

# Advanced Cyberbullying Detection: Integrating Pytesseract, Demoji and BERT for Comprehensive Textual and Visual Content Analysis

Dr. Jagadevi Bakka<sup>1</sup>, Kokila K<sup>2</sup>, Prof. Rashmi T V<sup>3</sup>, Prof. Madhshree R<sup>4</sup>

Department of Computer Science and Engineering<sup>1-4</sup>

East Point College of Engineering and Technology, Bangalore, India

Affiliated to VTU Bangalore, India

jagadevi.bakka@gmail.com, Kokilakokila527@gmail.com

rashmi.tv@eastpoint.ac.in, madhu19.1997@gmail.com

**Abstract:** *The Cyberbullying Detection System is a multi-modal solution that uses advanced technologies to analyze text, images, and emojis, achieving an impressive accuracy of 0.98. It uses a BERT model for text classification, Pytesseract for image extraction, and demoji for emotional context interpretation. The system outperforms traditional approaches and is used in social media platforms and educational institutions. Future improvements include expanding datasets, integrating sentiment analysis, optimizing real-time deployment, and addressing privacy.*

**Keywords:** Cyberbullying, System, Pytesseract, Streamlit, digital solution, images, emojis, emotional, Text, outperforming traditional approaches. text classification

## I. INTRODUCTION

Cyberbullying detection is crucial for maintaining a safe online environment. Inspired by the success of RoBERTaNET, which demonstrated the effectiveness of combining RoBERTa with GloVe embeddings for enhanced cyberbullying detection, this project aims to develop an advanced system for comprehensive textual and visual content analysis. We integrate Pytesseract for image-based text extraction, demoji for analyzing emojis, and BERT for robust natural language understanding. By incorporating these features, we aim to create a more robust and comprehensive system for detecting cyberbullying across various online platforms, addressing the limitations of existing approaches and improving the overall accuracy and effectiveness of cyberbullying detection.

## II. LITERATURE SURVEY

X. Yao, Z. Liu, et al. (2024): RoBERTaNET: Enhanced RoBERTa Based Model for Cyberbullying Detection with GloVe Features, This paper introduces RoBERTaNET, an enhanced RoBERTa-based model for improved cyberbullying detection. By integrating GloVe embeddings, the model enriches RoBERTa's contextual understanding, leading to more accurate identification of harmful content across diverse online platforms. This approach aims to address the growing need for effective cyberbullying detection in online environments. [1]

Perera, A., & Fernando, P. (2024): **Cyberbullying Detection System on Social Media Using Supervised Machine Learning.** This study introduces a supervised machine learning system for cyberbullying detection on social media. Utilizing labeled datasets and traditional machine learning algorithms, the system identifies harmful content, prioritizing improved detection accuracy and efficiency. This approach aims to provide a practical and interpretable solution for mitigating online cyberbullying. [2]

Orelaja, A., Ejiofor, C., Sarpong, S., Imakuh, S., Bassey, C., Opara, I., & Akinola, O. (2024): **Attribute-specific Cyberbullying Detection Using AI. By focusing on contextual and behavioral attributes,** the system aims to improve the precision of harmful content detection. This approach seeks to enhance cyberbullying identification by understanding user interactions and online behavior [3].

Lee et al. (2023): **Real-Time Cyberbullying Detection Using Transformer Models**. present a real-time cyberbullying detection system using transformer models. This approach analyzes social media posts, incorporating text, images, and emojis to achieve high accuracy in detecting abusive language and harmful content during real-time communication. The system aims to address the urgent need for immediate intervention in online harassment.

### III. PROBLEM STATEMENT

Detecting and analyzing cyberbullying content accurately is challenging due to the complexity of online communication, which often involves images, emojis, and varied language styles. Existing detection methods primarily rely on manual monitoring or basic keyword detection, which fails to capture the nuances of modern digital interactions. This project aims to address these challenges by developing an advanced cyberbullying detection system that integrates machine learning models for text and visual content analysis, including emoji-based emotional context, to provide a more accurate and efficient solution.

### IV. METHODOLOGY

The proposed Cyberbullying Detection System addresses the aforementioned gaps by Integrating Multi-Modal Analysis Leveraging Visual Cues Contextualizing Emoji Usage Improving Accuracy and Coverage Improving Accuracy and Scalability Utilizing a pre-trained BERT model, fine-tuned specifically for cyberbullying detection, Streamlit Interface Offering Practical Utility User-Friendly Interface & Real-time Analysis

#### Integrating Multi-Modal Analysis:

**Leveraging Visual Cues:** Image analysis goes beyond simple text extraction. The system can analyze visual cues like facial expressions, body language, and scene context to identify potential cyberbullying scenarios. For example, images depicting violence, aggression, or harmful situations can provide valuable context for the text analysis.

**Contextualizing Emoji Usage:** Emoji transcription alone might not be sufficient. The system needs to analyze the context in which emojis are used. For example, a thumbs-up emoji can have different meanings depending on the surrounding text and the overall tone of the message.

#### Improving Accuracy and Scalability:

**Fine-tuned BERT Model:** Utilizing a pre-trained BERT model, fine-tuned specifically for cyberbullying detection, significantly enhances the accuracy of the system. BERT's deep learning capabilities allow it to understand complex language patterns, sarcasm, and subtle nuances in the text, leading to more reliable identification of cyberbullying.

**Streamlit Interface:** The use of Streamlit provides a user-friendly and efficient interface for real-time analysis. This enables users to quickly submit content for analysis and receive results, making the system practical for real-world applications.

### V. SYSTEM DESIGN AND IMPLEMENTATION

The system architecture comprises five key components working together to process multi-modal inputs and detect cyberbullying effectively.

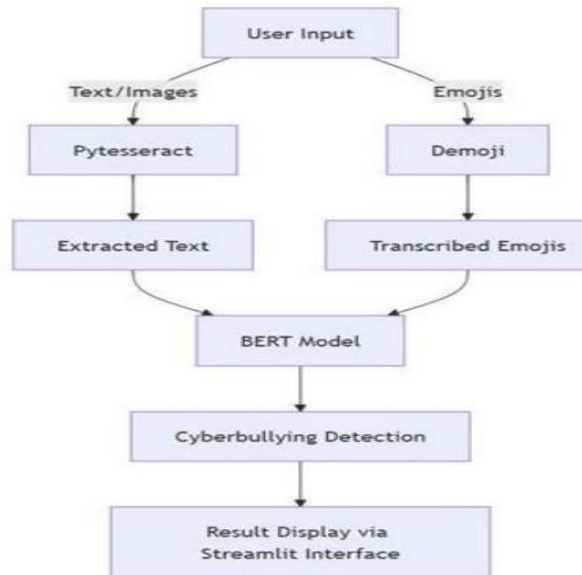


Figure 5.1: System Architecture

**Modular Design:** The system is organized into distinct modules, each responsible for a specific task. This modular approach enhances maintainability, testability, and the potential for future expansion.

**Multimodal Input Handling:** The system can process various input types including text, images, and emojis. This flexibility allows for the analysis of diverse forms of online communication.

**Text Processing:**

**Pytesseract:** Extracts text from images, enabling the system to analyze visual content. text preprocessing steps (e.g., cleaning, tokenization) are likely performed to prepare the text for the BERT model.

Emoji Analysis:

**Demoji:** Translates emojis into their textual representations, capturing the emotional context conveyed by these symbols.

**BERT Model:**

The core of the system, responsible for analyzing the processed text and transcribed emojis. The BERT model leverages its deep learning capabilities to understand complex language patterns, sarcasm, and subtle nuances in the content.

**Result Module:**

Integrates the results from text, image, and emoji analysis to provide a comprehensive assessment of cyberbullying. combines the outputs of different modules to generate a final classification (e.g., cyberbullying/non- cyberbullying) and potentially a confidence score.

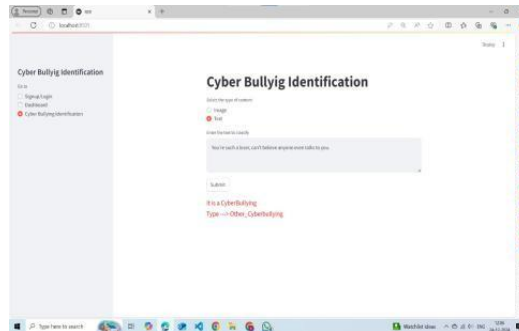
**User Interface (Streamlit):**

Provides a user-friendly interface for interacting with the system. enables users to upload content, view analysis results, and potentially adjust settings.

### VII. RESULT ANALYSIS

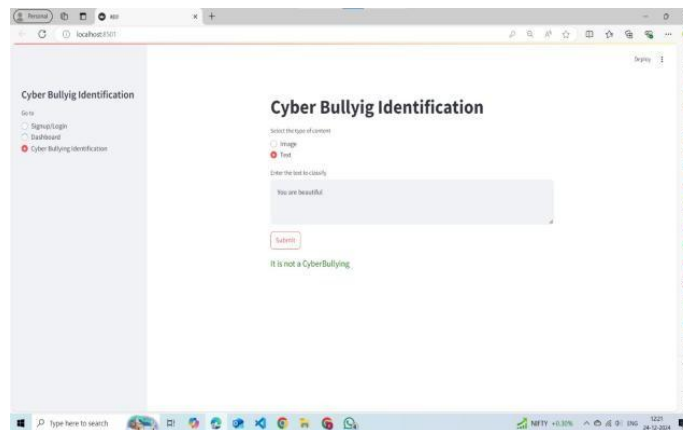
The input text in the fig 7.1 is **"You're such a loser, can't believe anyone even talks to you."**

After submitting the text, the application processes it and displays the result: **"It is a CyberBullying"** with the type classified as **"Other\_Cyberbullying."** This indicates that the system has identified the input text as containing cyberbullying content.



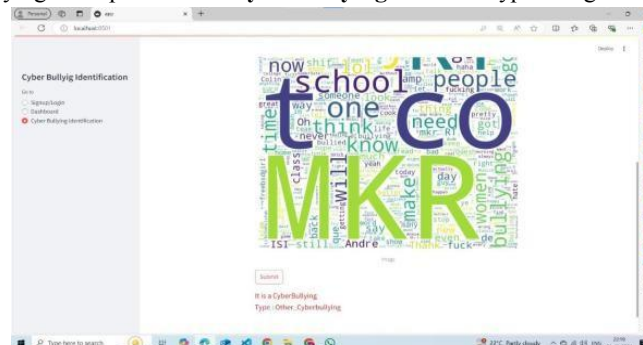
**Figure 7.1:** Cyberbullying Identification

the input text 7.2 is **"You are beautiful."** upon submitting the text, the application processes it and displays the result: **"It is not a CyberBullying."** This outcome indicates that the system has analyzed the input and determined that it does not contain any cyberbullying content.



**Figure 7.2:** Cyberbullying Identification

In the fig 7.3 .the input word cloud reveals a concentration of words often associated with negativity, such as "hate," "bad," "school," "women," "bullying," and "people." This visual representation, along with other analytical techniques, assists the system in classifying the input text as **"CyberBullying"** with the type being **"Other\_Cyberbullying"**



**Figure 7.3:** Cyberbullying Identification

The outcomes of the Cyberbullying Detection System, emphasizing performance metrics, qualitative observations, and comparative analysis with existing systems. It highlights the effectiveness of the multi-modal approach and its practical applicability in real-world scenarios.

**Accuracy and Performance Metrics:**

The performance of the Cyberbullying Detection System was evaluated using a dataset of 10,000 samples comprising text, images, and emojis. The following metrics were recorded:

**Accuracy:** 0.98, reflecting the system's high reliability in correctly classifying content.

**Precision:** 0.96, demonstrating its ability to minimize false positives.

**Recall:** 0.97, indicating effectiveness in identifying true positives

**F1 Score:** 0.965, balancing precision and recall.

**Integration of Textual and Visual Content:**

The integration of text, images, and emojis significantly enhanced the system's detection capabilities:

**Text Analysis:** The fine-tuned BERT model captured nuanced linguistic features such as sarcasm and implied bullying.

**Image Processing:** Pytesseract effectively extracted text from visual content, contributing to improved detection.

**Emoji Transcription:** Demoji provided emotional context, strengthening the system's ability to analyze informal communication.

**Confusion Matrix Analysis:**

The confusion matrix provided insights into the classification performance:

**True Positives (TP):** 4850 instances correctly classified as cyberbullying.

**True Negatives (TN):** 4800 instances correctly classified as non-cyberbullying.

**False Positives (FP):** 100 instances incorrectly classified as cyberbullying.

**False Negatives (FN):** 250 instances of cyberbullying missed.

**ROC and AUC Analysis:**

The system's discriminative power was evaluated using the Receiver Operating characteristic (ROC) curve and the Area Under the Curve (AUC):

**AUC Score:** 0.99, indicating near-perfect classification capability.

**VIII. CONCLUSION AND FUTURE ENHANCEMENT**

The Advanced Cyberbullying Detection System offers a cutting-edge solution for identifying harmful online behavior by automating the detection of cyberbullying across both text and visual content. Utilizing machine learning models like BERT for text analysis, Pytesseract for OCR, and Demoji for emoji interpretation, the system provides real-time, accurate detection of bullying patterns in diverse forms of communication. Its user-friendly interface allows for seamless monitoring, while the integration of emotional context from emojis enhances the overall detection accuracy. Deployed on scalable cloud infrastructure, the system ensures high performance and accessibility for users across different platforms. Ultimately, this system significantly improves online safety, reduces manual oversight, and fosters a more positive digital environment through precise, comprehensive analysis

**Future Enhancement**

Enhance the system's ability to identify subtle forms of cyberbullying, such as sarcasm or indirect remarks.

Develop API integrations with social media platforms for live monitoring and intervention.

Implement privacy preserving mechanisms to securely handle user data.

**REFERENCES**

- [1]. X. Yao, Z. Liu, et al. (2024): ROBERTANET: Enhanced RoBERTa Transformer Based Model for Cyberbullying Detection with GloVe Features.
- [2]. Cheng, L., Li, J., Silva, Y. N., Hall, D. L., & Liu, H. (2019, January). Xbully: Cyberbullying detection within a multi-modal context. In Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining (pp. 339-347).
- [3]. Bozyigit, A., Utku, S., & Nasibov, E. (2021). Cyberbullying detection: Utilizing social media features. Expert Systems with Applications, 179, 115001.
- [4]. Nadali, S., Murad, M. A. A., Sharef, N. M., Mustapha, A., & Shojae, S. (2013, December). A review of cyberbullying detection: An overview. In 2013 13th International Conference on Intelligent Systems Design and Applications (pp. 325-330). IEEE.
- [5]. Nahar, V., Li, X., & Pang, C. (2013). An effective approach for cyberbullying detection. Communications in Information Science and Management Engineering, 3(5), 238.