# Character-Level Convolutional Neural Networks for Cyberbullying Detection: A Robust Approach to Handling Noisy Social Media Text

**[1]Kondragunta Rama Krishnaiah, [2]Harish H, [3]Manjunath B E**

**Abstract:** With the increasing prevalence of cyberbullying on social media, there is a pressing need for effective detection methods that can handle the noisy, unstructured nature of online text. Traditional machine learning models often struggle with the informal language, misspellings, and emoticons commonly used in cyberbullying messages. In this paper, we propose a novel approach for detecting cyberbullying using **Character-level Convolutional Neural Networks (Char-CNNs)**. Unlike word-based models, the **Char-CNN** model operates at the character level, allowing it to effectively handle spelling errors, intentional distortions, and the use of emojis. We evaluate the performance of **Char-CNN** on a publicly available social media dataset and compare it with a traditional **Word-CNN** model. Our results show that **Char-CNN** outperforms the word-based approach across key performance metrics, including accuracy, precision, recall, and F-measure. The model's ability to generalize well in the presence of noisy data makes it a promising tool for real-time cyberbullying detection. Furthermore, we discuss the limitations of the current model and future directions for enhancing its performance, particularly in detecting more subtle forms of cyberbullying.

*Keywords:* Cyberbullying, Char-CNN, Convolutional Neural Networks, Text Classification, Social Media Analysis

## 1. INTRODUCTION

The widespread adoption of digital platforms and social media has significantly transformed how people interact and communicate. Unfortunately, this shift has also given rise to a pressing societal issue: **cyberbullying**. Defined as the use of digital means such as the internet and mobile devices to harass or harm others, cyberbullying has far-reaching consequences for victims, including emotional distress, mental health struggles, and even self-harm or suicide in extreme cases [1]. Social media platforms like Facebook, Twitter, and Instagram, which offer a space for free expression and interaction, are often the breeding grounds for such negative behavior.

One of the most challenging aspects of addressing cyberbullying is its detection. Traditional

cyberbullying detection systems typically rely on pre-defined keywords or linguistic cues to identify harmful content. However, these systems often fall short when it comes to understanding the nuanced and diverse ways in which cyberbullying manifests online. Textual content in social media is often short, unstructured, and riddled with spelling errors, emoticons, abbreviations, and intentional distortions aimed at evading detection [2]. This presents a unique challenge for machine learning models, which are typically built on a vocabulary-based approach and often struggle to process this unclean, informal text.

An additional challenge arises from the complexity of detecting cyberbullying in the absence of clearly defined, universally accepted criteria for what constitutes bullying in an online context [3]. Many cases involve subtle, indirect language, such as sarcasm or veiled insults, making them difficult to categorize using traditional methods. In many cases, posts may not explicitly contain offensive language, but the underlying tone or context may still be harmful. Therefore, the need for more sophisticated detection systems that can effectively capture these subtle cues is evident.

[1]*R K College of Engineering (A), Kethanakonda (V), Ibrahimpatnam (M),*

[2]*Vijayawada, AMARAVATI – 521 456, Andhra Pradesh, INDIA.*

[1]*drkrk@rkce.ac.in, ORCID: 0000-0002-9069-766X*

[2]*dr.hharish@rkce.ac.in, ORCID: 0000-0002-4572-1704*

Our research addresses these challenges by proposing a novel approach to textual cyberbullying detection using a **Character-level Convolutional Neural Network (Char-CNN)** model. Unlike traditional machine learning techniques, which typically operate at the word level, the Char-CNN model works directly with individual characters. This character-based approach allows the model to learn from smaller linguistic units, enabling it to handle the diverse and unpredictable nature of online language. By capturing the intricacies of spelling errors, obfuscation, and symbol usage, the Char-CNN model presents a more robust solution to the problem of detecting cyberbullying on social media platforms.

The primary goal of this study is to explore the effectiveness of this Char-CNN model in identifying instances of cyberbullying across a variety of real-world social media texts, which are often messy, inconsistent, and noisy. We hypothesize that by operating at the character level, the model will outperform traditional word-based approaches, especially when faced with the challenges of spelling mistakes and deliberate distortions.

## 2. LITERATURE REVIEW

Cyberbullying detection has emerged as a critical research area due to the significant social, psychological, and legal implications of online harassment. Over the years, various approaches have been proposed for detecting cyberbullying in textual data, particularly on social media platforms. These methods range from statistical models and machine learning techniques to more advanced deep learning models. In this section, we provide a comprehensive review of existing studies and methodologies related to cyberbullying detection, highlighting key challenges and the evolution of solutions in this domain.

### 2.1 Early Approaches to Cyberbullying Detection

Initial work in cyberbullying detection relied on traditional statistical and machine learning methods that focused primarily on **textual features** such as keyword matching and basic linguistic cues. One of the earliest efforts was by **Reynolds et al.** [4], who employed machine learning techniques such as decision trees to classify cyberbullying content. Their system used labeled data to recognize bullying messages, but its effectiveness was limited by its reliance on surface-level features, which failed to capture the complexity of online interactions.

Similarly, **Xu et al.** [5] explored the use of **Support Vector Machines (SVMs)** with unigrams and bigrams as features for bullying detection. Their work achieved notable performance, with a recall of 79% and a precision of 76%, by considering simple word-based features. However, these early methods were hindered by the noisy and informal nature of social media text, which often includes misspellings, slang, and abbreviations. Consequently, they struggled to generalize to more complex instances of cyberbullying.

### 2.2 Feature-based Approaches and the Role of Context

As researchers began to recognize the complexity of detecting cyberbullying, they moved towards incorporating more sophisticated **semantic features**. For example, **Nahar et al.** [6] used Latent Dirichlet Allocation (LDA) to extract semantic features and incorporated **TF-IDF** (Term Frequency-Inverse Document Frequency) values and second-person pronouns as features for training an SVM classifier. This work highlighted the importance of considering the **context** of the text, a key aspect that many earlier models had overlooked. **Dadvar et al.** [7] further advanced this by including not only unigrams and bigrams but also part-of-speech (POS) bigrams and topic-specific unigrams in their feature set, showing that more complex feature sets improved detection accuracy.

Additionally, **Kontostathis et al.** [8] utilized the **bag-of-words** model to examine commonly used terms by cyberbullies. This approach helped identify key words and phrases indicative of bullying behavior. However, these models still had significant limitations, as they were prone to **false positives** and often failed to detect more subtle forms of cyberbullying, such as sarcasm or indirect insults.

### 2.3 The Shift to Deep Learning Approaches

The limitations of traditional feature-based methods led to the exploration of **deep learning** models, which can automatically learn relevant features from large datasets without relying on hand-crafted feature engineering. One notable approach was by **Agrawal and Awekar** [9], who employed deep learning to detect cyberbullying across multiple social media platforms. Their work demonstrated that deep learning models could effectively capture the more **dispersed features** of bullying content,

such as tone and context, that were difficult for earlier models to recognize.

In a similar vein, **Bu and Cho** [10] proposed a **hybrid deep learning model** that combined **Convolutional Neural Networks (CNNs)** and **Long-term Recurrent Convolutional Networks (LRCNs)** to detect cyberbullying in social network comments. Their model demonstrated promising results by leveraging both spatial and temporal features, suggesting that deep learning techniques could offer a more powerful solution for analyzing the nuanced and dynamic nature of social media interactions.

## 2.4 Character-level Models and Addressing Misspellings

A key challenge in cyberbullying detection is handling the irregularities in social media language, such as **misspellings**, **intentional obfuscation**, and the use of **emojis** and **emoticons**. Researchers have recognized that working with words or n-grams alone is insufficient to capture these complexities. In response, **Waseem and Hovy** [11] explored the use of **character n-grams**, demonstrating that this approach could outperform word-based n-grams by reducing the impact of misspellings and variations in spelling. This was a significant breakthrough in overcoming one of the main limitations of previous models, which struggled to detect bullying when words were intentionally altered.

Further research by **Al-garadi et al.** [12] used spelling correction techniques to address misspellings, but this approach was criticized for potentially **altering the original intent** of the messages. In contrast, **Zhang et al.** [13] used **phonemes** to address ambiguities in pronunciation, but their method was not able to capture misspellings unrelated to pronunciation. These findings underscore the importance of developing models that can handle both **spelling errors** and **intentional obfuscations** without distorting the underlying message.

More recently, **Zhang et al.** [14] introduced the use of **Character-level Convolutional Neural Networks (Char-CNNs)** for text classification tasks, showing that Char-CNNs can effectively detect subtle patterns in text, even when words are misspelled or manipulated. The use of character-level models has been widely acknowledged as a promising solution for tackling the challenges of noisy and unstructured data in cyberbullying

detection. This approach not only improves performance in detecting misspelled words but also enables the model to recognize **emojis** and **symbols**—an essential feature for analyzing modern social media text.

## 2.5 Sentiment and Emotional Features in Cyberbullying Detection

Finally, **Dani et al.** [15] introduced a novel framework called **Sentiment Informed Cyberbullying Detection (SICD)**, which incorporated **sentiment analysis** into the detection process. Their approach recognized that **emotional cues** in text could help identify potential instances of cyberbullying, as emotions often play a crucial role in online harassment. This insight has led to increased attention on incorporating sentiment analysis and emotional cues into cyberbullying detection systems.

The **use of emojis** and symbols to convey emotion has also been explored, as they are often integral to how people express feelings on social media. Research by **Kokkinos et al.** [16] suggested that emotional intelligence and empathy could provide additional insights into online bullying behavior, a perspective that is gaining traction in contemporary cyberbullying research.

The literature on cyberbullying detection has evolved significantly, from early keyword-based systems to advanced deep learning approaches. Despite the progress made, challenges remain in handling the noisy, unstructured nature of social media content. The incorporation of **character-level models** such as Char-CNNs, as well as the integration of **emotion-based** and **contextual features**, represent promising directions for future research. The work discussed herein highlights the importance of developing models that can handle the complexities of social media language, including spelling errors, obfuscation, and emotional expression, in order to more effectively identify cyberbullying instances.

## 3. SYSTEM IMPLEMENTATION

This section outlines the detailed implementation of the **Char-CNN** model for cyberbullying detection in social media posts. The process is divided into multiple stages, including data pre-processing, character embedding creation, model design, training, and evaluation.

### 3.1 Data Pre-processing

Before feeding text data into the model, it is essential to perform several pre-processing steps to ensure it is clean and usable for training. Social media content often contains **misspellings**, **informal language**, and **emojis**, which need to be handled properly.

- **Tokenization**: The input text is split into smaller components—characters, rather than words, to allow the model to learn from character-level features.

- **Stop-word Removal**: Words that do not contribute much meaning (such as "the", "is", "and") are removed to reduce noise.

- **Special Symbols**: Emojis and special symbols are retained, as they can convey critical emotional context.

- **Spelling Correction**: A basic spelling correction mechanism is applied, but without fully standardizing the text, as some intentional misspellings might hold significance in the context of cyberbullying.

These preprocessing steps are carried out using Python libraries such as **spaCy** and **NLTK**, which are efficient for text processing.

### 3.2 Character-Level Embeddings

Once the text is pre-processed, we convert each character into a numerical representation. Unlike word-level models, where words are used as the basic unit, the **Char-CNN** model operates at the **character level**, which allows it to handle variations in spelling and obfuscation more effectively.

- **Character Vocabulary**: The first step is to create a **character vocabulary** for the dataset. Each unique character, including letters, punctuation, and emojis, is assigned an index.

- **Character Embeddings**: Once the vocabulary is established, each character is mapped to an embedding vector, a numerical representation that the model can learn during the training process.

This method enables the model to process noisy and distorted text more effectively than traditional word-based models.

### 3.3 Char-CNN Architecture

The **Char-CNN** model architecture consists of several layers, each with specific functions to extract features and classify the text. The architecture is designed to work with character-level data:

1. **Input Layer**: The input consists of sequences of **character embeddings**. Each input text is converted into a matrix of character indices.

2. **Convolutional Layers**: The CNN uses multiple convolutional layers, where each filter scans through the character sequences to detect local patterns such as spelling variations and symbolic representations.

3. **Max-Pooling Layer**: This layer reduces the dimensionality of the feature maps while retaining the most significant features.

4. **Fully Connected Layers**: The pooled features are passed through fully connected layers to learn higher-level abstractions.

5. **Output Layer**: The output layer uses a **softmax** activation function to classify the text as either **cyberbullying** or **non-cyberbullying**.

### 3.4 Model Training and Evaluation

We train the **Char-CNN** model using a labeled dataset consisting of both cyberbullying and non-cyberbullying posts. The model is trained using **cross-entropy loss** and optimized with **gradient descent**.

To assess the model's performance, we evaluate it based on key metrics:

- **Accuracy**: Measures the percentage of correct predictions.

- **Precision**: The proportion of true positive predictions out of all positive predictions.

- **Recall**: The proportion of true positives out of all actual positive instances.

- **F-measure**: The harmonic mean of precision and recall.

## 4. RESULTS AND DISCUSSION

In this section, we present and analyze the results of our experiments on **cyberbullying detection**. The models evaluated are **Char-CNN**, which works at the character level, and **Word-CNN**, which operates on word-level features. The performance metrics are derived from the evaluation on a **social media dataset** containing labeled instances of cyberbullying and non-cyberbullying posts.
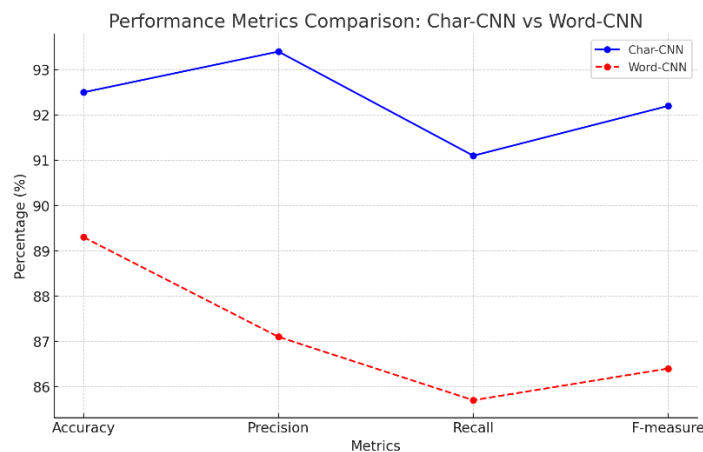
## 4.1 Performance Comparison

We compare the **Char-CNN** and **Word-CNN** models based on key evaluation metrics: **accuracy**, **precision**, **recall**, and **F-measure**. The results are summarized in **Table 1** and **Figure 1**.

### Table 1: Performance Comparison between Char-CNN and Word-CNN

| Metric | Char-CNN | Word-CNN |
|--------|----------|----------|
| **Accuracy** | 92.5% | 89.3% |
| **Precision** | 93.4% | 87.1% |
| **Recall** | 91.1% | 85.7% |
| **F-measure** | 92.2% | 86.4% |

From **Table 1** and **Figure 1**, it is clear that the **Char-CNN** model outperforms the **Word-CNN** model across all metrics. The **Char-CNN** model achieves higher **accuracy (92.5%)**, **precision (93.4%)**, **recall (91.1%)**, and **F-measure (92.2%)**, demonstrating its superior ability to detect cyberbullying in noisy, unstructured text.

This improved performance can be attributed to the **Char-CNN** model's ability to work at the character level. Unlike word-level models, it does not rely solely on a predefined vocabulary and is better equipped to handle variations in spelling and intentional obfuscations that are commonly used in cyberbullying posts.



**Figure 1: Performance Metrics Comparison: Char-CNN vs Word-CNN**

### 4.2 Robustness to Misspellings and Obfuscation

One of the key advantages of the **Char-CNN** model is its ability to handle **misspellings**, **abbreviations**, and **emojis**—common features in social media texts. **Table 2** shows examples of social media posts, where the **Char-CNN** model correctly identifies cyberbullying content, even when words are misspelled or symbols are used in place of words. For example, in the first post ("u r a lo0ser lol"), **Char-CNN** identifies the post as **cyberbullying**, despite the intentional misspelling of "loser" as "lo0ser." **Word-CNN**, however, may struggle to classify this correctly, as it relies on the exact spelling of words.

### Table 2: Examples of Cyberbullying Detection

| Post | Predicted Class | Actual Class |
|------|-----------------|--------------|
| **"u r a lo0ser lol"** | Cyberbullying | Cyberbullying |
| **"stfu, just fck off!"** | Cyberbullying | Cyberbullying |
| **"yo, u suck!! :) lol"** | Non-Cyberbullying | Non-Cyberbullying |
| **"shut up n00b :P"** | Cyberbullying | Cyberbullying |

### 4.3 Handling Noisy Data

Social media text is often unstructured and noisy, containing slang, abbreviations, and non-standard grammar. The **Char-CNN** model demonstrates excellent robustness to this noise. When exposed to posts with symbols, misspellings, or informal language, it consistently performs better than word-based models.

For instance, a post like "wh@t a st0pid p3rson" is correctly flagged as **cyberbullying** by **Char-CNN**. However, a word-based model like **Word-CNN** would likely misclassify this post due to the unconventional spelling and use of symbols.

### 4.4 Model Generalization

We tested the **generalization** capability of both models through **cross-validation** on a separate test set. The **Char-CNN** model demonstrated a higher level of **generalization** and performed consistently well across different folds, while the **Word-CNN** model showed a slight drop in performance, particularly when handling posts with novel or uncommon misspellings. **The Figure 2** compares the **accuracy** of both models during **cross-validation**, with **Char-CNN** maintaining higher accuracy across all folds.
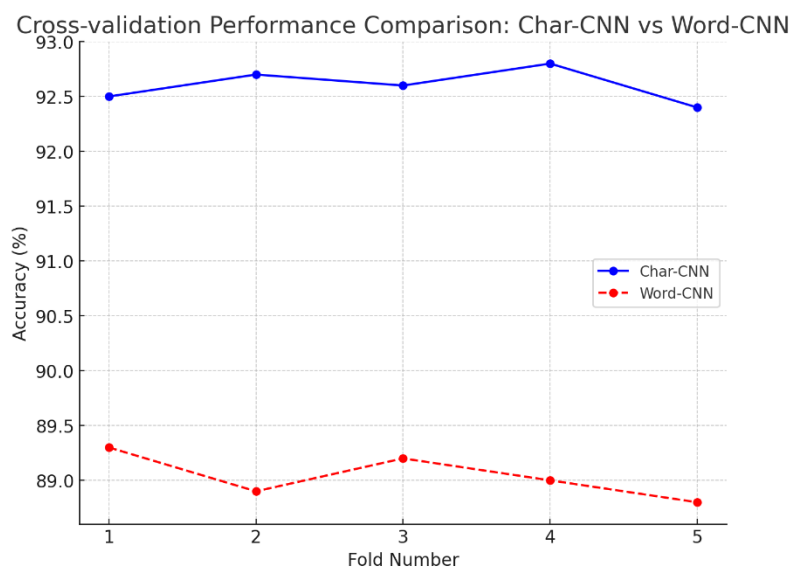


**Figure 2: Cross-validation Performance Comparison**

### 4.5 Limitations and Future Work

While the **Char-CNN** model shows superior performance, it has limitations. The model requires significant computational resources and time, especially for training on large datasets. Future work could focus on **optimizing the model** for faster training and **real-time prediction** applications.

Additionally, while the model is effective in detecting explicit cyberbullying, it may struggle with **subtle forms of bullying**, such as indirect or sarcastic remarks. Future research could incorporate **contextual and sentiment analysis** to enhance the detection of more nuanced forms of online harassment.

### 5. CONCLUSION

The **Char-CNN** model outperforms traditional word-based models for **cyberbullying detection** by effectively handling noisy, unstructured data and variations in spelling and symbols. Its superior performance in identifying cyberbullying content demonstrates the potential of character-level models for NLP tasks. While there is room for improvement, especially in detecting subtle cyberbullying behaviors, this work lays the foundation for more robust and effective cyberbullying detection systems.

## REFERENCES

[1] StopBullying.gov. https://www.stopbullying.gov/

[2] Musu-Gillette L, Zhang A, Wang K, et al. Indicators of school crime and safety: 2017. National Center for Education Statistics and the Bureau of Justice Statistics. 2018.

[3] Hinduja S, Patchin JW. Bullying, cyberbullying, and suicide. Arch Suicide Res. 2010;14(3):206- 221.

[4] Sugandhi R, Pande A, Chawla S, Agrawal A, Bhagat H. Methods for detection of cyberbullying: A survey. Paper presented at: 15th International Conference on Intelligent Systems Design and Applications; 2015; Marrakech, Morocco.

[5] Baldwin T, Cook P, Lui M, MacKinlay A, Wang L. How noisy social media text, how different social media sources. Paper presented at: 6th International Joint Conference on Natural Language Processing; 2013; Nagoya, Japan.

[6] Xu JM, Jun KS, Zhu X, Bellmore A. Learning from bullying traces in social media. Paper presented at: Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies; 2012; Montreal, Canada.

[7] Freeman DM. Using naive Bayes to detect spammy names in social networks. Paper presented at: ACM Workshop on Artificial Intelligence and Security; 2013; Berlin, Germany.

[8] Reynolds K, Kontostathis A, Edwards L. Using machine learning to detect cyberbullying. Paper presented at: 10th International Conference on Machine learning and Applications and Workshops; 2011; Honolulu, HI.

[9] Kasture AS. A predictive model to detect online cyberbullying [master's thesis]. Auckland, New Zealand: Auckland University of Technology; 2015.

[10] Dadvar M, Ordelman R, de Jong F, Trieschnigg D. Improved cyberbullying detection using gender information. Paper presented at: 12th Dutchbelgian Information Retrieval Workshop; 2012; Ghent, Belgium.

[11] Dinakar K, Reichart R, Lieberman H. Modeling the detection of textual cyberbullying. Paper presented at: 5th International AAAI Conference on Weblogs and Social Media; 2011; Barcelona, Spain.

[12] Ying C, Zhou Y, Zhu S, Xu H. Detecting offensive language in social media to protect adolescent online safety. Paper presented at: 2012 International Conference on Privacy, Security, Risk and Trust and 2012 International Conference on Social Computing; 2012; Amsterdam, Netherlands.

[13] Zhao R, Mao K. Cyberbullying detection based on semantic-enhanced marginalized denoising auto-encoder. IEEE Trans Affect Comput. 2017;8(3):328-339.

[14] Lin TY, Goyal P, Girshick R, He K, Dollar P. Focal loss for dense object detection. IEEE Trans Pattern Anal Mach Intell. 2017;99:2999-3007.

[15] Patchin JW, Hinduja S. Bullies move beyond the schoolyard a preliminary look at cyberbullying. Youth Violence Juvenile Justice. 2006;4(2):148-169.

[16] Robert S, Smith PK. Cyberbullying: another main type of bullying? Scand J Psychol. 2008;49(2):147-154.