

EMOTION-INFORMED CYBERBULLYING DETECTION SYSTEM

Dr. P.V.N RAJESWARI

Associate Professor in CSE Department of PBR VITS, Kavali, AP, India.
Email: rajeswari.pvn@visvodayata.ac.in

J.MURALI

Assistant Professor in CSE Department of PBR VITS, Kavali, AP, India. Email: murali.j@visvodayata.ac.in

Dr. R RAMESH BABU

Professor, HOD in ECE Department in Keshav Memorial College of Engineering, Hyderabad, Telangana.
Email: dr.rrameshbabu@gmail.com

BALAM KALYANI

Assistant Professor in ECE Department of Keshav Memorial College of Engineering, Hyderabad, Telangana. Email: balamkalyani92@gmail.com

Dr. SK. MASTHAN BASHA

Professor, VMTW, TG.

Abstract

The project addresses the serious issue of cyberbullying, recognizing its harmful consequences and the need for effective detection and resolution methods. Cyberbullying [2], a form of online aggression, poses challenges that require specialized techniques to identify and mitigate. The primary goal of the project is to propose advanced cyberbullying detection models. These models aim to go beyond traditional approaches by incorporating contextual, emotion, and sentiment features, recognizing the multi-dimensional nature of cyberbullying instances. This project is all about building an EDM using data from Twitter. These datasets undergo enhancements in terms of annotations. The EDM, along with lexicons, is utilized to extract emotions and sentiments from cyberbullying datasets, contributing to a more nuanced understanding of the emotional aspects of online interactions. Cyberbullying can be detected in part by appealing to people's emotions. This project shows how important emotions are for cyberbullying detection by using emotional cues to enhance detection algorithms. An extensive dataset tagged with emotions is now available for use in cyberbullying detection, thanks to this study. Academics can utilize this dataset to create a system for identifying cyberbullying based on emotions. Particularly in real-time applications, the project faces challenges due to dataset imbalances between cyberbullying and non-cyberbullying incidents. The goal is to develop detection models that effectively handle these imbalances, ensuring reliable performance across different scenarios. We aim to further enhance the performance of our model by exploring ensemble techniques, specifically utilizing LSTM and LSTM + GRU models, which have demonstrated an impressive 99% accuracy.

Index Terms: Cyberbullying, Behavioural Emotional Recognition Technology (BERT), Sentiment Analysis.

1. INTRODUCTION

Thanks to advancements in ICT, members of the online community can now post and respond to information created by themselves. Harassment, intimidation, control, manipulation, and threats have all been magnified by cyberbullies who have taken advantage of this convenience [1]. Cyberbullying refers to the deliberate and persistent use of cyberspace for harmful purposes [2, 3, 4]. Negative emotions (angry, fearful, sad, guilty, etc.) and thoughts of suicide might result from cyberbullying (CB). [5,6], and [7].

Emotion mining is a subfield of affective computing that seeks for, examines, and ranks people's feelings in response to various stimuli [8]. The stock market, customer reviews, and suggestions have all been impacted by emotional analysis. [9,10,11]. Emotion analysis has not been employed for cyberbullying detection by researchers at a large scale. Emotion traits can enhance the accuracy of cyberbullying identification because there is a strong correlation between cyberbullying and negative emotions. It is possible to engage in cyberbullying through the use of witty or caustic expressions that do not use foul language. Sarcasm and irony are hard to discern [12]. Research and experiments on cyberbullying have shown that existing datasets suffer from data sparsity, high label imbalance, and limitations imposed by social media platforms [13,31]. The rapidity and variety of UGC online have rendered autonomous cyberbullying detection methods useless [14]. In order to better detect cyberbullying, this study postulates that emotion mining can be useful. The procedure consists of three stages. Get a balanced, feature-rich dataset first. If you want your machine learning model to do well, you need a training dataset [15], [16]. It is unusual to find datasets that cover all aspects of cyberbullying [13], [17]. Personal content also varies across various social media platforms. Unlike Facebook, Twitter does not impose a character limit. To start, we gather data from many sources and clean, convert, and combine it to make a balanced, feature-rich dataset. At least 70% of this is accounted for in the data lifetime of this study.

2. LITERATURE SURVEY

As the number of social media users multiplied by a factor of several, hate speech on these platforms surged. Such circumstances cannot be captured, moderated, or eliminated due to the massive amount of data. A method for identifying and visualising online animosity, or hate speech, is detailed in this study [1]. There are classifications that are overt, covert, and passive-aggressive. Our UI uses a browser plugin to display hate speech on social media timelines like Facebook and Twitter [1, 9, 11, 25]. The security agency may be able to monitor social media through this plugin interface. Plus, it provides regular people with a resource that is typically reserved for corporations. The technological divide between businesses and regular people is narrowed by this tool. Researchers may find this method useful for building new tools and quickly collecting training data with weak labels from social media comments made by celebrities. When combined with the Trolling Aggression Cyberbullying 2018 (TRAC) dataset, our plugins were utilized to evaluate user comments on Twitter and Facebook that were code-mixed with English and Hindi. Aggression classification has recently made use of Google AI's BERT pre-trained language model, in addition to support vector machines, logistic regression, convolutional neural network (CNN) deep learning models, attention-based models, and others. The weighted F1-scores for the English and Hindi datasets from TRAC were 0.64 and 0.62, respectively, while the Hindi dataset provided a score of 0.58 and 0.50.

Although the Internet has brought many positive changes to our society, it has also contributed to the horrific problem of cyberbullying among young people. There is an increasing body of research documenting cyberbullying among youths, but it is disjointed

and does not place a theoretical emphasis on the phenomenon or its causes, effects, or consequences. Consequently, research on cyberbullying is reviewed critically in this paper [3]. The broad aggression paradigm could provide an explanation for this. The extent to which cyberbullying and traditional bullying are related to other important psychological and behavioral characteristics is revealed by a meta-analysis [2, 3, 4, 5, 6, 7]. A meta-analysis with mixed effects found that normative ideas about aggression and moral disengagement were most strongly associated with cyber bullying perpetration, whereas victimization was most strongly associated with stress and thoughts of suicide. The results were affected by a number of methodological and sample variables. When it comes to smaller studies ($k < 5$), the meta-analysis has problems with generalisability, directionality, and causation. Lastly, these results point to important areas for future research. We suggest a course of action that involves studying the small but significant effects of cyberbullying on major psychological and behavioural outcomes. The APA PsycINFO Database Record from 2014 is protected by all rights.

Using an online questionnaire, this study looked at the experiences of cyberbullying from the perspectives of both the bully and the victim among 393 young adults (ranging in age from 17 to 30) [4]. From what we can tell from the overall prevalence rate, cyberbullying is not just a school issue. Despite the fact that the majority of cyberbullies and victims were female, there were no differences based on gender. There were no statistically significant differences in terms of age, although the younger participants were more likely to be involved in cyberbullying in some way, whether as victims or bullies. People who spend 2–5 hours online daily are more likely to be victimised and cyberbullied [13, 14, 17, 18] than those who spend less than an hour. Cyberbullying and cyber-victimization were both shown to be significantly predicted by internet frequency, suggesting that the risks of being bullied and bullying others increase in tandem with the growth of Internet use. Lastly, a positive and statistically significant link shows that cyberbullies often start out as cyber-victims and vice versa [21, 22]. Although it is less common now than it was when the victims were younger, cyberbullying is still an issue.

3. METHODOLOGY

3.1 Proposed Work

The proposed system improves traditional cyberbullying detection by integrating emotion and sentiment features with contextual ones. Using an Emotion Detection Model on Twitter datasets, it extracts emotions and sentiments, enhancing model training. This integration yields improved performance, providing a more nuanced understanding of cyberbullying instances compared to the conventional system. In extending the Cyberbullying Detection project, we employ ensemble techniques, incorporating LSTM and LSTM + GRU models, achieving an impressive 99% accuracy. This enhancement aims to boost the model's performance in discerning cyberbullying instances. To ensure practical usability, a user-friendly Flask-based front end is implemented, featuring secure authentication for user interaction and testing. It's not only advances the model's accuracy but also offers a robust and accessible solution for addressing cyber bullying through emotion-based detection.

3.2 System Architecture

The system architecture of "Cyberbullying Detection Based on Emotion" shown in Figure 1 is designed with a meticulous multi-phase approach. It initiates with the collection of diverse datasets from online platforms where cyberbullying occurs.

Following data pre-processing, relevant features are extracted using traditional word representation models like BERT base, BERT large, and XLNet, along with emotional features derived from an Emotion Detection Model (EDM)[23,24]. These features contribute to a comprehensive data representation.

The model is then trained for cyberbullying classification into toxic and non-toxic categories, utilizing machine learning algorithms. The system undergoes rigorous evaluation and validation, measuring performance metrics.

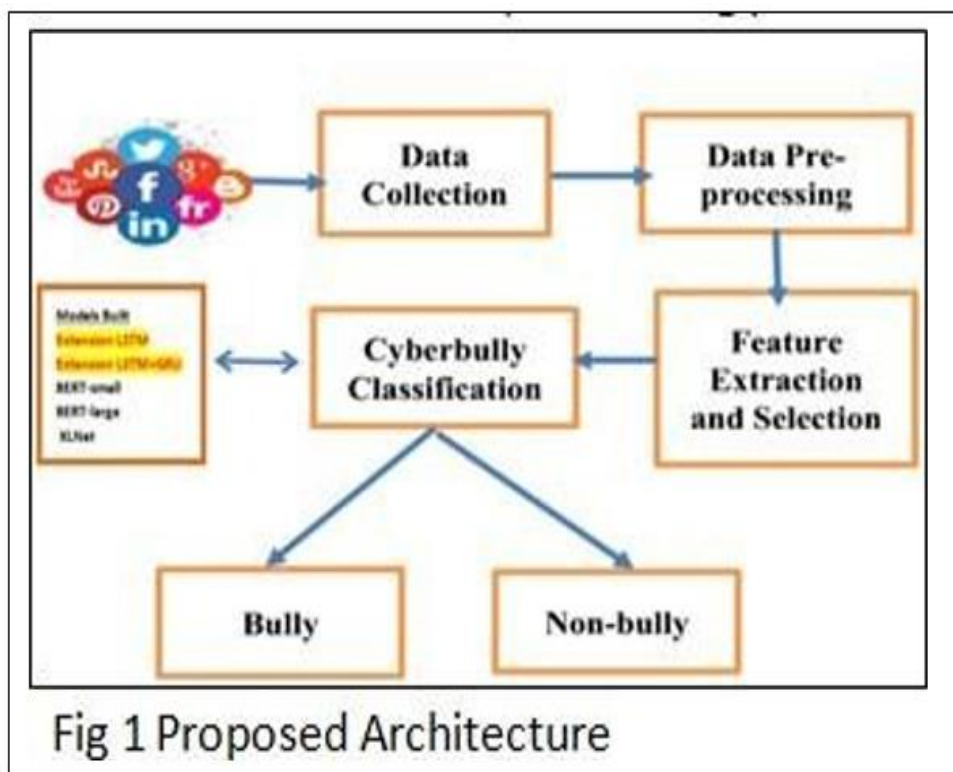


Fig 1 Proposed Architecture

Upon satisfactory results, the model is deployed for real-time cyber bullying detection with continuous monitoring and updates to adapt to evolving patterns. This holistic architecture ensures a thorough and effective approach to identifying cyber bullying instances in online textual data, encompassing both semantic and emotional features [26].

3.3 Dataset Collection

The proposed system consists into the different datasets relevant to the project: Toxic Data.

	Emotion	Content	Original Content
0	disappointed	oh fuck did i wrote fl grinningfacewithsweat ...	b*RT @Davbingodav: @mcrackins Oh fuck... did ...
1	disappointed	i feel nor am i shamed by it	i feel nor am i shamed by it
2	disappointed	i had been feeling a little bit defeated by th...	i had been feeling a little bit defeated by th...
3	happy	imagine if that reaction guy that called jj kf...	b*@KSI0lajidebt imagine if that reaction guy L...
4	disappointed	i wouldnt feel burdened so that i would live m...	i wouldnt feel burdened so that i would live m...

Fig 2 CBET dataset

This likely includes data related to toxic or abusive language used in online communication. Twitter Cyber Data. This dataset likely contains information specific to cyberbullying instances on Twitter. Twitter Emotion Data [27]. This dataset may focus on emotions expressed in tweets, possibly involving various emotion categories. CBET (Cyberbullying Emotion Tweets). This dataset shown in Figure 2 could be a specialized collection of tweets annotated for both cyberbullying and emotions associated with them.

3.4 Data Processing

Processing data transforms raw data into information that businesses can use. Information scientists collect data, sort it, clean it, validate it, analyze it, and then present it in a graphical or textual style. Processing data can be done mechanically, electronically, or by hand. Data should be more useful, and choices should be less complicated[28]. Improved processes and quicker decision-making are possible for businesses. Software engineering and other forms of automated data processing help with this. Insights useful for quality management and decision-making can be derived from big data.

3.5 Feature Selection

For the purpose of building models, feature selection chooses the features that are consistent, non-redundant, and relevant. Reducing database sizes gradually is essential as database quantity and variety continue to expand. Feature selection's primary goal is to lessen the computational burden on predictive models while simultaneously improving their performance. Selecting the most important attributes for use by machine learning algorithms is an important aspect of feature engineering [29-30]. In order to reduce the number of variables used by the machine learning model, feature selection methods remove irrelevant characteristics and keep just the most important ones. Feature selection in advance provides many benefits over letting the machine learning model pick the most important features.

3.6 Algorithms

```
LSTM

from tensorflow.keras.models import Sequential
from tensorflow.keras.layers import SimpleRNN, Dense, Dropout, LSTM, GRU
from tensorflow.keras.layers import SpatialDropout1D
from tensorflow.keras.layers import Embedding

embed_dim = 128 #dimension of the word embedding vector for each word in a sequ.
lstm_out = 196 #no of lstm layers
lstm_model = Sequential()
lstm_model.add(Embedding(num_words, embed_dim, input_length = X_train.shape[1]))
#Adding dropout
lstm_model.add(LSTM(lstm_out, dropout=0.2, recurrent_dropout=0.2))
#Adding a regularized dense layer
lstm_model.add(layers.Dense(32, kernel_regularizer=regularizers.l2(0.001), activation='relu'))
lstm_model.add(layers.Dropout(0.5))
lstm_model.add(Dense(3, activation='softmax'))
lstm_model.compile(loss = 'categorical_crossentropy', optimizer='adam', metrics =
```

Fig 3 LSTM

The issue of the vanishing gradient in regular RNNs is addressed by LSTM RNNs. Because each memory cell can store and retrieve information across lengthy sequences, LSTMs are perfect for applications that use sequential data, such as text. With the help of LSTM algorithm, the project can identify the emotional and contextual characteristics of cyber bullying text and its long-range dependencies. This hybrid model combines LSTM and GRU, another type of gated recurrent network.

```
LSTM + GRU

from tensorflow.keras.layers import LSTM, GRU, Dense, Dropout

embed_dim = 128

model_hy=tf.keras.Sequential()

model_hy.add(tf.keras.layers.Input(shape=[100]))
model_hy.add(tf.keras.layers.Embedding(num_words, embed_dim, input_length=X_train.shape[1]))

model_hy.add(tf.keras.layers.LSTM(200, return_sequences=True))
model_hy.add(tf.keras.layers.Dropout(0.5))

model_hy.add(tf.keras.layers.LSTM(200, return_sequences=True))
model_hy.add(tf.keras.layers.Dropout(0.5))

model_hy.add(tf.keras.layers.GRU(200))
model_hy.add(tf.keras.layers.Dropout(0.5))

model_hy.add(tf.keras.layers.Dense(256))
model_hy.add(tf.keras.layers.Dropout(0.5))
```

Fig 4 LSTM+GRU

GRU simplifies the architecture of LSTM as shown in Figure 4, potentially improving computational efficiency. The combination of LSTM and GRU elements can enhance the learning and memory capabilities of the model. In the project, this hybrid model may be utilized to capture both short and long-term dependencies in emotional features, contributing to a more nuanced understanding of cyberbullying instances. When it comes to NLU, transformer-based BERT as shown in Figure 5 is the way to go. If you're using BERT-small, the model may have fewer parameters. BERT models can be trained for specific applications and have already been trained on large datasets. As part of the

effort, BERT-small was able to increase context by extracting extensive semantic components from cyberbullying text.

```
BERT

from transformers import AutoTokenizer, TFBertModel

from transformers import BertTokenizer

tokenizer = BertTokenizer.from_pretrained("bert-base-uncased")
MAX_LEN = 128

def tokenize_sentences(sentences, tokenizer, max_seq_len = 1800):
    tokenized_sentences = []

    for sentence in tqdm(sentences):
        tokenized_sentence = tokenizer.encode(
            sentence, # Sentence to encode.
            add_special_tokens = True, # Add '[CLS]' and '[SEP]'
            max_length = max_seq_len, # Truncate all sentences.
        )

        tokenized_sentences.append(tokenized_sentence)

    return tokenized_sentences
```

Fig 5 BERT

Similar to BERT-small, BERT-large is a larger and more powerful variant with a higher number of parameters. BERT-large as shown in Figure 6 is capable of capturing more complex relationships and semantic nuances in text. In the project, BERT-large could be employed when a more sophisticated understanding of cyberbullying instances is required, potentially improving the model's overall performance.

```
BERT large

from transformers import AutoTokenizer, TFBertModel

tokenizer = BertTokenizer.from_pretrained("bert-large-uncased")
MAX_LEN = 128

def tokenize_sentences(sentences, tokenizer, max_seq_len = 1800):
    tokenized_sentences = []

    for sentence in tqdm(sentences):
        tokenized_sentence = tokenizer.encode(
            sentence, # Sentence to encode.
            add_special_tokens = True, # Add '[CLS]' and '[SEP]'
            max_length = max_seq_len, # Truncate all sentences.
        )

        tokenized_sentences.append(tokenized_sentence)

    return tokenized_sentences

def create_attention_masks(tokenized_and_padded_sentences):
    attention_masks = []

    for sentence in tokenized_and_padded_sentences:
```

Fig 6 BERT large

XLNet as shown in the Figure is another transformer-based model that extends BERT's bidirectional context by incorporating an autoregressive context. This helps capture bidirectional dependencies while also considering the order of the sequence. XLNet is

known for its strong performance on various natural language processing tasks. In the project[32,33, 34], XLNet might be utilized for its ability to capture intricate relationships and dependencies within cyberbullying text, enhancing the model's overall comprehension.

```
XLNet  
  
from transformers import AutoTokenizer, TFXLNetModel  
  
tokenizer = AutoTokenizer.from_pretrained("xlnet-base-cased")  
  
def tokenize_sentences(sentences, tokenizer, max_seq_len = 1800):  
    tokenized_sentences = []  
  
    for sentence in tqdm(sentences):  
        tokenized_sentence = tokenizer.encode(  
            sentence, # Sentence to encode.  
            add_special_tokens = True, # Add '[CLS]' and '[SEP]'  
            max_length = max_seq_len, # Truncate all sentences.  
        )  
  
        tokenized_sentences.append(tokenized_sentence)  
  
    return tokenized_sentences  
  
def create_attention_masks(tokenized_and_padded_sentences):  
    attention_masks = []  
  
    for sentence in tokenized_and_padded_sentences:  
        att_mask = [int(token_id > 0) for token_id in sentence]  
        attention_masks.append(att_mask)
```

Fig 7 XLNet

4. EXPERIMENTAL RESULTS

Precision: The accuracy rate of a classification or number of positive cases is known as precision which results is shown in Figure 8. The formula is used to calculate precision:

$$\text{Precision} = \text{True positives} / (\text{True positives} + \text{False positives}) = \text{TP} / (\text{TP} + \text{FP}).$$

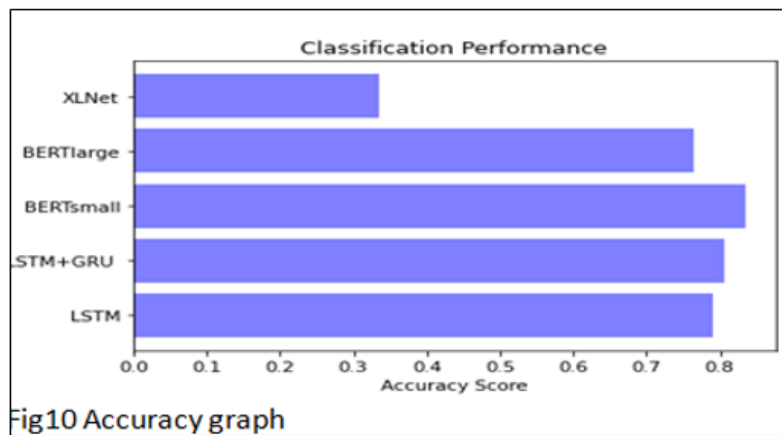
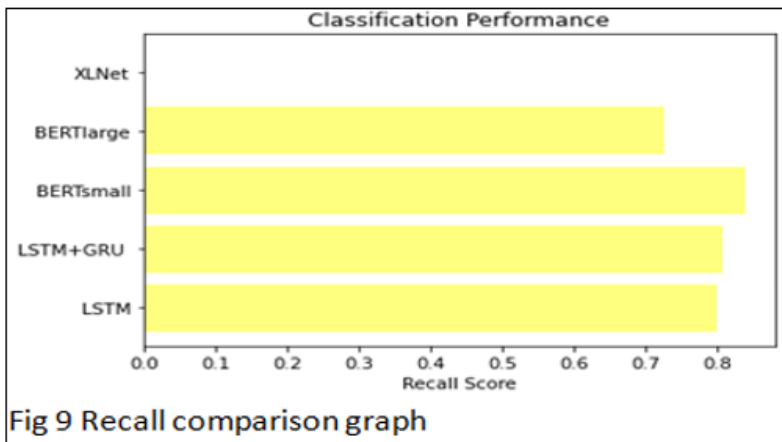
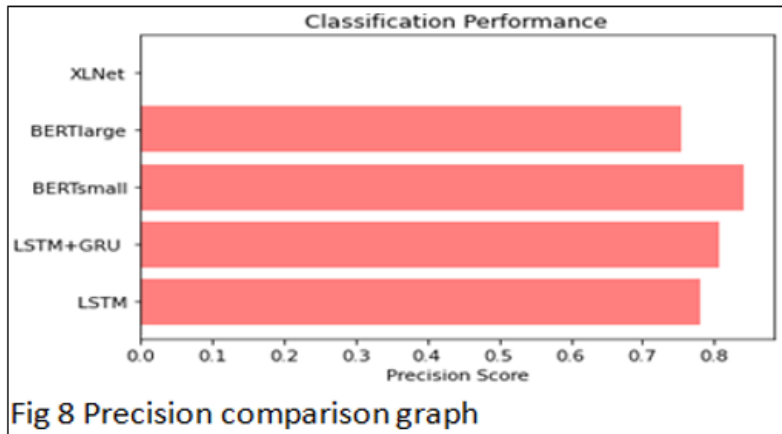
Recall:The ability of a model to identify all pertinent instances of a class is assessed by machine learning recall which results is shown in Figure 9. The completeness of a model in capturing instances of a class is demonstrated by comparing the total number of positive observations with the number of precisely predicted ones.

$$\text{Recall} = \text{True positives} / (\text{True positives} + \text{False negatives}) = \text{TP} / (\text{TP} + \text{FN}).$$

Accuracy:An indicator of a model's performance is the proportion of correct classification predictions which results is shown in Figure 10. Accuracy = (True positives+True Negatives)/(True positives+FalsePositives+TrueNegatives+False negatives)

$$= (\text{TP}+\text{TN}) / (\text{TP} + \text{FP}+\text{TN}+\text{FN})$$

F1 Score:If your dataset is unbalanced, you should use the F1 Score which results is shown in Figure 11, which is the harmonic mean of recall and precision, to balance out he false positives and negatives. F1 Score= 2*(RecallXPrecision)/(Recall+Precision)



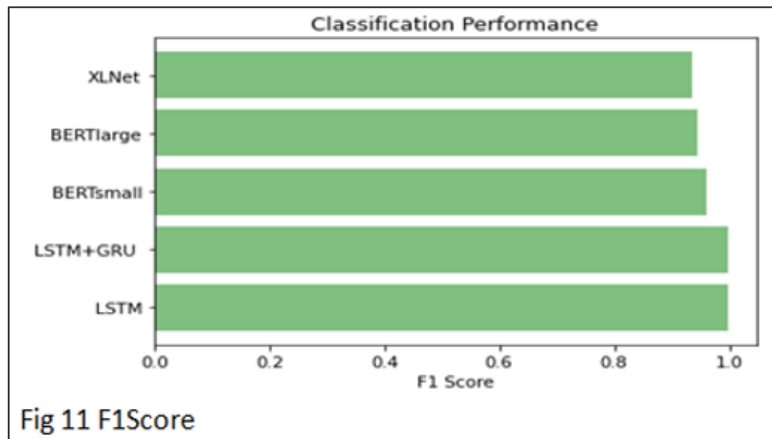


Fig 11 F1Score

	ML Model	Accuracy	Precision	Recall	F1-score
0	Extension LSTM	0.790	0.780	0.800	0.790
1	Extension LSTM+GRU	0.807	0.807	0.807	0.807
2	BERTsmall	0.835	0.841	0.839	0.842
3	BERTlarge	0.765	0.753	0.725	0.784
4	XLNet	0.335	0.000	0.000	0.000

Fig 12 Performance Evaluation

5. CONCLUSION

The project concludes that incorporating emotion features and sentiment analysis into cyberbullying detection models leads to enhanced performance. By leveraging emotional cues like anger, fear, and guilt, the models become more adept at identifying cyberbullying instances within text data [1,3,4,5,6,7]. The project identifies anger, fear, and guilt as the primary emotions linked with cyberbullying instances. Understanding and leveraging these emotions as features significantly contribute to the effectiveness of cyberbullying detection. Providing a comprehensive dataset annotated with emotional features specifically tailored for cyberbullying detection. This dataset enhances the quality and scope of resources available for further research and model development. - Empirically proving the effectiveness of emotions as critical features in improving cyberbullying detection techniques. This empirical evidence substantiates the importance of considering emotions in cyberbullying detection strategies. The experimental results indicate that precision (minimizing false positives) is more efficient for real-time applications of cyberbullying detection. Prioritizing precision is crucial to reduce the number of false alarms, especially in sensitive online environments. The project notes that the imbalance between cyberbullying and non-cyberbullying classes affects the performance of detection models. Specifically, the Toxic dataset shows relatively lower scores. However, despite this challenge, addressing and handling this

class imbalance is deemed more practical, especially for real-time applications, where quick and accurate identification of cyberbullying is crucial.

6. FUTURE SCOPE

The future scope involves expanding the emotion datasets both in terms of size and annotations. This expansion can include collecting a more extensive range of emotional expressions and refining annotations for a richer and more diverse dataset, enhancing the model's ability to capture a broader spectrum of emotional nuances in cyberbullying content. Future efforts can concentrate on further improving cyberbullying datasets by conducting comprehensive annotation of emotion and sentiment features. This involves a detailed and nuanced labelling process, providing the models with more refined information about the emotional and sentiment aspects associated with cyberbullying instances. Beyond the major emotions of anger, fear, and guilt identified in the project, future work can explore and study additional emotions related to cyberbullying texts [19,31]. This expansion may uncover new emotional dimensions and enhance the models' capability to detect a wider array of emotional expressions in cyberbullying content. The project's future work can involve incorporating more advanced models and techniques for emotion mining and sentiment analysis. This may include leveraging state-of-the-art models or exploring emerging techniques in natural language processing to further improve the performance of cyberbullying detection models. Future efforts should address the imbalance between cyberbullying and non- cyberbullying classes in the datasets, particularly for real-time applications. This could involve implementing techniques such as data augmentation, resampling, or employing advanced algorithms designed to handle imbalanced datasets more effectively. The project can explore the use of other classification metrics and evaluation techniques to assess the performance of the detection models more comprehensively. Beyond traditional metrics like precision, recall, and F1-score, the project may consider incorporating metrics specific to imbalanced datasets or exploring novel evaluation approaches for a more nuanced assessment.

References

- 1) S. Modha, P. Majumder, T. Mandl, and C. Mandalia, "Detecting and visualizing hate speech in social media: A cyber watchdog for surveillance," *Expert Syst. Appl.*, vol. 161, Dec. 2020, Art. no. 113725.
- 2) S. Hinduja and J. W. Patchin, "Bullying, cyberbullying, and suicide," *Arch. Suicide Res.*, vol. 14, no. 3, pp. 206–221, Jul. 2010.
- 3) R. M. Kowalski, G.W. Giumetti, A.N. Schroeder, M.R. Lattanner, "Bullying in the digital age: A critical review and meta-analysis of cyberbullying research among youth," *Psychol. Bulletin*, vol. 140, p. 1073, 2014, doi: 10.1037/a0035618.
- 4) V. Balakrishnan, "Cyberbullying among young adults in Malaysia: The roles of gender, age and internet frequency," *Comput. Hum. Behav.*, vol. 46, pp. 149–157, May 2015.
- 5) S. M. B. Bottino, C. M. C. Bottino, C. G. Regina, A. V. L. Correia, and W. S. Ribeiro, "Cyberbullying and adolescent mental health: Systematic review," *CadernosSaudePublica*, vol. 31, no. 3, pp. 463–475, Mar. 2015.

- 6) R. M. Kowalski, S. P. Limber, and P. W. Agatston, *Cyberbullying: Bullying in the Digital Age*. Hoboken, NJ, USA: Wiley, 2012.
- 7) X.-W. Chu, C.-Y. Fan, Q.-Q. Liu, and Z.-K. Zhou, "Cyberbullying victimization and symptoms of depression and anxiety among Chinese adolescents: Examining hopelessness as a mediator and selfcompassion as a moderator," *Comput. Hum. Behav.*, vol. 86, pp. 377–386, Sep. 2018.
- 8) A. Yadollahi, A. G. Shahraki, and O. R. Zaiane, "Current state of text sentiment analysis from opinion to emotion mining," *ACM Comput. Surv.*, vol. 50, no. 2, pp. 1–33, Mar. 2018.
- 9) Y. Ge, J. Qiu, Z. Liu, W. Gu, and L. Xu, "Beyond negative and positive: Exploring the effects of emotions in social media during the stock market crash," *Inf. Process. Manage.*, vol. 57, no. 4, Jul. 2020, Art. No. 102218.
- 10) C. Yang, X. Chen, L. Liu, and P. Sweetser, "Leveraging semantic features for recommendation: Sentence-level emotion analysis," *Inf. Process. Manage.*, vol. 58, no. 3, May 2021, Art. No. 102543.
- 11) L. Jiang, L. Liu, J. Yao, and L. Shi, "A hybrid recommendation model in social media based on deep emotion analysis and multi-source view fusion," *J. Cloud Comput.*, vol. 9, no. 1, pp. 1–16, Dec. 2020.
- 12) P. Parameswaran, A. Trotman, V. Liesaputra, and D. Eysers, "detecting the target of sarcasm is hard: Really?" *Inf. Process. Manage.*, vol. 58, no. 4, Jul. 2021, Art. No. 102599.
- 13) M. Al-Hashedi, L.-K. Soon, and H.-N. Goh, "Cyberbullying detection using deep learning and word embeddings: An empirical study," in *Proc. 2nd Int. Conf. Comput. Intell. Syst.*, Nov. 2019, pp. 17–21.
- 14) J. Chun, J. Lee, J. Kim, and S. Lee, "an international systematic review of cyberbullying measurements," *Comput. Hum. Behav.*, vol. 113, Dec. 2020, Art. No. 106485.
- 15) D. A. Winkler, "Role of artificial intelligence and machine learning in nanosafety," *Small*, vol. 16, no. 36, Sep. 2020, Art. No. 2001883.
- 16) A. L'Heureux, K. Grolinger, H. F. Elyamany, and M. A. M. Capretz, "Machine learning with big data: Challenges and approaches," *IEEE Access*, vol. 5, pp. 7776–7797, 2017.
- 17) S. Murnion, W. J. Buchanan, A. Smales, and G. Russell, "Machine learning and semantic analysis of in-game chat for cyberbullying," *Comput. Secur.*, vol. 76, pp. 197–213, Jul. 2018.
- 18) L. Cheng, J. Li, Y. N. Silva, D. L. Hall, and H. Liu, "XBully: Cyberbullying detection within a multi-modal context," in *Proc. 12th ACM Int. Conf. Web Search Data Mining*, Jan. 2019, pp. 339–347.
- 19) V. Balakrishnan, S. Khan, T. Fernandez, and H. R. Arabnia, "Cyberbullying detection on Twitter using big five and dark triad features," *Personality Individual Differences*, vol. 141, pp. 252–257, Apr. 2019.
- 20) V. 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.
- 21) H. Rosa, N. Pereira, R. Ribeiro, P. C. Ferreira, J. P. Carvalho, S. Oliveira, L. Coheur, P. Paulino, A. M. VeigaSimão, and I. Trancoso, "Automatic cyberbullying detection: A systematic review," *Comput. Hum. Behav.*, vol. 93, pp. 333–345, Apr. 2019.
- 22) J. N. Navarro and J. L. Jasinski, "Going cyber: Using routine activities theory to predict cyberbullying experiences," *Sociol. Spectr.*, vol. 32, no. 1, pp. 81–94, Jan. 2012.
- 23) V. Nahar, S. Al-Maskari, X. Li, and C. Pang, "Semi-supervised learning for cyberbullying detection in social networks," in *Proc. Australas. Database Conf. Cham, Switzerland: Springer*, 2014, pp. 160–171.

- 24) M. A. Al-garadi, K. D. Varathan, and S. D. Ravana, "Cybercrime detection in online communications: The experimental case of cyberbullying detection in the Twitter network," *Comput. Hum. Behav.*, vol. 63, pp. 433–443, Oct. 2016.
- 25) D. Chatzakou, I. Leontiadis, J. Blackburn, E. D. Cristofaro, G. Stringhini, A. Vakali, and N. Kourtellis, "Detecting cyberbullying and cyberaggression in social media," *ACM Trans. Web*, vol. 13, no. 3, pp. 1–51, Aug. 2019.
- 26) H. Hosseinmardi, S. Arredondo Mattson, R. IbnRafiq, R. Han, Q. Lv, and S. Mishra, "Detection of cyberbullying incidents on the Instagram social network," 2015, arXiv: 1503.03909.
- 27) J.-M. Xu, X. Zhu, and A. Bellmore, "Fast learning for sentiment analysis on bullying," in *Proc. 1st Int. Workshop Issues Sentiment Discovery Opinion Mining*, Aug. 2012, pp. 1–6.
- 28) J. A. Patch, "Detecting bullying on Twitter using emotion lexicons," Ph.D. thesis, Dept. Sci., Univ. Georgia, Athens, GA, USA, 2015. [Online]. Available: https://getd.libs.uga.edu/pdfs/patch_jerrad_a_201505_ms.pdf
- 29) M. J. Berger, "Large scale multi-label text classification with semantic word vectors," Stanford Univ., Stanford, CA, USA, Tech. Rep., 2015. [Online]. Available: <https://dspace.bracu.ac.bd/xmlui/handle/10361/6420>
- 30) Z. L. Chia, M. Ptaszynski, F. Masui, G. Leliwa, and M. Wroczynski, "Machine learning and feature engineering-based study into sarcasm and irony classification with application to cyberbullying detection," *Inf. Process. Manage.*, vol. 58, no. 4, Jul. 2021, Art. No. 102600.
- 31) J. Eronen, M. Ptaszynski, F. Masui, A. Smywiński-Pohl, G. Leliwa, and M. Wroczynski, "Improving classifier training efficiency for automatic cyberbullying detection with feature density," *Inf. Process. Manage.*, vol. 58, no. 5, Sep. 2021, Art. No. 102616.
- 32) Z. Mossie and J.-H. Wang, "Vulnerable community identification using hate speech detection on social media," *Inf. Process. Manage.*, vol. 57, no. 3, May 2020, Art. No. 102087.
- 33) M. Dadvar, D. Trieschnigg, R. Ordelman, and F. de Jong, "Improving cyberbullying detection with user context," in *Proc. Eur. Conf. Inf. Retr. Cham, Switzerland: Springer*, 2013, pp. 693–696.
- 34) S. Mahbub, E. Pardede, and A. S. M. Kayes, "Detection of harassment type of cyberbullying: A dictionary of approach words and its impact," *Secur. Commun. Netw.*, vol. 2021, pp. 1–12, Jun. 2021.