools and expanding its scope to handle multimedia content such as images and videos. This approach provides a scalable and efficient solution to mitigate the growing issue of cyberbullying and foster safer online environments. With the exponential of online growth communication, cyberbullying has emerged as a critical societal issue, affecting mental health and safety. This study proposes a real-time cyberbullying detection framework leveraging natural language processing (NLP) and deep learning techniques. Using a robust dataset, the system analyzes social media text to identify harmful language architectures patterns. Advanced like bidirectional LSTMs and transformer-based models, such as BERT, are utilized to capture semantic and contextual nuances. Data augmentation and class balancing

techniques address dataset imbalances, enhancing model generalization. Experimental evaluations demonstrate the system's ability to achieve high precision and recall rates, outperforming traditional machine learning approaches. Future work will focus on multilingual capabilities and multimedia content analysis, fostering a safer digital environment.

**KEYWORDS:** Cyberbullying Detection, Natural Language Processing (NLP), Machine Learning, Text Classification, Sentiment Analysis, Deep Learning, BERT, Real-Time Analysis, Abusive Language Identification, Data Preprocessing, Word Embeddings, TF-IDF, Class Imbalance Handling, Data Augmentation, Explainable AI (XAI).

### **1.INTRODUCTION**

Cyberbullying is a growing concern in the digital age, especially with the increasing use of social media platforms, messaging apps, and other online communication channels. Defined as the use of electronic devices or online platforms to harass, harm individuals, embarrass, or cyberbullying can take many forms. including threatening messages, spreading rumors, harassment, impersonation, and exclusion from online groups. Unlike traditional bullying, which typically occurs in physical spaces like schools or workplaces, cyberbullying can occur 24/7 and can reach a global audience, making it difficult for victims to escape or find support. The emotional and psychological toll of cyberbullying can lead to serious

consequences, including depression, anxiety, self-harm, and in extreme cases, suicide.

Given the pervasive nature of cyberbullying and its impact on mental health, it is critical to develop efficient methods for detecting and preventing it. Traditionally, detecting instances of cyberbullying relied on manual reporting by victims or bystanders, which is often insufficient due to the high volume of content generated on social media platforms and the reluctance of individuals to report such incidents. As a result, there is a growing interest in leveraging machine learning (ML) techniques to automatically identify instances of cyberbullying from text, images, and other multimedia content. By using natural language processing (NLP) and sentiment analysis, machine learning algorithms can process large volumes of data and identify harmful interactions more effectively and efficiently.

This paper aims to explore the application of machine learning techniques in the detection of cyberbullying. By using algorithms such as decision trees, support vector machines (SVM), deep learning models, and neural networks, it is possible to create automated that can analyze digital systems conversations and flag instances of cyberbullying for further intervention. The research will explore various models, evaluate their performance, and propose an optimal solution for detecting cyberbullying in real-time.

### **2.RELATED WORK**

Numerous studies have been conducted in the field of cyberbullying detection, with

researchers exploring various machine learning approaches to identify bullying behavior in text, images, and videos. One of the earliest efforts involved the use of keyword-based detection techniques, where researchers manually defined a set of offensive words or phrases and created models to identify these terms in text-based interactions. However, this approach had significant limitations, as it failed to capture the context of the conversation, often leading to high false-positive rates.

In more recent years, researchers have turned to machine learning algorithms to improve the accuracy of cyberbullying detection. For instance, a study by Zhang et al. (2018) utilized a support vector machine (SVM) to classify harmful text from social media platforms. Their model incorporated a range of features, including the frequency of specific words, sentiment, and syntactic patterns. The results showed a significant improvement over keyword-based detection methods, but the model was still limited by the availability of labeled data and the inherent complexity of natural language.

Further advancements in machine learning for cyberbullying detection have been made through the use of deep learning algorithms, particularly recurrent neural networks (RNN) and convolutional neural networks (CNN). For example, in 2019, Wang et al. proposed the use of an RNN-based model to detect cyberbullying in text messages. Their model focused on learning the sequential dependencies in the text, allowing it to better understand the context and detect instances of bullying that involved insults, threats, or other forms of aggression. The study found that deep learning-based models outperformed traditional machine learning models, achieving higher accuracy and reducing false-positive rates.

Another important area of research involves the detection of cyberbullying in multimedia content, such as images and videos. While the majority of cyberbullying research has focused on text-based interactions, there is increasing interest in detecting visual forms cyberbullying, such of as image manipulation, offensive photos, and videos. A study by Goh et al. (2020) explored the use of CNNs to detect harmful content in images shared on social media platforms. model was able to identify Their manipulated or offensive images, achieving impressive results in terms of both detection accuracy and processing speed.

In addition to text and image analysis, researchers have also looked into combining multiple modalities of data for a more comprehensive cyberbullying detection system. In 2021, Liu et al. proposed a multimodal approach that incorporated both text and image features using a hybrid deep learning model. This model demonstrated an improved ability to detect cyberbullying across different types of content, showcasing the potential of combining various data sources for more accurate detection.

# **3.LITERATURE SURVEY**

The detection of cyberbullying using machine learning has been the subject of extensive research, particularly in the last decade. Several studies have demonstrated the potential of machine learning algorithms

to detect and mitigate cyberbullying, each adopting different approaches to handle the complexity of natural language and the diverse forms of bullying behavior. One of the key challenges in this field is understanding the nuances of language, such as sarcasm, implicit threats, and indirect bullying, which makes it difficult for traditional rule-based systems to effectively classify harmful content.

Researchers have employed a variety of machine learning techniques to address these challenges. In a study by Sriram et al. (2017), a deep learning-based approach using Long Short-Term Memory (LSTM) proposed networks was to identify cyberbullying in Twitter posts. The model was trained on a large dataset of annotated tweets and was able to detect various types of cyberbullying, including name-calling, threats, and exclusion. The LSTM network outperformed traditional methods like decision trees and logistic regression, demonstrating the power of deep learning in this domain.

Another study by Gupta et al. (2018) focused on the application of convolutional neural networks (CNN) to classify bullyingrelated content in social media posts. By treating text data as a sequence of characters rather than words, CNNs were able to capture local patterns in the text and identify potentially harmful content. This approach demonstrated strong performance in detecting subtle forms of bullying, such as passive-aggressive comments or subtle insults. Despite the success of deep learning models, one of the key challenges is the need for large, high-quality datasets to train these models. Many existing datasets suffer from issues such as class imbalance, where the number of instances of cyberbullying is much smaller than the number of nonbullying instances, leading to biased predictions. To address this, recent research focused on data augmentation has techniques, such paraphrasing as or generating synthetic examples of cyberbullying, to create more balanced datasets for training machine learning models.

In a study by Sharma et al. (2019), the authors proposed a hybrid machine learning model that combined supervised learning techniques with unsupervised anomaly detection methods to identify cyberbullying in online text. This approach was able to detect previously unseen forms of bullying and adapt to new types of harmful behavior, making it more flexible than traditional supervised models.

### 4.METHODOLOGY

The methodology for detecting using machine cyberbullying learning involves several steps, including data collection. pre-processing, feature extraction, model training, and evaluation. The first step in the process is data collection, where a large dataset of social media posts, text messages, or other online content is gathered. This dataset is typically annotated with labels indicating whether a post contains cyberbullying content or not.

Next, the collected data undergoes preprocessing to prepare it for feature extraction. This may involve cleaning the data to remove irrelevant information, handling missing values, and normalizing the text by converting it to lowercase, removing stopwords, and applying stemming or lemmatization to reduce words to their base forms.



Once the data is cleaned and prepared, the next step is to extract features that will help the machine learning model distinguish between bullying and non-bullying content. Common features include word frequency, sentiment, syntactic structures, and context. For deep learning models, embeddings such as Word2Vec or GloVe can be used to represent words as vectors, capturing their semantic meaning.

The extracted features are then used to train a machine learning model. Various algorithms can be employed, including decision trees, SVM, random forests, or deep learning models such as CNNs, RNNs, or LSTMs. The model is trained on a labeled dataset, with the goal of learning patterns in the data that can be used to predict whether a new post contains cyberbullying content. After the model is trained, it is evaluated using a separate test set to assess its performance. Evaluation metrics such as accuracy, precision, recall, and F1-score are used to measure the model's ability to correctly identify cyberbullying while minimizing false positives and negatives.

### **5.PROPOSED SYSTEM**

The proposed system for detecting cyberbullying is based on a deep learning approach that leverages a combination of convolutional neural networks (CNN) and long short-term memory (LSTM) networks. This hybrid model is designed to analyze the text data from social media posts, comments, and messages, and identify harmful content with high accuracy.

The system first collects data from various social media platforms and pre-processes it by removing irrelevant information and normalizing the text. The cleaned data is then input into the CNN-LSTM hybrid model, where the CNN component is responsible for extracting local features such as word patterns and phrase structures, while the LSTM component captures the sequential dependencies and context of the text.

The model is trained using a large labeled dataset of social media posts that contain both bullying and non-bullying content. Once the model is trained, it is deployed in a real-time environment, where it continuously monitors social media platforms for new posts and analyzes them for potential cyberbullying content.

In addition to text analysis, the system incorporates a sentiment analysis module that helps determine the emotional tone of a post. Posts with a negative sentiment, such as anger or frustration, are more likely to contain bullying behavior, and the system uses this information to improve the accuracy of its predictions.

# 6.IMPLEMENTATION

The implementation of the proposed system involves several key components, including data collection, data pre-processing, model training, and deployment. The first step in the implementation process is to collect data from



various social media platforms using web scraping techniques or APIs provided by platforms like Twitter, Facebook, or Reddit.

Once the data is collected, it is preprocessed by cleaning the text, removing stop words, and applying stemming or lemmatization. The next step is to extract features from the text, including word frequencies, sentiment scores, and syntactic structures.

The CNN-LSTM model is then implemented using a deep learning framework like TensorFlow or PyTorch. The model is trained on the pre-processed dataset, and its performance is evaluated using a test set of labeled data. Once the model achieves satisfactory performance, it is deployed in a environment. where real-time it monitors social media continuously platforms for new posts and classifies them as either bullying or non-bullying content.

# 7.RESULT AND DISCUSSION

The system is evaluated based on several performance metrics, including accuracy, precision, recall, and F1-score. The results show that the CNN-LSTM hybrid model outperforms traditional machine learning models in detecting cyberbullying in text. The system is able to identify harmful content with high accuracy while minimizing false positives and negatives.

In addition to text-based detection, the sentiment analysis module helps improve the accuracy of the system by considering the emotional tone of posts. Posts with negative sentiment are more likely to contain cyberbullying behavior, and this additional layer of analysis enhances the model's ability to detect harmful content.

datasets.



# **8.CONCLUSION**

In conclusion, the use of machine learning techniques, particularly deep learning models, offers a promising solution for detecting cyberbullying in real-time. The proposed **CNN-LSTM** hybrid model demonstrates performance strong in identifying bullying content in social media posts, and the sentiment analysis module further improves its accuracy. The system provides an automated, scalable solution for addressing the growing problem of cyberbullying, and it has the potential to be deployed across various online platforms to protect individuals from online harassment.

# **9.FUTURE SCOPE**

The future scope of cyberbullying detection using machine learning includes several areas for improvement. One area is the expansion of the system to handle other types of content, such as images and videos, which may also contain harmful or offensive material. Additionally, the integration of multi-modal data, such as combining text, images, and sentiment analysis, could further enhance the accuracy of detection. Finally, there is potential to refine the model by using more advanced natural language processing techniques and exploring the use of generative models for data augmentation to address class imbalance issues in training

Vol.15, Issue No 1, 2025

# **10.REFERENCES**

- Sriram, P., & Lakkaraju, H. (2017). "Cyberbullying detection using deep learning techniques". Journal of Social Media Studies, 5(2), 89-105.
- Gupta, A., & Verma, R. (2018). "Detection of cyberbullying using machine learning on social media platforms". International Journal of Computer Science and Technology, 9(4), 101-114.
- Sharma, R., & Tiwari, S. (2019). "Cyberbullying detection through hybrid machine learning models". Journal of Machine Learning Research, 8(4), 345-360.
- Zhang, X., & Lee, W. (2020). "Cyberbullying detection in social media using deep learning models". Social Computing and Social Media, 10(1), 134-148.
- Wang, J., & Lu, X. (2019). "Detecting cyberbullying in text using recurrent neural networks". Journal of Artificial Intelligence and Social Media, 15(3), 88-100.
- Goh, C., & Chen, D. (2020). "Detecting visual cyberbullying using deep learning on images". Journal of Visual Communication Technology, 14(2), 77-91.

- Liu, Y., & Zhang, S. (2021). "Multimodal cyberbullying detection using deep learning approaches". Journal of Multimedia Technology, 22(4), 122-137.
- Choi, M., & Park, S. (2017). "A survey of machine learning techniques for cyberbullying detection". Proceedings of the International Conference on Computational Intelligence, 13, 241-255.
- Ali, K., & Singh, A. (2020). "Machine learning-based detection of cyberbullying in text messages". Journal of Computer Science and Technology, 14(3), 56-70.
- 10. O'Neill, A., & Mathews, J. (2019).
  "Using convolutional neural networks for social media content analysis in cyberbullying detection". International Journal of Machine Learning, 12(5), 142-158.
- 11. Zhang, M., & Lee, H. (2018).
  "Sentiment analysis-based detection of cyberbullying in social media posts".
  Social Media Analysis Journal, 17(4), 103-115.
- 12. Faris, R., & Smith, K. (2020)."Cyberbullying and machine learning: An empirical study of detection systems". Journal of Online Behavior and Ethics, 7(2), 75-89.
- 13. Gupta, M., & Sharma, S. (2019)."Detecting cyberbullying using hybrid models in online platforms". Journal of Digital Communication, 13(6), 244-258.

- 14. Raj, S., & Patel, A. (2018). "Automatic detection of cyberbullying behavior in text using machine learning". Journal of Cybersecurity and Privacy, 10(1), 29-44.
- Tiwari, R., & Soni, P. (2021). "Social media monitoring for cyberbullying detection using NLP techniques". Journal of Social Media Research, 8(2), 116-130.
- 16. Kumar, V., & Rao, S. (2020). "Analysis of cyberbullying detection using data-driven machine learning techniques". Data Science and Applications, 15(1), 25-39.
- 17. Sharma, V., & Kumar, M. (2020). "A hybrid approach for cyberbullying detection in social networks". International Journal of Advanced Computing, 22(2), 83-97.
- Yang, H., & Chen, G. (2019). "Deep learning for cyberbullying detection in social media: A review". Journal of Computer Vision and Pattern Recognition, 18(3), 50-61.
- Chen, W., & Zhao, L. (2021).
   "Cyberbullying detection in online platforms using machine learning techniques". International Journal of Information Security, 29(4), 121-134.
- 20. Patel, A., & Singh, H. (2020). "Realtime cyberbullying detection using deep neural networks". Journal of Artificial Intelligence Research, 14(2), 76-91.