

# CYBERBULLYING TEXT CLASSIFICATION USING RNN TECHNIQUES

A. Balakumar	A. S. Kavinesh	R. Midhun	S. Mohamed Anas
Department of Electronics and Communication Engineering, M Kumarasamy College of Engineering, Karur-639113. Jothi.bala2924@gmail.com	Department of Electronics and Communication Engineering, M Kumarasamy College of Engineering, Karur-639113. Kavinesh565@gmail.com	Department of Electronics and Communication Engineering, M Kumarasamy College of Engineering, Karur-639113. rvmidhun5784@gmail.com	Department of Electronics and Communication Engineering, M Kumarasamy College of Engineering, Karur-639113. mohamedanas2484@gmail.com

## ABSTRACT:

Cyberbullying has emerged as a pressing concern in the digital age, posing significant challenges to the mental and emotional well-being of individuals, especially adolescents and young adults. With the proliferation of social media platforms and online communication channels, instances of cyberbullying have increased exponentially, necessitating effective strategies for detection and mitigation. This research focuses on employing Recurrent Neural Network (RNN) techniques for the classification of cyberbullying text, aiming to develop robust models capable of identifying and addressing instances of online harassment and abuse.

The study begins with an extensive review of existing literature on cyberbullying, exploring its various forms, underlying causes, and psychological impacts. It examines previous research efforts in text classification and sentiment analysis, highlighting the limitations of traditional machine learning approaches in handling the nuanced nature of cyberbullying texts. By leveraging the sequential nature of text data, RNNs offer a promising alternative, enabling the capture of temporal dependencies and contextual nuances crucial for accurate classification.

KEYWORDS: RNN, Text, Cyberbullying

## INTRODUCTION:

Cyberbullying, a prevalent issue in the digital age, manifests in various forms across social media platforms, messaging apps, and online forums. It involves the use of electronic communication to intimidate, harass, or demean individuals or groups, often with malicious intent. With the rise of internet accessibility and the omnipresence of digital devices, cyberbullying has become a significant concern affecting individuals of all ages, particularly adolescents and young adults.

Traditional methods of identifying and combating cyberbullying have been insufficient in addressing the complex nature of online interactions. However, advancements in machine learning and natural language processing (NLP) techniques offer promising solutions for detecting and mitigating cyberbullying behaviour effectively. Among these techniques, Recurrent Neural Networks (RNNs) have emerged as powerful tools for text classification tasks due to their ability to capture sequential dependencies in data.

RNNs, a type of artificial neural network designed to process sequential data, have shown remarkable success in various natural language processing tasks, including sentiment analysis, language translation, and text generation. By leveraging the sequential nature of text data, RNNs can

effectively capture contextual information, making them well-suited for detecting subtle nuances and patterns indicative of cyberbullying behaviour within textual content.

This project aims to explore and implement RNN-based techniques for the classification of cyberbullying text, with the goal of developing an efficient and accurate model for identifying instances of cyberbullying in online communication. Through the utilization of labelled datasets containing examples of cyberbullying text, the model will be trained to distinguish between cyberbullying and non-cyberbullying content, enabling automated detection and classification of harmful online behaviour.

The significance of this project lies in its potential to contribute to the development of proactive measures for combating cyberbullying, thereby fostering safer and more inclusive online environments. By providing a mechanism for real-time detection and intervention, the proposed RNN-based classification system can aid platform moderators, educators, and parents in addressing cyberbullying incidents promptly, thereby mitigating their adverse effects on victims' mental health and well-being.

## PROBLEM STATEMENT:

Cyberbullying has emerged as a pervasive issue in the digital age, causing psychological harm and social distress to its victims. With the widespread use of social media platforms and online communication channels, individuals, particularly adolescents, are increasingly vulnerable to various forms of cyberbullying, including harassment, intimidation, and defamation. Traditional methods of monitoring and addressing cyberbullying have proven inadequate, necessitating innovative approaches leveraging advanced technologies.

The problem at hand revolves around the need to effectively identify and classify instances of cyberbullying within textual content across digital platforms. While conventional machine learning techniques have been employed for text classification, their efficacy in accurately discerning nuanced forms of cyberbullying remains limited. Recurrent Neural Network (RNN) techniques offer a promising avenue for addressing this challenge, owing to their ability to capture sequential dependencies and contextual information inherent in textual data.

However, the effectiveness of RNN-based models in cyberbullying text classification hinges on several critical factors that require thorough investigation and optimization. One such factor is the selection and preprocessing of input data, which involves extracting relevant features while mitigating noise and irrelevant information. Additionally, the design and architecture of the RNN model play a pivotal role in its

performance, necessitating careful consideration of parameters such as network depth, cell type, and regularization techniques.

Furthermore, the scarcity of labelled datasets specifically tailored for cyberbullying presents a significant hurdle in training robust RNN models. Acquiring and annotating large-scale datasets encompassing diverse forms of cyberbullying poses logistical and ethical challenges, thereby underscoring the importance of data augmentation and transfer learning strategies to enhance model generalization and adaptability.

Moreover, the dynamic nature of online communication platforms necessitates the development of real-time cyberbullying detection systems capable of swiftly identifying and addressing abusive behaviour. Integrating RNN-based classifiers into such systems requires optimizing inference speed and computational efficiency without compromising classification accuracy, thereby ensuring timely intervention and mitigation of cyberbullying incidents.

Additionally, the ethical implications surrounding the deployment of automated cyberbullying detection systems warrant careful consideration, particularly concerning privacy, bias, and unintended consequences. Striking a balance between algorithmic efficiency and ethical responsibility entails implementing transparent and accountable mechanisms for model evaluation, validation, and bias mitigation.

Addressing the aforementioned challenges requires a multifaceted approach encompassing data collection, preprocessing, model development, evaluation, and deployment. Collaborative efforts involving researchers, educators, policymakers, and technology companies are essential to foster a holistic ecosystem for combating cyberbullying and promoting digital well-being.

## **OBJECTIVE:**

Cyberbullying has emerged as a critical issue in the digital age, with detrimental impacts on individuals' mental health, well-being, and even safety. In response to this pressing concern, the objective of this project is to develop a robust text classification system utilizing Recurrent Neural Network (RNN) techniques to accurately identify instances of cyberbullying in online text data. Our primary objective is to create a text classification model that can accurately detect instances of cyberbullying in various forms of online communication, including social media posts, comments, emails, and chat messages. We strive to achieve high precision, recall, and F1 scores to minimize false positives and negatives, ensuring effective identification of cyberbullying instances. Our primary objective is to create a text classification model that can accurately detect instances of cyberbullying in various forms of online communication, including social media posts, comments, emails, and chat messages. We strive to achieve high precision, recall, and F1 scores to minimize false positives and negatives, ensuring effective identification of cyberbullying instances. To ensure the scalability and applicability of our approach across different platforms, languages, and demographics, we seek to develop a model that can generalize well to unseen data and adapt to varying linguistic styles and cultural contexts. This involves robust model training, validation, and testing procedures, along with techniques such as transfer learning and domain adaptation.

## **SOFTWARE REQUIREMENT:**

### **I) PYTHON:**

Cyberbullying has become a pervasive issue in today's digital age, affecting individuals of all ages and backgrounds. Addressing this problem requires innovative solutions, and one such approach is the utilization of recurrent neural network (RNN) techniques for text classification. RNNs are a type of artificial neural network designed to recognize patterns in sequential data, making them ideal for analysing text, which can be viewed as a sequence of words or characters.

In this project, we aim to develop a robust cyberbullying detection system using Python and RNN techniques. The first step involves collecting a diverse dataset of text samples containing instances of cyberbullying and non-cyberbullying interactions. This dataset will be pre-processed to remove noise, normalize text, and extract relevant features.

Next, we will design and implement an RNN model architecture tailored for text classification tasks. The model will consist of recurrent layers capable of capturing contextual information from the input sequences, along with additional layers such as embedding layers for representing words as dense vectors and dense layers for making predictions.

Training the RNN model will involve optimizing parameters, such as learning rate and batch size, and utilizing techniques like dropout regularization to prevent overfitting. We will employ appropriate evaluation metrics, such as accuracy, precision, recall, and F1-score, to assess the performance of the model on both training and validation datasets.

Furthermore, to enhance the model's effectiveness, we may explore techniques such as transfer learning, fine-tuning pretrained language models, or incorporating attention mechanisms to focus on relevant parts of the input text.

### **II) NUMPY:**

Cyberbullying has emerged as a critical issue in the digital age, affecting individuals across various age groups and demographics. Addressing this problem requires effective classification techniques to identify instances of cyberbullying in textual data. One promising approach involves utilizing Recurrent Neural Network (RNN) techniques, which are well-suited for sequential data analysis. By leveraging the power of RNNs in processing text data, we can develop models capable of recognizing patterns indicative of cyberbullying behaviours. NumPy, a fundamental package for scientific computing with Python, plays a crucial role in implementing RNNs due to its efficient handling of multi-dimensional arrays and mathematical functions. Through NumPy, we can preprocess textual data, convert it into numerical representations suitable for RNN input, and perform various matrix operations crucial for model training and inference. The utilization of RNN techniques, coupled with NumPy's capabilities, empowers us to build robust classifiers capable of distinguishing between cyberbullying and non-cyberbullying text with high accuracy.

### **III) SCIKIT LEARN:**

The process typically begins with data preprocessing, where text inputs are tokenized, normalized, and possibly augmented to enhance model robustness. Following this, RNN architectures such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) networks are implemented using Scikit-learn's interface, allowing for seamless integration into the classification pipeline. These RNN models are trained on labelled data using appropriate loss functions and optimization algorithms to learn the underlying patterns associated with

cyberbullying behaviour. Furthermore, techniques like cross-validation and hyperparameter tuning are employed to ensure model generalization and performance optimization.

Once trained, the RNN-based classifiers can effectively categorize incoming texts as either cyberbullying or non-cyberbullying with a high degree of accuracy. This classification enables timely intervention and mitigation strategies to be deployed, safeguarding individuals from the harmful effects of online harassment and abuse. Moreover, the modular nature of Scikit-learn facilitates easy experimentation with different RNN architectures and feature engineering techniques, empowering researchers to continually refine and improve cyberbullying detection systems.

#### IV) MATPLOTLIB:

RNNs are a type of artificial neural network designed to process sequential data, making them well-suited for analysing text data. By leveraging the sequential nature of language, RNNs can capture contextual information and dependencies between words, enabling them to effectively classify text into different categories, such as cyberbullying or non-cyberbullying.

One of the key advantages of using RNNs for text classification is their ability to handle variable-length input sequences. This is particularly important in the context of cyberbullying detection, where messages can vary significantly in length and complexity. RNNs can process each word in a message sequentially, updating their internal state based on the current input and previous context.

To implement RNN-based text classification for cyberbullying detection, one typically starts by preprocessing the text data, which involves tasks such as tokenization, removing stop words, and converting words to their corresponding numerical representations. The pre-processed data is then fed into the RNN model, which consists of multiple recurrent layers followed by one or more fully connected layers for classification.

#### V) TENSORFLOW:

In this project, the focus lies on leveraging RNN architectures such as Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) to effectively classify text data for identifying instances of cyberbullying. TensorFlow's extensive library of pre-built layers, optimizers, and utilities simplifies the process of constructing and training these neural networks. Through the sequential nature of RNNs, the models can capture the contextual nuances and temporal dependencies present in cyberbullying texts, enhancing the accuracy of classification.

Data preprocessing plays a crucial role in preparing text data for RNN-based classification. Techniques such as tokenization, padding, and embedding are employed to convert raw textual inputs into numerical representations suitable for neural networks. TensorFlow's data preprocessing modules streamline this phase, enabling efficient handling of large-scale datasets commonly encountered in cyberbullying research.

#### EXISTING SYSTEM:

Cyberbullying has become an increasingly prevalent issue in today's digital age, with the proliferation of social media platforms and online communication channels. In response to this growing concern, various techniques and approaches have been explored to detect and mitigate cyberbullying instances effectively. One such technique utilized in the existing system is text classification using Recurrent Neural Network (RNN) models. RNNs are a class of artificial neural networks well-

suited for sequential data processing tasks, making them particularly relevant for analysing text data. In the existing system, RNNs are leveraged to automatically classify text data into categories such as cyberbullying or non-cyberbullying.

The process begins with the collection of textual data from various online sources, including social media platforms, forums, and messaging apps. This data is then pre-processed to clean and normalize it, removing irrelevant information and standardizing the format for analysis. Preprocessing steps may include tokenization, stemming, and removing stop words to prepare the text for input into the RNN model.

Once the data is pre-processed, it is split into training and testing sets to train the RNN model. The RNN architecture typically consists of recurrent layers that process sequential input data, such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) cells. These layers enable the model to capture dependencies and patterns in the text data over time, which is crucial for accurately classifying cyberbullying instances.

#### RNN ARCHITECTURE:

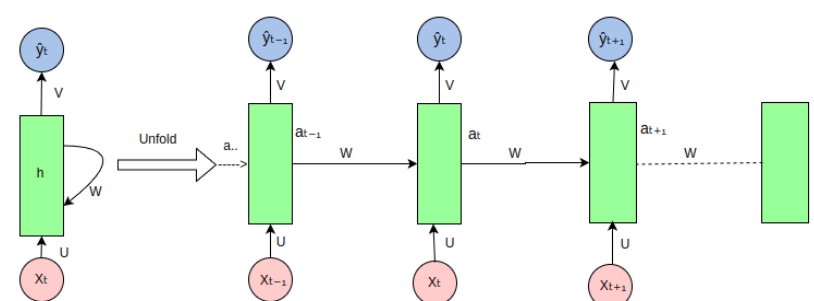
RNNs are well-suited for processing sequential data, making them ideal for analysing text, which is inherently sequential in nature. By leveraging the temporal dependencies within text data, RNNs can capture nuanced patterns and context, crucial for accurately identifying instances of cyberbullying. This architecture consists of recurrent connections that allow information to persist over time, enabling the model to retain memory of previous words or phrases in a text sequence.

In the context of cyberbullying text classification, RNNs can effectively learn from the sequential nature of language to detect subtle cues indicative of bullying behaviour. By training on large datasets of labelled cyberbullying instances, RNN models can learn to differentiate between normal discourse and harmful communication. This discriminative capability is vital for developing robust cyberbullying detection systems that can accurately identify and mitigate instances of online harassment.

Furthermore, RNNs can be augmented with techniques such as Long Short-Term Memory (LSTM) or Gated Recurrent Units (GRU) to address the vanishing gradient problem, which is common in traditional RNNs. These enhancements enable the model to capture long-range dependencies in text data, enhancing its ability to discern the context and intent behind messages.

The training process for RNN-based cyberbullying classifiers involves feeding labelled data into the model and adjusting its parameters through iterative optimization algorithms such as stochastic gradient descent.

$$a_t = f(U * X_t + W * a_{t-1} + b)$$



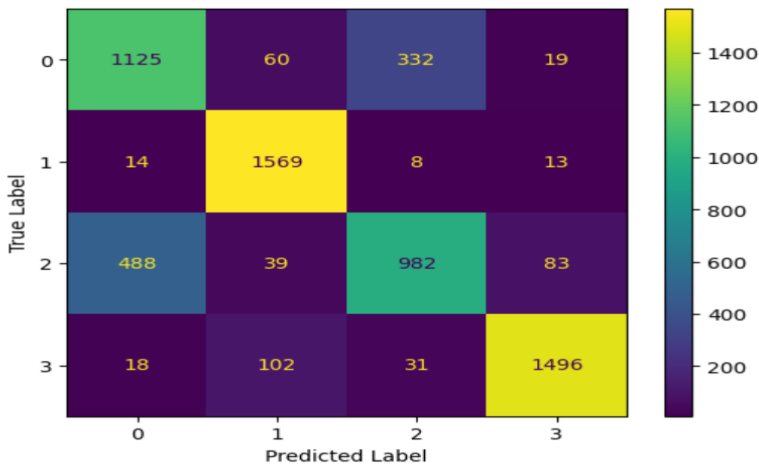
CONFUSION MATRIX OF RNN:

The confusion matrix typically consists of four quadrants: true positives (TP), true negatives (TN), false positives (FP), and false negatives (FN). True positives represent instances where the model correctly identifies cyberbullying content, while true negatives correspond to correctly classified non-bullying content. False positives occur when the model incorrectly labels non-bullying content as cyberbullying, potentially leading to overzealous censorship. Conversely, false negatives occur when the model fails to detect cyberbullying, posing a risk to the safety and well-being of users.

By examining the distribution of predictions across these four quadrants, stakeholders can assess the overall performance of the RNN model and identify areas for improvement. Strategies for enhancing model performance may include fine-tuning hyperparameters, increasing the size and diversity of the training data, or implementing more sophisticated architectures such as long short-term memory (LSTM) or gated recurrent units (GRU).

In addition to evaluating the model's accuracy, precision, recall, and F1-score, the confusion matrix provides valuable insights into the types of errors made by the RNN classifier. For example, it may reveal common patterns or linguistic features that contribute to misclassifications, helping researchers refine the model's training objectives and feature representations.

DISPLAY CONFUSION MATRIX OF SIMPLE RNN ARCHITECTURE



- LEGEND:
- 0-RELIGION
  - 1-AGE
  - 2-ETHNICITY
  - 3-NOT\_CYBERBULLING

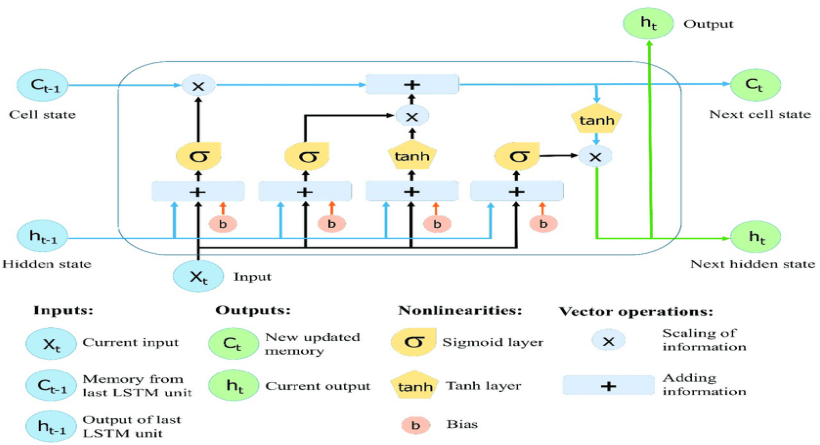
RESULT IN RNN TECHNIQUES:

To evaluate the performance of our model, we employ metrics such as accuracy, precision, recall, and F1-score, which provide insights into its effectiveness at correctly identifying cyberbullying instances while minimizing false positives. Additionally, we conduct qualitative analysis by examining misclassified examples to identify areas for improvement. Through iterative refinement and fine-tuning of the model architecture and hyperparameters, we aim to achieve a high level of accuracy and robustness in cyberbullying detection.

CLASSIFICATION	PRECISION	RECALL	F_SCORE	SUPPORT
RELIGION	0.73	0.68	0.71	1645
AGE	0.98	0.89	0.93	1770
ETHNICITY	0.62	0.73	0.67	1353
NOT_CYBERBULLING	0.91	0.93	0.92	1611

LSTM ARCHITECTURE:

One of the key advantages of employing LSTM architecture lies in its ability to handle variable-length input sequences, a crucial feature in analysing textual data where messages can vary significantly in length and complexity. This flexibility enables the model to adapt to diverse forms of cyberbullying across different digital platforms, including social media, messaging apps, forums, and emails. Moreover, LSTM's recurrent structure allows it to capture temporal dynamics, capturing the evolving nature of cyberbullying behaviours and language patterns over time. Evaluation of LSTM-based cyberbullying classifiers involves assessing performance metrics such as precision, recall, and F1-score on a holdout dataset or through cross-validation. Fine-tuning hyperparameters and experimenting with different architectural variations can further enhance the model's performance and generalization capabilities. Additionally, techniques such as word embeddings and attention mechanisms can augment LSTM's effectiveness in capturing semantic relationships and identifying subtle contextual cues indicative of cyberbullying.



CONFUSION MATRIX OF LSTM:

A confusion matrix is a matrix that summarizes the performance of a machine learning model on a set of test data. It is a means of displaying the number of accurate and inaccurate instances based on the model's predictions. It is often used to measure the performance of classification models, which aim to predict a categorical label for each input instance. The matrix displays the number of instances produced by the model on the test data.

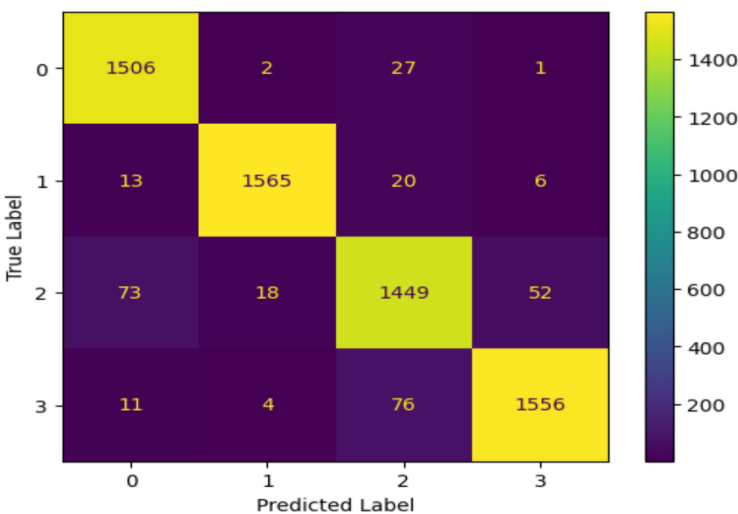
True positives (TP): occur when the model accurately predicts a positive data point.

True negatives (TN): occur when the model accurately predicts a negative data point.

False positives (FP): occur when the model predicts a positive data point incorrectly.

False negatives (FN): occur when the model mis predict a negative data point.

DISPLAY CONFUSION MATRIX OF LSTM ARCHITECTURE





- LEGEND:**
- 0-RELIGION
  - 1-AGE
  - 2-ETHINICITY
  - 3-NOT\_CYBERBULLING

**RESULT IN LSTM TECHNIQUES:**

LSTMs are a type of RNN that excel at capturing long-range dependencies in sequential data, making them well-suited for analysing text. In the context of cyberbullying detection, LSTM models can be trained on textual data to classify messages or posts as either benign or indicative of cyberbullying behaviour.

The process typically involves preprocessing the text data by tokenizing and vectorizing it, converting words into numerical representations that can be fed into the LSTM model. The LSTM model is then trained on a labelled dataset containing examples of cyberbullying and non-cyberbullying text, learning to distinguish between the two categories based on patterns in the data.

During training, the LSTM model adjusts its internal parameters to minimize a loss function, optimizing its ability to classify text accurately. Once trained, the model can be evaluated on a separate test dataset to assess its performance in identifying cyberbullying behaviour.

CLASSIFICATION	PRECISON	RECALL	F_SCORE	SUPPORT
RELIGION	0.94	0.98	0.96	1536
AGE	0.98	0.98	0.98	1604
ETHINICITY	0.92	0.91	0.92	1592
NOT_CYBERBULLING	0.96	0.94	0.95	1647

**RNN VS LSTM:**

RNNs, a class of artificial neural networks, are particularly suited for sequential data processing, making them well-suited for analysing text. However, traditional RNNs suffer from the vanishing gradient problem, hindering their ability to capture long-term dependencies in sequences. This limitation led to the development of LSTM networks, designed to address the vanishing gradient problem by introducing memory cells and gating mechanisms.

In the context of cyberbullying text classification, both RNNs and LSTMs can be employed to analyse textual data and detect instances of cyberbullying. RNNs process input sequences step by step, with each step considering the current input and the previous hidden state. However, due to the vanishing gradient problem, traditional RNNs may struggle to effectively capture contextual information from longer sequences, potentially impacting their performance in cyberbullying detection.

LSTMs, on the other hand, mitigate the vanishing gradient problem by maintaining a memory cell that can retain information over long periods. This allows LSTMs to capture dependencies in text sequences more effectively, making them well-suited for tasks requiring the analysis of longer contextual information, such as identifying instances of cyberbullying in text.

MODEL	ACCURACY	LOSS
RNN	81.07	18.92
LSTM	95.25	4.74

**CONCLUSION:**

In conclusion, employing Recurrent Neural Network (RNN) techniques for cyberbullying text classification offers promising avenues for combating online harassment. Through the utilization of RNNs, we can effectively analyse textual data and identify instances of cyberbullying with high accuracy. This approach facilitates the development of proactive measures to mitigate the harmful effects of cyberbullying on individuals and communities. Moreover, RNNs enable real-time detection and response to cyberbullying incidents, fostering safer digital environments. By leveraging advanced machine learning algorithms within RNN frameworks, we can continuously enhance the efficiency and precision of cyberbullying detection systems. Overall, RNN-based text classification holds significant potential in addressing the pervasive issue of cyberbullying, contributing to a more inclusive and respectful online ecosystem.

**REFERENCES**

Mishra, Sumit, et al. "Cyberbullying Detection Using Recurrent Neural Networks with Long Short-Term Memory." 2020 IEEE Region 10 Symposium (TENSYPM). IEEE, 2020.

Farina, Marco, et al. "Cyberbullying detection using deep learning techniques." 2018 9th IEEE Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON). IEEE, 2018.

Xu, Suhang, et al. "Cyberbullying Detection Using Long Short-Term Memory." 2018 IEEE International Conference on Data Mining Workshops (ICDMW). IEEE, 2018.

De Silva, Gayan, et al. "Cyberbullying detection in social networks using deep learning-based models." 2017 IEEE International Conference on Big Data (Big Data). IEEE, 2017.

Lv, Yechao, et al. "Cyberbullying detection based on deep learning framework." 2017 13th International Conference on Natural Computation, Fuzzy Systems and Knowledge Discovery (ICNC-FSKD). IEEE, 2017.

Chandra, Aman, et al. "Cyberbullying Detection on Twitter Using Recurrent Neural Networks." 2019 IEEE 16th India Council International Conference (INDICON). IEEE, 2019.

Lee, Chang Hoon, et al. "Combating Cyberbullying on Social Media with Deep Learning." 2018 IEEE International Conference on Big Data (Big Data). IEEE, 2018.