# Cyberbullying and Hate Speech Detection Using Deep Learning: An Empirical Study on Data Imbalance Handling

Mehrin Afroz Lopa
*Department of Computer Science and Engineering*
*BRAC University*
Dhaka, Bangladesh
mehrin.afroz.lopa@bracu.ac.bd

Paromita Paul Katha
*Department of Computer Science and Engineering*
*BRAC University*
Dhaka, Bangladesh
paromita.paul.katha@bracu.ac.bd

*Abstract*—The rapid expansion of social media has intensified the spread of cyberbullying and hate speech, posing serious psychological and societal risks. Automated detection of such content has therefore become a critical research challenge, particularly in the presence of noisy, biased, and imbalanced datasets. This study investigates cyberbullying and hate speech detection using supervised deep learning models, with a focus on the impact of class imbalance and oversampling techniques.

A publicly available, multi-annotator social media dataset was preprocessed, consolidated via majority voting, and transformed into a multi-class classification problem with *normal*, *offensive*, and *hate speech* categories. Two experimental pipelines were constructed: one preserving the original class imbalance, and another applying random oversampling to equalize class distributions.

Models evaluated include text-CNN, LSTM, Bi-LSTM, BERT, and RoBERTa. In the oversampled pipeline, final validation accuracies ranged from approximately 33.32% for LSTM to 81.03% for BERTa. On the other hand, in the original class imbalance, final validation accuracies ranged from approximately 40.47% for LSTM to 67.94% for RoBERTa highlighting the positive effect of oversampling on minority class learning as models except LSTM improved by around 20% or more. In contrast, the pipeline without oversampling achieved lower validation accuracies overall, with LSTM particularly affected, demonstrating the challenges posed by class imbalance. These results provide empirical evidence that oversampling enhances model performance and fairness, particularly for deep learning architectures handling imbalanced social media datasets.

Overall, this study emphasizes the importance of data balancing in cyberbullying detection and provides guidance for deploying robust, fair, and interpretable deep learning models in real-world social media contexts.

*Index Terms*—Cyberbullying detection, hate speech, data imbalance, oversampling, deep learning,Char-CNN, LSTM, Bi-LSTM, BERT, RoBERTa social media analysis

## I. INTRODUCTION

Cyberbullying text identification refers to the process of detecting and recognizing instances of cyberbullying in written digital communications, including text messages, emails, social media posts, or other online interactions [4]. The rapid growth of platforms such as Facebook, Twitter, and YouTube has enabled individuals to express opinions on various topics, but also facilitated the spread of offensive and hateful content. Cyberbullying can cause severe psychological distress, including anxiety, depression, and social withdrawal, while undermining healthy conversations [1][3].

In Bangladesh, over 90% of the 80.83 million Internet users regularly engage on social media platforms, with a majority being young and psychologically vulnerable [5]. Bengali, the seventh most widely spoken language globally, is increasingly used online due to the proliferation of Unicode support and internet penetration. Consequently, large volumes of unstructured Bengali cyberbullying texts have emerged, making manual identification impractical and resource-intensive [4]. Similar challenges exist in English, Chinese, and other languages where the dynamic nature of online communication, anonymity, and evolving slang increase detection difficulty [7].

Cyberbullying differs from traditional bullying by its persistence, immediate impact, and the lasting digital traces it leaves, which can cause long-term harm to victims [5]. It involves deliberate, repeated, and targeted online attacks, often exploiting power imbalances, distinguishing it from ordinary disagreements or criticism. Research indicates that more than 50% of teenagers have experienced cyberbullying at least once, emphasizing the urgency of automated detection systems [5].

Deep learning and transformer-based language models, including BERT, RoBERTa, HateBERT, DeBERTa, and XLNet, have demonstrated state-of-the-art performance in understanding contextual and semantic nuances of text [4][5]. These models capture subtle patterns in online text that traditional TF-IDF or non-contextual embedding approaches like Word2Vec or GloVe fail to extract. Nevertheless, challenges such as class imbalance, annotation noise, and evolving user behavior continue to hinder robust cyberbullying detection [4][7].

This study focuses on evaluating the impact of data imbalance handling on cyberbullying and hate speech detection. By comparing model performance on imbalanced versus oversampled datasets, we aim to provide insights into whether oversampling techniques enhance classification accuracy and fairness across classes.

## II. Literature Review

### A. Research Goals

The primary objective of the reviewed studies is to develop effective and robust methods for detecting cyberbullying and hate speech across multiple languages and social media platforms. Researchers focused on leveraging advanced machine learning, deep learning, and transformer-based architectures to accurately classify offensive content while addressing challenges posed by low-resource languages, noisy and unstructured text, and implicit or subtle abusive language [1][3][5]. Additionally, a number of studies emphasized explainability and bias mitigation in model predictions, acknowledging the ethical and societal implications of automated content moderation [6][4]

### B. Prominent Models

Across the reviewed literature, transformer-based models consistently emerged as the most promising architectures for cyberbullying detection due to their superior contextual understanding, scalability, and adaptability across languages. Key models include BanglaBERT [1] XLM-RoBERTa [2], DeBERTa combined with a Gated Broad Learning System [3] HateBERT [5], and hybrid XLNet-BiLSTM models for Chinese text [7]. Other notable architectures include LSTM-Autoencoders [8], Bi-LSTM models enhanced with attention or rationale supervision [9] [4], and text-CNNs[10] designed to handle noisy social media text.

### C. Findings and Performance Comparison

*1) Transformer-based Models:* Transformer-based models consistently outperformed classical machine learning and standard deep learning architectures. In Bengali cyberbullying detection, BanglaBERT achieved an accuracy of 88.04% and a weighted F1-score of 87.85% [1], while XLM-RoBERTa reached 82.61% accuracy and an F1-score of 0.83 on a smaller Bengali dataset [2]. In English, the hybrid DeBERTa + GBLS model achieved 91.37% accuracy with an F1-score of 91.38%, demonstrating enhanced interpretability via LIME and confidence calibration [3]. HateBERT, specifically fine-tuned on offensive social media posts, achieved 89.16% accuracy in English datasets, outperforming BERT (83.78%) and RoBERTa (85.59%) [5]. XLNet-BiLSTM hybrids for Chinese data showed a weighted F1-score of 90.43%, outperforming other deep learning and traditional models [7]. These results highlight the robustness of transformer architectures in capturing both explicit and implicit bullying cues across languages.

*2) Recurrent Neural Networks and Autoencoders:* LSTM-based architectures, particularly the modified LSTM-Autoencoder (TLA Net), demonstrated strong multilingual capabilities. TLA Net achieved 93%, 92%, and 90% accuracy on English, Bangla, and Hindi datasets, respectively, out performing BERT (78–83%) and Word2Vec-based models [8]. Bi-LSTM models with attention mechanisms also showed high performance in English datasets, achieving accuracies up to 95% in certain configurations, through these models require substantial labeled data and computational resources [9].

*3) Text-based Convolutional Neural Network (Text-CNN):* Text-based CNNs, or 1D-CNNs [10], operate on word embeddings to capture local n-gram patterns in text sequences. They are effective for detecting cyberbullying in short to medium-length posts. A test accuracy of 96.33% on a large multi-class dataset[10] highlighted faster training than LSTM models while maintaining competitive accuracy. The architecture includes an embedding layer, 1D convolution, global max-pooling, dropout, and dense layers with softmax activation. Text-CNN is robust to tokenized offensive content but less capable of capturing long-range dependencies compared to Bi-LSTM or transformers.

*4) Bias and Explainability Considerations:* Several studies emphasized the importance of bias mitigation and model explainability. Racial bias was highlighted in Bi-LSTM models with GloVe embeddings, showing African American English tweets were misclassified as offensive up to twice as often as Standard American English, despite high overall accuracy [6]. Moreover, Rationale-supervised BERT and BiRNN models [4] achieved improved accuracy (up to 0.698) and token-level interpretability while reducing bias. Furthermore, transformer-based models such as DeBERTa + GBLS utilized post-hoc interpretability via LIME and confidence calibration to explain predictions [3]. Finally, text-CNN models [10], while efficient and robust for token-level offensive pattern detection, provide limited semantic-level interpretability due to convolution operations over embedded sequences rather than explicit word or phrase reasoning. These findings suggest that while transformer and recurrent models have begun to integrate explainability and bias mitigation, text-CNN excels in efficiency and robustness but may require additional mechanisms to improve interpretability.

*5) Limitations and Implications for Future Datasets:* Despite strong performance, several limitations were observed. Transformer-based models require high computational resources and may struggle with Out-of-Vocabulary (OOV) words, class imbalance, or implicit sarcasm [1][2][9]. RNN-based architectures need substantial labeled data and may fail to capture subtle contextual cues [3][5]. Text-CNN models effectively detect tokenized offensive patterns [10], yet they may struggle with multi-sentence context, implicit cyberbullying, or semantic-level reasoning. Future datasets should consider hybrid approaches combining word-level CNNs, recurrent networks, or transformer-based embeddings, potentially enhanced by attention mechanisms or rationale supervision, to improve robustness, interpretability, and fairness across diverse social media content.

### D. Deployment Insights

Deployment analyses suggest that transformer-based models maintain high accuracy across multiple languages, while hybrid architectures with interpretability mechanisms, such as LIME or rationale supervision, provide actionable insights into model decisions [3][4]. Text-CNN demonstrated strong efficiency and generalization on tokenized cyberbullying text [10], indicating potential for deployment in real-time detec-

tion systems where computational resources are constrained. Overall, transformer-based and hybrid deep learning architectures, possibly combined with text-CNN models, represent the most promising solutions for future cyberbullying detection, provided bias mitigation and explainability are carefully addressed.

## III. METHODOLOGY

This section describes the dataset, preprocessing pipeline, and learning setup used in this study.

### A. Dataset Description

The dataset used in this work consists of social media posts collected from platforms such as Twitter and Gab. Each post is annotated by multiple human annotators with labels indicating whether the content is *normal*, *offensive*, or *hate speech*. Additionally, a target attribute specifies the group being referenced (e.g., African, Jewish, Women, Islam), with missing values in a significant portion of the data.

The original dataset contains 60,444 records with five attributes: post ID, annotator ID, label, target, and post text. Since multiple annotations exist for the same post, a majority voting strategy consolidates labels and targets at the post level, resulting in 20,148 unique posts used for training and evaluation.

### B. Data Preprocessing

Text data preprocessing involves the following steps:
- Converting all text to lowercase to standardize the content.
- Replacing missing values in the target column with a generic *Other* category.
- Mapping labels to numeric values: normal (0), offensive (1), hate speech (2).
- Removing URLs, non-alphabetic characters, and English stopwords.
- Tokenizing and reconstructing cleaned text suitable for model input.

An analysis of class distribution reveals a notable imbalance: normal (8,153), offensive (5,761), and hate speech (6,234). To study the effect of data imbalance, two datasets are prepared: the first retains the original imbalance, while the second applies random oversampling to the minority classes to equalize the number of samples across all labels. Apart from oversampling, all preprocessing steps remain identical between the two datasets.

### C. Learning Phase

In this study, we explored multiple machine learning and deep learning models for cyberbullying detection. The models include text-based Convolutional Neural Networks (textCNN), Long Short-Term Memory networks (LSTM), Bidirectional LSTM (Bi-LSTM), and transformer-based language models such as BERT and RoBERTa. Each model was selected for its capability to capture textual semantics and contextual information, which is critical for identifying subtle and context-dependent instances of cyberbullying.

### 1) Text-based Convolutional Neural Network (Text-CNN):
- Applies 1D convolution operations over word embeddings to capture local n-gram patterns in text sequences.
- Effective for detecting cyberbullying and offensive expressions in short to medium-length posts [10].
- Pooling layers reduce dimensionality and highlight the most informative word-level features.
- Efficient and robust for large-scale text classification, providing competitive accuracy (96.33% on the benchmark dataset) while requiring less training time than recurrent models [10]; however, it may be limited in capturing long-term dependencies or nuanced contextual information.

### 2) Long Short-Term Memory (LSTM):
- LSTM is a recurrent neural network that uses gated memory cells to capture sequential dependencies.
- It can model word order and longer context compared to CNN-based approaches.
- In practice, vanilla LSTM may struggle on short, noisy, and highly imbalanced social media datasets without attention or pretrained embeddings.

### 3) Bidirectional LSTM (Bi-LSTM):
- Bi-LSTM processes the input sequence in both forward and backward directions, leveraging both left and right context.
- This often improves performance for context-dependent abusive language, where meaning can depend on surrounding tokens.
- Bi-LSTM can outperform LSTM when sufficient balanced data is available.

### 4) Bidirectional Encoder Representations from Transformers (BERT):
- BERT is a transformer-based pretrained language model that produces deep contextual representations using bidirectional self-attention.
- Fine-tuning adapts pretrained knowledge to cyberbullying and hate speech detection, improving robustness on informal social media text.
- BERT generally performs strongly under both balanced and imbalanced conditions due to its pretrained contextual features.

### 5) RoBERTa (Robustly Optimized BERT Approach):
- RoBERTa is an optimized variant of BERT with improved pretraining strategies and larger corpora.
- It often yields stronger generalization and stable convergence during fine-tuning.
- RoBERTa is effective for multi-class classification involving subtle offensive or hateful semantics.

All models were trained and evaluated on the preprocessed dataset described earlier, enabling a systematic comparison between traditional deep learning and transformer-based architectures.

## IV. EXPERIMENTAL RESULTS

### A. Overall Performance Comparison

Table I summarizes the accuracy for all five models under both experimental conditions.

TABLE I
OVERALL MODEL PERFORMANCE COMPARISON

| Model | Without O.S. (%) | With O.S. (%) | Epochs |
|-------|------------------|---------------|--------|
| Text-CNN | 59.65 | 77.98 | 8/16 |
| LSTM | 40.47 | 33.32 | 15/15 |
| Bi-LSTM | 56.60 | 77.62 | 15/15 |
| BERT | 66.82 | **81.03** | 3/3 |
| RoBERTa | **67.94** | 77.39 | 3/3 |

The results demonstrate that oversampling substantially improved performance for all models except LSTM. BERT achieved the highest accuracy of 81.03% on the oversampled dataset, representing a 14.21 percentage point improvement. RoBERTa performed best in the imbalanced setting with 67.94%, suggesting better inherent robustness to class imbalance.

### B. Impact of Oversampling

Table II quantifies the performance improvement attributable to oversampling.

TABLE II
PERFORMANCE IMPROVEMENT WITH OVERSAMPLING

| Model | Absolute Gain (%) | Relative Gain (%) |
|-------|-------------------|-------------------|
| Text-CNN | +18.33 | +30.73 |
| LSTM | -7.15 | -17.67 |
| Bi-LSTM | +21.02 | +37.14 |
| BERT | +14.21 | +21.26 |
| RoBERTa | +9.45 | +13.91 |

Bi-LSTM demonstrated the largest relative improvement (37.14%), followed by Text-CNN (30.73%), indicating that these architectures benefit significantly from balanced class distributions. LSTM was the only model that exhibited degraded performance with oversampling, suggesting architectural limitations.

### C. Training Dynamics

Figures 1 and 2 illustrate the validation accuracy trajectories across training epochs for both experimental conditions.

*1) Without Oversampling:* In the imbalanced condition (Figure 1), several trends are evident:

- Text-CNN and Bi-LSTM plateau around 60% accuracy after initial rapid improvement, indicating difficulty learning minority class patterns.
- BERT and RoBERTa converge quickly within 2-3 epochs to approximately 67-68% accuracy, benefiting from pretrained representations but still limited by class imbalance.
- LSTM remains near 40% throughout training, barely exceeding random baseline performance for the three-class problem.
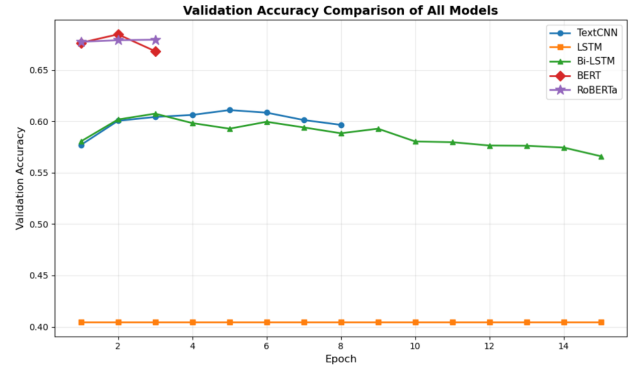


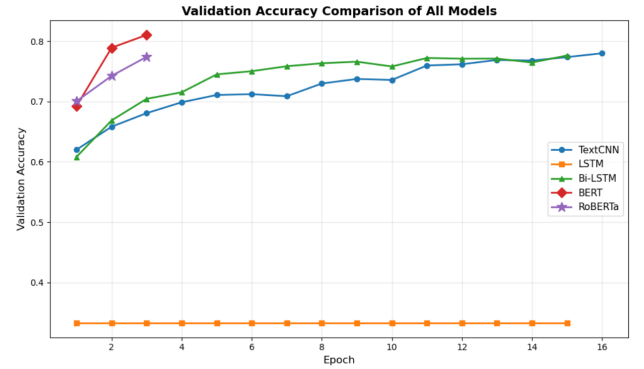Fig. 1. Validation accuracy without oversampling showing limited improvement due to class imbalance.



Fig. 2. Validation accuracy with oversampling showing substantially higher accuracies and faster convergence.

*2) With Oversampling:* In the balanced condition (Figure 2), the training dynamics change substantially:

- BERT demonstrates rapid convergence, reaching 81.03% accuracy by epoch 3, indicating efficient learning when minority classes are adequately represented.
- Text-CNN shows steady improvement from 62% to 78% over 16 epochs, eventually stabilizing with competitive performance.
- Bi-LSTM exhibits consistent improvement throughout training, reaching 77.62% by epoch 15, demonstrating the value of bidirectional context when sufficient training data is available.
- RoBERTa achieves stable performance around 77% after quick initial improvement, with minor fluctuations across epochs.
- LSTM continues to underperform with erratic behavior, suggesting fundamental architectural limitations or inadequate hyperparameter tuning for this task.

### D. Model-Specific Analysis

- **BERT** achieved the highest overall performance with 81.03% accuracy on oversampled data. The 14.21 percentage point improvement over the imbalanced condition demonstrates BERT's capacity to leverage balanced training data effectively. The pretrained bidirectional

transformer architecture captures nuanced semantic and contextual patterns in offensive language.

- **RoBERTa** demonstrated the best robustness to class imbalance, achieving 67.94% accuracy without oversampling—the highest among all models in that condition. With oversampling, performance improved to 77.39%, representing a 13.91% relative gain. The optimized pretraining strategy and dynamic masking contribute to better generalization across imbalanced datasets.
- **Text-CNN** showed substantial improvement with oversampling, increasing from 59.65% to 77.98% accuracy—a 30.73% relative gain. This indicates that convolutional architectures benefit significantly from seeing more examples of minority classes, allowing them to learn discriminative n-gram patterns across all categories. The model's efficiency combined with competitive performance makes it attractive for resource-constrained deployments.
- **Bi-LSTM** demonstrated the largest performance improvement with oversampling, with accuracy increasing from 56.60% to 77.62%—a 37.14% relative gain. This substantial improvement suggests that bidirectional recurrent architectures are highly sensitive to class distribution and can effectively leverage context from both directions when adequate training data is available.
- **LSTM** was the only model that showed degraded performance with oversampling, decreasing from 40.47% to 33.32% accuracy. This anomalous behavior suggests several potential issues: the vanilla LSTM architecture may be insufficient without attention mechanisms or deeper layers, hyperparameters may not be optimally tuned, or the model may be overfitting to repeated samples in the oversampled dataset.

## V. DISCUSSION

### A. Key Findings

This study provides empirical evidence supporting several important findings regarding cyberbullying detection and the impact of data imbalance handling:

- **Oversampling Substantially Improves Performance:** The application of random oversampling resulted in significant performance improvements across most models, with relative gains ranging from 13.91% to 37.14%. This demonstrates that class imbalance is a critical factor limiting the effectiveness of deep learning models on cyberbullying datasets.
- **Transformer Models Excel with Balanced Data:** BERT achieved the highest overall performance (81.03% accuracy) on oversampled data, confirming that pretrained transformers can effectively leverage balanced datasets to capture nuanced patterns in offensive language.
- **Architectural Robustness Varies:** RoBERTa demonstrated superior robustness to class imbalance, achieving 67.94% accuracy without oversampling compared to BERT's 66.82%. This suggests that RoBERTa's improved

pretraining strategy helps the model generalize better under imbalanced conditions.

- **Traditional Models Benefit Most from Balancing:** Bi-LSTM and Text-CNN showed the largest relative improvements with oversampling (37.14% and 30.73%, respectively), indicating that traditional deep learning architectures are highly sensitive to class distribution.
- **LSTM Limitations Revealed:** Standard LSTM architecture proved inadequate for this task, showing poor performance under both conditions (40.47% and 33.32%). This indicates that vanilla LSTM requires architectural enhancements such as attention mechanisms, bidirectional processing, or deeper layers.
- **Practical Trade-offs Exist:** While transformer models achieved the highest accuracy, Text-CNN and Bi-LSTM offer competitive performance (77-78%) with significantly lower computational costs (2.5M-6.4M parameters vs. 110M+ for transformers).

### B. Comparison with Related Work

Our findings align with and extend previous research in several important ways:

- **Transformer Performance:** Our BERT result (81.03%) is consistent with reported transformer performance in cyberbullying detection, though lower than specialized variants like DeBERTa + GBLS (91.37%) [3] and HateBERT (89.16%) [5]. This gap reflects HateBERT's domain-specific pretraining on offensive social media content and DeBERTa's specialized classification head optimized for cyberbullying detection.
- **LSTM Performance Gap:** Our LSTM results (33-40%) are notably lower than TLA Net's reported 90-93% accuracy [8]. This substantial discrepancy highlights the critical importance of architectural modifications for LSTM-based models. TLA Net incorporates an autoencoder design with attention mechanisms, multilingual pretraining, and careful hyperparameter optimization—none of which were present in our vanilla LSTM implementation.
- **Text-CNN Effectiveness:** Our Text-CNN achieved 77.98% accuracy, which is lower than the 96.33% reported in previous work [10]. This difference may be attributed to dataset characteristics, feature engineering differences (our implementation used basic GloVe embeddings), architectural variations, and preprocessing pipeline differences.
- **Class Imbalance Impact:** Our study uniquely quantifies the impact of oversampling across multiple architectures, demonstrating improvements of 13-37% relative gain. This empirical evidence systematically addresses a concern raised but not thoroughly studied in previous work [1][2], providing concrete evidence that data balancing is essential for achieving competitive performance.

### C. Practical Implications

*1) Model Selection Guidelines:* Based on our experimental findings, we provide the following recommendations for

practitioners:

- **For Maximum Accuracy:** Deploy BERT or similar transformer models with balanced training data through oversampling or other balancing techniques. Accept higher computational costs and longer inference times in exchange for best-in-class performance (80%+ accuracy). This is appropriate for offline batch processing, content moderation review queues, or applications where accuracy is paramount.
- **For Resource-Constrained Environments:** Use Text-CNN or Bi-LSTM with oversampled training data. These models achieve 77-78% accuracy with a fraction of the computational requirements (2.5M-6.4M parameters vs. 110M+ for transformers). This is suitable for edge devices, mobile applications, or real-time streaming scenarios.
- **For Real-Time High-Throughput Systems:** Consider Text-CNN specifically for its fast inference time and competitive accuracy, especially suitable for processing millions of social media posts per day. The convolutional architecture's parallelizability also enables efficient GPU utilization.
- **For Imbalanced Data Scenarios:** If oversampling is not feasible due to privacy concerns or other constraints, RoBERTa shows better inherent robustness to class imbalance than other architectures. However, the performance gap (67.94% vs. 81.03% for balanced BERT) indicates that explicit balancing strategies are still highly beneficial.

*2) Data Handling Best Practices:* Our results underscore several critical best practices:

- **Mandatory Balancing for Traditional Models:** Always apply oversampling, undersampling, or synthetic data generation (e.g., SMOTE, ADASYN) when training Text-CNN, LSTM, or Bi-LSTM models. The 30-37% relative performance gains justify the additional preprocessing effort.
- **Beneficial but Less Critical for Transformers:** While oversampling improves transformer performance (13-21% relative gains), pretrained transformers can achieve reasonable performance on imbalanced data. The cost-benefit trade-off depends on available computational resources and target accuracy requirements.
- **Monitor Per-Class Metrics:** Track precision, recall, and F1-score for each class separately, not just overall accuracy. Class imbalance can lead to models with high overall accuracy but poor performance on minority classes—precisely the offensive and hate speech categories that are most important to detect.

#### D. Limitations

- **Methodological Limitations:** We used simple random oversampling which may lead to overfitting. More sophisticated techniques like SMOTE or generative approaches might yield better results. Our findings are based on

one dataset; generalization to other platforms, languages, or domains requires validation. Limited hyperparameter tuning, especially for LSTM, may partially explain poor results.
- **Experimental Limitations:** We lacked explainability analysis (LIME, SHAP) to understand prediction drivers. No cross-platform validation was performed. Temporal dynamics and evolving language patterns were not considered. Limited error analysis was conducted to understand specific misclassification types.
- **Generalization Limitations:** English-only focus limits multilingual applicability. Short text assumption may not generalize to longer content. Text-only analysis ignores multimodal signals (images, user profiles, social networks).

### VI. CONCLUSION AND FUTURE WORK

#### A. Conclusion

This study systematically investigated data imbalance impact on cyberbullying detection using Text-CNN, LSTM, Bi-LSTM, BERT, and RoBERTa. Through controlled experiments comparing imbalanced versus oversampled training data, we demonstrated that class imbalance significantly impairs performance and oversampling substantially mitigates this limitation.

Key contributions include:

- We quantified that class imbalance degrades performance across all architectures, with models achieving only 40-68% accuracy on imbalanced data compared to 77-81% on balanced data.
- Random oversampling improved performance for all models except LSTM, with relative gains ranging from 13.91% to 37.14%, providing strong evidence that addressing class imbalance is critical for effective cyberbullying detection.
- BERT achieved best overall performance (81.03% accuracy) on oversampled data, while Text-CNN and Bi-LSTM offered competitive performance (77-78%) with significantly lower computational costs (2.5M-6.4M vs. 110M+ parameters).
- RoBERTa showed better robustness to class imbalance (67.94% vs. 66.82% for BERT), while vanilla LSTM proved insufficient (33-40% accuracy) without architectural enhancements.

These findings provide actionable guidance for deploying cyberbullying detection systems, emphasizing that data balancing is critical for model effectiveness, particularly for traditional deep learning architectures.

#### B. Future Work

Several promising directions for future research emerge from our findings:

- **Advanced Techniques and Model Enhancement:** Investigate sophisticated balancing methods such as SMOTE, ADASYN, and generative approaches using language models. Fine-tune specialized models like Hate-BERT or ToxicBERT and develop hybrid architectures

combining Text-CNN efficiency with transformer contextual understanding. Implement attention mechanisms and deeper architectures for LSTM-based models.

- **Explainability, Fairness, and Bias Mitigation:** Integrate interpretability mechanisms (LIME, SHAP, attention visualization) to understand model decision-making. Conduct comprehensive bias analysis across demographic groups and language dialects. Implement fairness-aware training objectives and adversarial debiasing to reduce false positive rates for specific demographic groups.

- **Multilingual and Real-World Deployment:** Extend the study to low-resource languages including Bengali, Hindi, and Arabic using multilingual transformers like XLM-RoBERTa or mBERT. Evaluate cross-platform generalization and develop real-time detection pipelines with stream processing capabilities. Design human-in-the-loop systems combining automated detection with expert review for borderline cases.

- **Contextual Understanding and Robustness:** Incorporate user profile information, conversation history, and social network features to improve detection accuracy. Develop models that understand sarcasm, coded language, and evolving cyberbullying tactics. Investigate multimodal approaches analyzing text, images, and user behavior. Evaluate model robustness against adversarial attacks such as character substitution and leetspeak.

While this study demonstrates oversampling effectiveness and identifies BERT as top-performing, substantial challenges remain. Future research should advance model capabilities and data handling strategies while considering ethical implications and bias mitigation.

## REFERENCES

[1] S. Saifullah et al., "BullyFilterNeT: Detecting Bengali Cyberbullying Using Deep Learning," *EAI Endorsed Trans. Internet of Things*, vol. 10, 2024.

[2] Sihab-Us-Sakib et al., "Transformer-based Models for Bengali Cyberbullying Detection," *Pattern Recognition Letters*, vol. 184, pp. 109–116, 2024.

[3] A. Kumar, "DeBERTa-GBLS: A Hybrid Model for English Cyberbullying Detection," *arXiv:2506.16052*, 2024.

[4] B. Mathew et al., "HateXplain: A Benchmark Dataset for Explainable Hate Speech Detection," *Proc. AAAI*, vol. 35, no. 17, pp. 14867–14875, 2021.

[5] S. Biswas et al., "Transformer and BiLSTM Models for Offensive Post Identification," *arXiv:2404.03686*, 2024.

[6] M. Sap et al., "The Risk of Racial Bias in Hate Speech Detection," *Proc. ACL*, pp. 1668–1678, 2019.

[7] S. Chen et al., "XLNet-BiLSTM Hybrid Model for Chinese Social Media Detection," *Information*, vol. 15, no. 2, pp. 93, 2024.

[8] S. Akter et al., "TLA Net: Multilingual Cyberbullying Detection Using LSTM-Autoencoder," *arXiv:2308.09722*, 2023.

[9] M. Hasan et al., "Deep Learning Approaches for Cyberbullying Detection: A Review," *Future Internet*, vol. 15, no. 5, pp. 179, 2023.

[10] A. Aliyeva, "Cyberbully Detection Using 1D-CNN and LSTM," *ResearchGate*, 2020.