

Intel Shield

SHRUTHI B L¹, RAMYA H P², SANGEETHA C R³, SWATHI R⁴, YASHODHA S⁵

¹Assistant Professor, Department of Computer Science and Engineering, Rajiv Gandhi Institute of Technology, Bangalore

^{2, 3, 4, 5}Department of Computer Science and Engineering, Rajiv Gandhi Institute of Technology, Bangalore, India

Abstract- Cyberbullying is a growing threat in the digital age, particularly among teenagers and young adults on social media platforms. Traditional moderation methods such as manual review and keyword filtering often fail to capture the nuanced and evolving nature of abusive behaviour, especially across multimedia formats. This paper proposes Intel Shield, an AI-powered mobile application that leverages deep learning models like CNN, RNN, and LSTM to detect cyberbullying across text, images, and audio. The system includes a dynamic reputation score that penalizes users exhibiting repeated abusive behaviour, encouraging accountability. Real-time alerts, multilingual support, and seamless integration with popular communication platforms ensure timely intervention and a safer digital environment. With its modular, plug-and-play design, Intel Shield offers a scalable and effective solution to address cyberbullying across various online ecosystems.

Index Terms- Cyberbullying Detection, Deep Learning, CNN, RNN, LSTM, Multimodal Content Analysis, Reputation Scoring, Real-Time Monitoring, AI in Social Media, Digital Safety, -Play Architecture, Multilingual Support

I. INTRODUCTION

Social media changed how people now communicate. Yet with that, dangers quickly grew worldwide. Cyberbullying became common in digital interaction. Young users face emotional harm every day. Harassment spreads fast through posts or comments. Victims often struggle to report online abuse. Traditional moderation tools are slow and limited. They miss sarcasm, slang, and mixed language. Also, they do not support image detection.

Cyberbullying often hides behind creative content layers. Old systems use keywords and rule-based filters, leveraging contextual embedding and attention mechanisms to identify subtle patterns of abuse. For instance, hybrid approaches combining DE BERT a with Gated Broad Learning Systems (GBLS) have achieved high accuracy while integrating explainable AI techniques like LIME and SHAP to enhance model transparency. Such explainability is critical in ensuring that moderation systems are trusted by users and platform administrators. Research on multilingual and codemixed language detection, such as Hinglish-specific models using MURIL, highlights the importance of adapting detection systems to diverse linguistic contexts.

Another emerging trend is the development of severity-aware and user-specific models that classify not only the presence of cyberbullying but also its intensity. By incorporating user demographics, behavioural data, and psychological traits, these systems can provide more targeted interventions. Privacy preserving methods, such as block chain and federated learning frameworks combined with deep networks like LSTM and Deep Belief Networks (DBN), are also being explored to secure sensitive user information while enabling distributed model training.

Finally, researchers are focusing on real-time, streaming based detection systems capable of identifying and blocking harmful messages before they cause harm. These systems integrate large language models for continuous feature extraction and deploy explainable dashboards for content moderators. Despite these advancements, challenges remain in handling sarcasm, evolving slang, and cross-platform detection. Addressing these gaps will be crucial for building an effective, explainable, and privacy-aware Cyber Bully Blocker system.

They fail to understand context or emotions. Text-only models miss memes and voice messages. Manual reviews are time-consuming and not scalable. False positives and negatives reduce system trust. Language diversity makes filtering even more complex. Moderators cannot review everything in real time. Thus, detection methods need to evolve quickly. Real threats must be found and stopped. New systems must work across many formats. Intel Shield, a powerful mobile solution. It uses AI to scan harmful content. Built with CNN, RNN, and LSTM models. Intel Shield detects abuse in text, audio, visuals. A reputation score flags repeat online offenders. Multilingual support broadens its social media scope. It integrates easily with online communication platforms. The system encourages accountability in online spaces. Intel Shield makes digital communities safer and respectful.

II. LITERATURE SURVEY

Early methods used keyword lists for filtering. Simple matching flagged harmful or offensive content. They lacked understanding of slang or sarcasm. These systems failed with spelling and abbreviations. Context and emotional tone were completely ignored. Accuracy was low, and false rates high. Manual updates were needed for slang detection. These techniques were easy to bypass often. They worked best on structured formal content.

Unstructured social posts caused frequent detection failure.

Machine learning models were later introduced widely. Naïve Bayes and SVM were commonly used. They required pre-processing of all training data. Features like n-grams, TF-IDF were manually chosen. Still, they lacked strong semantic understanding capacity. Most models worked only for English language. They failed to detect sarcasm or humour. Scalability issues occurred with huge content streams. Multimedia and multilingual inputs were not supported. They showed poor generalization to unseen examples.

Deep learning improved detection with higher accuracy. LSTM captures emotion in message-based conversation text. CNN detects abuse in visual online content. RNN models handle long contextual

conversation data. Multilingual content and memes are often ignored.

Multimodal systems ensure safer, smarter cyber environments. They detect text, audio, and image content. Reputation scores flag repeat offenders in networks. Intel Shield offers scalable, multilingual, plug-and-play detection.

Cyberbullying detection has evolved significantly with the rise of advanced Natural Language Processing (NLP) and Machine Learning (ML) techniques. Early approaches relied on traditional classifiers such as Support Vector Machines (SVM), Naïve Bayes (NB), and Random Forests (RF) combined with textual features like TF-IDF or bag-of-words. While these methods achieved reasonable accuracy, they often failed to capture the contextual and semantic nuances of online harassment, especially in short, informal, or codemixed social media texts. Recent studies show that deep learning architectures such as Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), and Bi-LSTM outperform classical models by learning complex linguistic patterns and long-range dependencies in text.

III. METHODOLOGY

The development of the proposed diabetic retinopathy (DR) detection system follows a structured, end-to-end machine learning pipeline designed to ensure accuracy, transparency, and real-world deploy ability. The methodology comprises five core phases: data acquisition and preprocessing, data augmentation and class balancing, model architecture design, integration of explainable AI techniques, and system deployment via a web application interface.

A. Research Design

The study follows a mixed-methods approach, combining quantitative analysis of social media data with qualitative insights from literature and expert opinions. This ensures a comprehensive understanding of cyberbullying patterns, causes, and prevention methods.

B. Data Collection

Primary Data

Collected through online surveys targeting students, teachers, and parents to understand awareness levels, common platforms for cyberbullying, and its emotional impact. Conducted semi-structured interviews with school counsellors and cybercrime officials.

Secondary Data

Extracted from publicly available social media datasets (e.g., Twitter, Instagram). Reviewed academic journals, government reports, and NGO publications on cyberbullying cases and trends.

C. Data Pre-processing

Removed irrelevant or duplicate entries. Cleaned text data (removing special characters, links, emoji's). Tokenized and stemmed/lemmatized words for NLP processing. Applied stop-word removal to improve model performance

D. Data Analysis

Quantitative Analysis

Used statistical methods to identify patterns such as frequency of incidents, most affected age groups, and high-risk online platforms. Performed sentiment analysis on social media comments to detect bullying language.

Qualitative Analysis

Thematic analysis of interviews to extract common themes (e.g., types of harassment, coping mechanisms). Compared findings with existing literature to validate patterns. System Development and Deployment.

E. Testing and Validation

The model was tested with noisy text, misspellings, slang, and mixed-language inputs to check reliability

IV. EXISTING SYSTEM

Keyword filters match abusive terms without understanding context. Manual moderation is time-consuming and not easily scalable. Most systems fail to detect sarcasm or slang. Text only models ignore memes, images, and voice clips. False positives and negatives reduce trustworthiness greatly. Support for

regional or mixed languages is missing. Traditional systems cannot handle today's online abuse patterns.

A. Happens Across Multiple Digital Communication Platforms

Cyberbullying occurs through social media and messaging apps. It spreads via posts, comments, and direct messages. Victims often receive hate through anonymous user accounts. Platforms include Instagram, WhatsApp, Facebook, and gaming chats. Images, memes, and videos carry harmful intent too. Cyberbullies hide identity, making them hard to trace. Teenagers and children are frequent and silent victims. It affects their mental and emotional health deeply.

B. Hard to Detect by Traditional Filter Methods

Old systems use keyword filters to detect abuse. These fail when slang or sarcasm is used. Misspelled words bypass detection in many common cases. Multimedia bullying isn't detected with text-based tools. Memes and voice clips often carry hidden threats. Context, tone, and emotion aren't analysed accurately. Many cases are overlooked or falsely flagged wrongly. Advanced tools must understand full meaning and intent.

C. Effects Mental Health and Social Wellbeing Severely

Cyberbullying leads to stress, fear, and depression. Victims feel isolated, ashamed, and emotionally exhausted. Constant harassment lowers confidence and trust in people. Some may withdraw from school or social life. Sleep disorders and anxiety are common lasting effects. Teens may hesitate to report or seek help. Long-term exposure leads to serious mental health issues. Timely detection and support can save young lives.

V. PROPOSED SYSTEM

Intel Shield is an AI-powered cyberbullying detection tool.

It scans text, images, and audio content. Deep learning models improve accuracy and context recognition. CNN, RNN, LSTM help detect abusive behaviour. Reputation scores decrease with repeated harmful user actions. Multilingual input support ensures broad platform compatibility.

The proposed system, Cyber Bully Blocker, is designed to automatically detect, alert, and prevent cyberbullying incidents on digital platforms. It aims to create a safer online environment by using advanced text analysis and real-time monitoring techniques.

The system will function in the following way:

1. Data Monitoring – Continuously scan text-based content from chats, social media posts, comments, and emails for potential bullying behaviour.
2. Natural Language Processing (NLP) Engine – Use machine learning and NLP algorithms to analyse messages for harmful language, threats, harassment, hate speech, and offensive slang, even when disguised with symbols, misspellings, or code words.
3. Content Classification – Categorize messages into “Safe,” “Suspicious,” or “Bullying” based on severity levels.
4. Real-Time Alerts – Notify users, moderators, or parents instantly when harmful content is detected. This enables quick intervention before the situation escalates.
5. User-Friendly Dashboard – Provide a dashboard for administrators to review flagged content, check statistics, and take necessary actions.
6. Self-Help Resources – Suggest coping strategies, helpline contacts, and educational materials to victims for mental and emotional support.
7. Privacy & Security – Ensure that all data is processed securely, without storing personal information unnecessarily, to protect user privacy.
8. Continuous Learning – The system will improve detection accuracy over time by learning from new cyberbullying patterns and updating its keyword database and AI model.

CONCLUSION

The Intel Shield Cyber Bully Blocker was developed as a smart and proactive solution to one of the most pressing challenges of the digital age — cyberbullying. With the rise of social media, instant messaging apps, and online gaming platforms, harmful interactions have become easier to spread and harder to control. This system addresses the problem by providing real-time detection, alerting, and prevention of abusive language and harmful content.

The core strength of Intel Shield lies in its Natural Language Processing (NLP) engine, which can identify bullying not only in direct abusive terms but also in disguised forms such as slang, abbreviations, misspellings, and symbolic representations. The system categorizes the severity of detected content, enabling a swift and proportionate response — from simple warnings to immediate moderator or guardian alerts.

Another key advantage is its continuous learning capability, which allows the system to adapt to evolving cyberbullying patterns and language trends. By integrating a user-friendly dashboard, it also empowers teachers, parents, moderators, and platform administrators to review flagged content, track trends, and take preventive actions effectively. Importantly, Intel Shield operates with a privacy-first design, ensuring that personal information is protected, and only the necessary data is processed for detection purposes. In addition, it provides self-help resources, such as mental health support contacts and safety guidelines, to assist victims in coping with their experiences.

Overall, the Intel Shield Cyber Bully Blocker is not just a technological safeguard — it is a social responsibility tool that promotes respect, empathy, and safe digital communication. Its deployment can significantly reduce the frequency and impact of cyberbullying, making online spaces more secure for all, especially vulnerable groups like children and teenagers.

Future Scope:

With further development, Intel Shield can be integrated into popular social media APIs, gaming platforms, and educational portals to provide universal cyberbullying protection. Multi-language support, voice-message analysis, and AI-driven sentiment tracking are potential upgrades that can make the system even more effective in the fight against online harassment.

REFERENCES

- [1] Zhang, Z., Robinson, D., and Tepper, J. (2018). Detecting hate speech on Twitter using deep

- learning. European Semantic Web Conference, pp. 745–760.
- [2] Fortuna, P., and Nunes, S. (2018). A survey on automatic detection of hate speech in text. ACM Computing Surveys, 51(4), pp. 1–30.
- [3] Waseem, Z., and Hovy, D. (2016). Hateful symbols or hateful people? Predictive features for hate speech detection on Twitter. NAACL Student Research Workshop, pp. 88–93.
- [4] Badjatiya, P., Gupta, S., and Varma, V. (2017). Deep learning for hate speech detection in tweets. WWW Companion, pp. 759–760.
- [5] Al-Garadi, M. A., et al. (2020). Cyberbullying on social media: A review of detection techniques. IEEE Access, 8, pp. 96471–96487
- [6] R. R. Selvaraju et al., “Grad-CAM: Visual Explanations from Deep Networks via Gradient-Based Localization,” in *Proc. IEEE Int. Conf. on Computer Vision (ICCV)*, 2017, pp. 618–626. [For XAI reference]
- [7] J. S. Malik, H. Qiao, G. Pang, and A. van den Hengel, “Deep learning for hate speech detection: A comparative study,” 2022, arXiv:2202.09517.
- [8] M. H. Obaid, S. K. Guirguis, and S. M. Elkaffas, “Cyberbullying detection and severity determination model,” IEEE Access, vol. 11, pp. 97391–97399, 2023, doi: 10.1109/ACCESS.2023.3313113. [9] V.
- [9] Balakrishnan, S. Khan, and H. R. Arabnia, “Improving cyberbullying detection using Twitter users’ psychological features and machine learning,” Comput. Secur., vol. 90, Mar. 2020, Art. no. 101710, doi 10.1016/j.cose.2019.101710.
- [10] J. D. Angelis* and G. Perasso, “Cyberbullying detection through machine learning: Can technology help to prevent Internet bullying?” Int. J. Manage. Humanities, vol. 4, no. 11, pp. 57–69, Jul. 2020, doi: 10.35940/ijmh.k1056.0741120.
- [11] S. Prashar and S. Bhakar, “Real time cyberbullying detection,” Int. J. Eng. Adv. Technol., vol. 9, no. 2, pp. 5197–5201, Dec. 2019, doi: 10.35940/ijeat.b4253.129219.
- [12] J. Yadav, D. Kumar, and D. Chauhan, “Cyberbullying detection using pretrained BERT model,” in Proc. Int. Conf. Electron. Sustain.
- [13] Commun. Syst. (ICESC), Jul. 2020, pp. 1096–1100, doi: 10.1109/ICESC48915.2020.9155700.
- [14] B. A. H. Murshed, J. Abawajy, S. Mallappa, M. A. N. Saif, and H. D. E. Al-Ariki, “DEA-RNN: A hybrid deep learning approach for cyberbullying detection in Twitter social media platform,” IEEE Access, vol. 10, pp. 25857–25871, 2022, doi: 10.1109/ACCESS.2022.3153675.