

# **Analyzing the Impact of Hate Speech Using NLP Techniques**

**Enroll. No – 22103177, 22103302, 22103061**

**Name of Student – Dev Agarwal, Alokik Garg, Swapnil Pandey**

**Name of Supervisor – Dr Sayani Ghosal**



**December-2024**

**Submitted in partial fulfilment of the degree of Bachelors of Technology**

**In**

**Computer Science Engineering**

**DEPARTMENT OF COMPUTER SCIENCE ENGINEERING &  
INFORMATION TECHNOLOGY**

**JAYPEE INSTITUTE OF INFORMATION TECHNOLOGY, NOIDA**

(I)

**TABLE OF CONTENTS**

<b>Chapter No.</b>	<b>Topics</b>	<b>Page No.</b>
<b>Chapter - 1</b>	<b>Introduction</b>	<b>8-12</b>
	1.1 General introduction	8
	1.2 Problem statement	8
	1.3 Significance/Novelty of the problem	9
	1.4 Empirical study	9
	1.5 Brief description of solution approach	9-10
	1.6 Comparison of existing approaches of the problem	10-12
<b>Chapter - 2</b>	<b>Literature Survey</b>	<b>13-22</b>
	2.1 Summary of papers studied	13-21
	2.2 Integrated summary of the literature studied	21-22
<b>Chapter - 3</b>	<b>Requirement Analysis and Solution Approach</b>	<b>22-28</b>
	3.1 Overall Description of the Project:	22-25
	I. Core Components of the Project:	22-
	3.2 Development Process:	26-28
	3.2.1 Functional Requirements:	26-27
	3.2.2 Non-Functional Requirements	27
	3.3 Solution Approach:	27-28
<b>Chapter - 4</b>	<b>Modelling and Implementation Design</b>	<b>29-38</b>
	4.1 Design Diagram	29
	4.2 Implementation Details	29-37
	4.2.1 Multilabel Classification:	29-32

	4.2.2 Multiclass Classification:	32-37
	4.2.3 Combined Embeddings – BERT and GPT-Neo	38
	4.3 Risk mitigation and analysis	38
<b>Chapter – 5</b>	<b>Testing</b>	<b>39-40</b>
	5.1 Testing Plan	39
	5.2 Component Testing	39
	5.3 Dataset Validation	39
	5.4 Model Evaluation	39
	5.5 Test Results	40
<b>Chapter – 6</b>	<b>Conclusion</b>	<b>41-42</b>
	6.1 Findings	41
	6.2 Conclusions	41
	6.3 Future Scope	42
	<b>References</b>	<b>43-44</b>

**(II)**

**DECLARATION**

I/We **Dev Agarwal, Alokik Garg and Swapnil Pandey** declare that this submission is my/our own work and that, to the best of my knowledge and belief, it contains no material previously or written by another person nor material which has been accepted for the award of any other degree or diploma of the university or other institute of higher learning, except where due acknowledgement has been made in the text

Place: Noida, India

Date: Nov 20, 2024

Student Name	Enrollment No.	Signature
Alokik Garg	22103302	
Dev Agarwal	22103177	
Swapnil Pandey	22103061	

(III)

## CERTIFICATE

This is to certify that the work titled “**Analyzing the Impact of Hate Speech Using NLP Techniques**” submitted by Dev Agarwal, Alokik Garg and Swapnil Pandey in partial fulfillment for the award of degree of B.Tech in Computer Science of Jaypee Institute of Information Technology, Noida has been carried out under my supervision. This work has not been submitted partially or wholly to another University or Institute for the award of this or other degree or diploma

Signature of Supervisor:

Name of Supervisor: Dr. Sayani Ghosal

Designation: Assistant Professor (Senior Grade)

Date: Nov 20, 2024

(IV)

## ACKNOWLEDGEMENT

We would like to express our sincere gratitude to all those who have contributed to the successful completion of this project. We would like to convey our heartfelt gratitude to Dr. Sayani Ghosal, who provided us with this wonderful opportunity to work on this project and guided us in its completion. We learned many things not just theoretically but the implementation too from the very basics, including how to work with Embeddings and ML Models(GPT, Bert, etc), Libraries like Transformers(Hugging Face) and Pytorch.

Date: Nov 20, 2024

Student Name	Enrollment No.	Signature
Alokik Garg	22103302	
Dev Agarwal	22103177	
Swapnil Pandey	22103061	

## (V)

### SUMMARY

1. **Hate Speech Impact Classification:** This project introduces a comprehensive model for analyzing and classifying the impact of hate speech across various categories such as emotional distress, provoking violence, and individual harassment. By leveraging advanced machine learning techniques, the project offers insights into the broader consequences of online hate speech.
2. **Data Collection and Preprocessing:** The project utilizes a diverse dataset(hate-speech) sourced from Hugging Face, ensuring the analysis is grounded in real-world data. Rigorous preprocessing steps, including tokenization and label categorization, prepare the data for effective model training, addressing challenges like class imbalance and noise.
3. **Model Development and Evaluation:** Using state-of-the-art models such as BERT, GPT-Neo embeddings, along with attention-based multilabel classification and LSTM networks, the project achieves robust classification results. Models are evaluated on multiple metrics, including accuracy, precision, and recall, to ensure reliability and effectiveness in identifying hate speech impacts.
4. **Predictive Analytics for Social Impact:** The project goes beyond classification by predicting the social consequences of hate speech, providing valuable insights into its potential effects on individuals and communities. The predictive component adds significant value to understanding the real-world relevance of detected hate speech.
5. **Comparison with Existing Approaches:** This project advances beyond traditional hate speech detection systems by focusing on the specific impacts of hate speech, offering a more nuanced understanding of its effects. It bridges the gap between detection and social impact prediction, contributing to more targeted interventions.
6. **Ethical Considerations and Data Security:** Ethical considerations are prioritized throughout the project, particularly in the handling of sensitive data. Secure data practices are implemented to protect the integrity and confidentiality of the dataset and model outputs.
7. **Future Work and Continuous Improvement:** Future work includes refining the models for higher accuracy and extending the system to analyze multilingual hate speech.

**Date:** 20/11/2024

**Signature of Supervisor:**

## **Chapter – 1: Introduction**

### **1.1 General Introduction**

In response to the growing concerns surrounding online toxicity, we introduce a pioneering initiative—HateSpeech Impact Analyzer. Positioned at the intersection of technology and social responsibility, this project aims to address the far-reaching effects of hate speech on individuals and communities. By leveraging advanced machine learning techniques, the project seeks to predict potential social impact of Hate - Speech, such as emotional distress, the incitement of violence, and individual harassment. Through sophisticated models, including BERT, and GPT-Neo embeddings, this system offers valuable insights into the consequences of online harmful speech. The goal of this initiative is not only to contribute to the field of natural language processing but also to foster a safer, more responsible online environment. This project empowers both individuals and organizations to assess and address the impact of hate speech effectively, aiming to create a more informed, responsive, and empathetic digital landscape.

### **1.2 Problem Statement:**

Existing models for hate speech analysis primarily focus on detecting whether a post contains hate speech or not, often neglecting the far-reaching consequences of such speech on individuals and society. While these detection systems are valuable for identifying harmful content, they fail to provide a comprehensive understanding of the broader impacts of hate speech. Specifically, current systems do not account for the psychological toll that hate speech takes on its victims, such as emotional distress, feelings of isolation, or anxiety. Moreover, they overlook the potential for hate speech to incite violence or provoke harmful actions, nor do they address how it can contribute to individual harassment and long-term societal harm.

This limitation leaves a critical gap in understanding the severity and full implications of online hate speech. Most models treat hate speech detection in a generalized manner, focusing solely on classification without delving into the specific effects it may have on people or communities. This project addresses these shortcomings by shifting the focus from detection to understanding the deeper impact of hate speech, aiming to provide a more nuanced and contextualized perspective on its harmful effects, such as emotional distress, the provocation of violence, and individual harassment.



### 1.3 Significance/Novelty of the Problem

The proposed **Hate Speech Impact Analysis** project holds significant importance in addressing a critical gap in the current landscape of online safety and digital well-being. While existing systems focus on detecting hate speech, they often fail to understand or predict its deeper effects on individuals and communities. By shifting the focus from mere identification to analyzing the psychological and societal impacts, this project offers a novel approach to studying the far-reaching consequences of hate speech.

The significance lies in the ability to assess how online hate speech contributes to emotional distress, provokes violence, and fuels individual harassment, enabling a more comprehensive understanding of its harm. This approach not only deepens the conversation around online toxicity but also opens avenues for more targeted interventions and strategies to combat its impact. By providing insights into these consequences, the project seeks to create a safer, more informed digital environment, making it a crucial step forward in the ongoing effort to tackle the harmful effects of hate speech in the online world.

### 1.4 Empirical Study:

The **Hate Speech Impact Analysis** project implements various machine learning models to assess the real-world effects of hate speech, focusing on emotional distress, provocation of violence, and individual harassment. Using the Hate-Speech dataset from Hugging Face, the project utilizes BERT and GPT-Neo embeddings to extract features from text. These embeddings are processed through various multilabel and multiclass classification models to predict the potential impact of hate speech.

The models were trained on labelled data, addressing challenges like class imbalance and noise through preprocessing techniques. These models provide valuable insights into the impact of hate speech, offering a foundation for future applications aimed at real-time impact assessment.

### 1.5 Brief Description of the solution approach:

The **Hate Speech Impact Analysis** project uses machine learning models to assess the psychological and societal effects of hate speech. Instead of detecting hate speech, the focus is on understanding its impact, such as emotional distress, provocation of violence, and individual harassment.

Key Features:

- **Data Collection:** Utilizes the Hate-Speech dataset from Hugging Face, focusing on posts that are labelled with potential impacts.

- **Modelling Approach:** Implements BERT and GPT-Neo embeddings to process text, followed by multilabel and multiclass classification models to predict the impact of hate speech based on emotional distress caused, provocation of violence, and individual harassment.
- **Preprocessing:** Addresses data noise and class imbalance to ensure accurate predictions.
- **Future Application:** The models, once deployed, can provide real-time insights into the impact of hate speech, offering a valuable tool for individuals, organizations, and policymakers to respond effectively.

## 1.6 Comparison of the Existing Approaches to the Problem Found

Existing approaches to analyzing hate speech typically focus on detection rather than assessing its impact. While these detection systems identify harmful content, they fall short in evaluating the consequences of hate speech on individuals and communities. Below is a comparison between traditional approaches and the approach used in the **Hate Speech Impact Analysis** project:

### 1. Hate Speech Detection Systems:

- **Existing Approach:**
  - Relies on classifiers to detect whether a post contains hate speech or not.
  - Focuses primarily on identifying offensive language, without considering the effects on victims.
- **Comparison with Hate Speech Impact Analysis:**
  - **Advantages of the Proposed Approach:**
    - Unlike traditional systems, this project goes beyond detection and assesses the emotional distress, provocation of violence, and harassment caused by hate speech.
    - Provides a deeper, more nuanced understanding of the real-world consequences of hate speech, enhancing the ability to respond appropriately.

## 2. Sentiment Analysis Systems:

- **Existing Approach:**

- Sentiment analysis models typically gauge the emotional tone (positive, negative, neutral) of a text, but they do not differentiate between the specific impact of hate speech.
- These systems are unable to classify whether a text is inciting harm or causing emotional damage.

- **Comparison with Hate Speech Impact Analysis:**

- **Advantages of the Proposed Approach:**

- The project uses specialized models to classify the impact of hate speech into categories like emotional distress, provoking violence, and harassment, providing actionable insights.
    - It captures a more detailed picture of the psychological effects on individuals, offering a more specific and targeted analysis.

## 3. Manual Analysis Methods:

- **Existing Approach:**

- Relying on human annotators to manually assess the content of posts and determine their potential impact.
- This process is slow, subjective, and prone to inconsistencies.

- **Comparison with Hate Speech Impact Analysis:**

- **Advantages of the Proposed Approach:**

- The project automates the impact assessment using machine learning models, ensuring consistency and reducing human bias.
    - By processing large datasets, it can evaluate a broader range of content efficiently, providing a more comprehensive understanding of hate speech impacts.

#### 4. **General Hate Speech Categorization:**

- **Existing Approach:**

- Existing systems categorize hate speech into broad categories, such as hate speech, offensive, or neutral, but they do not delve into the specific types of harm caused by the speech.

- **Comparison with Hate Speech Impact Analysis:**

- **Advantages of the Proposed Approach:**

- This project introduces a more detailed classification system, assessing whether the speech provokes violence, causes emotional distress, or leads to harassment.
- It enhances understanding of how different types of hate speech affect individuals, offering more precise intervention strategies.

Through this comparison, it is clear that the **Hate Speech Impact Analysis** project provides a more comprehensive and insightful solution by focusing on the psychological and societal effects of hate speech rather than just its identification.

## Chapter-2 Literature Survey

### 2.1 Summary of the papers studied:

#### Research Paper 1:

[1] Zhou, Y., Yang, Y., Liu, H., Liu, X., & Savage, N. (2020). Deep learning-based fusion approach for hate speech detection. *IEEE Access*, 8, 128923-128929.

This paper addresses the increasing problem of hate speech online and the limitations of existing detection systems. Traditional manual detection methods are time-consuming, while automated systems often fail to capture the nuanced nature of hate speech. The authors propose a deep learning-based fusion approach that combines multiple classifiers to improve detection performance.

#### Section 1- Introduction:

The introduction highlights the rise of hate speech online and the inadequacy of traditional detection systems in addressing this issue. Current models struggle to detect subtle variations in hate speech due to a lack of nuance understanding. The paper emphasizes the need for more advanced models, focusing on deep learning techniques and fusion methods to improve the accuracy and context-awareness of hate speech detection.

#### Section 2 - Related Work:

In this section, the authors define hate speech and discuss its harmful societal impact, including discrimination and inequality. They then review existing detection methods, categorizing them into traditional machine learning and deep learning approaches. The paper explains key models like ELMo, BERT, CNN, and RNN, detailing their relevance in understanding the complex language of hate speech. The authors argue that while individual models perform well, combining them into a fusion model could significantly enhance classification accuracy.

#### Section 3 - Methodology:

This section introduces the framework and fusion rules for combining the classifiers. The authors discuss how different classifiers—ELMo, BERT, and CNN—are combined to enhance the detection of hate speech. The methodology also analyzes the factors influencing fusion results, aiming to optimize the performance of the combined models.

#### Section 4 - Experiments and Results Analysis:

The authors use the SemEval 2019 Task 5 dataset, which includes Twitter messages that contain hate speech directed at immigrants and women. The dataset is divided into hate (42%) and non-hate (58%) categories. The paper compares the performance of ELMo, BERT, and CNN models in

detecting hate speech, evaluating them based on accuracy and F1-score. The results show that while each model performs well on its own, the fusion of classifiers leads to a noticeable improvement in detection performance.

#### Section 5 - Conclusion:

The paper concludes that combining ELMo, BERT, and CNN classifiers using a fusion approach significantly enhances hate speech detection accuracy. The fusion method, while simple, improves classification performance at a minimal cost. However, the paper acknowledges that the fusion process could be improved by integrating classifiers earlier in the process, rather than after classification. The authors suggest that future work could explore early cooperation of classifiers for even better performance.

#### **Research Paper 2:**

**[2] Aluru, S. S., Mathew, B., Saha, P., & Mukherjee, A. (2020). Deep learning models for multilingual hate speech detection. arXiv preprint arXiv:2004.06465.**

This paper addresses the challenges and solutions for detecting hate speech in multiple languages, extending beyond English-focused datasets. It explores deep learning models for detecting hate speech in a multilingual context. The paper focuses on multilingual hate speech detection, conducting a large-scale analysis across 9 languages from 16 different sources, aiming to address limitations in current hate speech models that mostly focus on English.

#### Section 1 - Introduction:

With the rise of online social media, the dissemination of hate speech has increased rapidly, fueled by bad actors leveraging the platform for harmful purposes. The paper highlights the importance of developing multilingual hate speech detection systems, as most existing models are limited to English. The authors stress the necessity of models that can handle the complexities and nuances of various languages.

#### Section 2 - Related Work:

The paper delves into the complexities of hate speech, which intersects with freedom of expression and the rights of individuals and minorities. It reviews early hate speech detection methods such as dictionary look-ups and bag-of-words models, which were limited in scope. As datasets grew larger, more complex models like deep learning and graph embedding techniques gained traction. The paper references models such as CNN-GRU and BERT for improving detection accuracy, and notes that multilingual approaches remain an under-explored area.

### Section 3 - Dataset Description:

The authors examine 16 publicly available hate speech datasets across 9 languages: Arabic, Portuguese, English, German, Indonesian, Italian, Polish, Spanish, and French. These datasets are sourced from various platforms, providing a diverse range of hate speech examples from different cultural and linguistic contexts.

### Section 4 - Experiments:

The paper discusses the experimental setup, where the datasets for each language were combined, and a stratified train/validation/test split was used (70%/10%/20%). Three key models were tested:

- MUSE + CNN-GRU: This combines MUSE embeddings with CNN-GRU to detect hate speech.
- Translation + BERT: Input sentences are translated into English and then fed into the BERT model.
- LASER + LR: LASER embeddings are used as input to a Logistic Regression model.

These models were evaluated based on their F1-scores and accuracy across various languages, providing insights into their performance.

### Section 5 - Results:

The results compare the performance of the models in both monolingual and multilingual scenarios. The authors found that the multilingual models, especially those using embeddings like MUSE and LASER, showed improved performance compared to traditional methods. They also discuss how models like BERT performed well in handling language translation, while CNN-GRU demonstrated strong results when combined with MUSE embeddings.

### Section 6 - Discussion and Error Analysis:

The authors discuss interpretability in hate speech models, exploring how model decisions can be better understood and trusted. They also provide an error analysis, focusing on common mistakes made by the models, such as confusion between subtle hateful and non-hateful speech. Suggestions for addressing these errors include refining embeddings and improving model training across multiple languages.

### Section 7 - Conclusion:

The paper concludes that multilingual hate speech detection is an essential but challenging problem. While the models used show promising results, there is room for improvement in terms of accuracy

and interpretability. The authors suggest that more robust models could be developed by focusing on better language-specific features and increasing the quality of training datasets.

### **Research Paper 3:**

**[3] Das, M., Saha, P., Mathew, B., & Mukherjee, A. (2022). Hatecheckhin: Evaluating hindi hate speech detection models. arXiv preprint arXiv:2205.00328.**

This paper addresses the challenges of detecting hate speech in Hindi, a low-resource language. The authors introduce Hatecheckhin, a framework to evaluate Hindi hate speech detection models. By comparing existing models, the study assesses their performance on various datasets and highlights the importance of better evaluation tools for non-English languages.

#### **Section 1 - Introduction:**

The rise of online hate speech has led to growing interest in detection systems, but most research focuses on high-resource languages like English. There is a gap in robust models for low-resource languages, especially Hindi. The authors present Hatecheckhin, a framework for evaluating Hindi hate speech detection systems.

#### **Section 2 - Related Work:**

The authors review existing research on hate speech detection, emphasizing models for English. While there have been significant strides in detecting hate speech in high-resource languages, there is limited research and fewer resources for Hindi. The paper reviews previous methods, including classical machine learning models and deep learning-based models.

#### **Section 3 - Dataset Description:**

The authors use publicly available Hindi datasets for hate speech detection, which include a variety of social media posts. These datasets consist of labeled examples of hate speech and non-hate speech, providing a foundation for training and testing models.

#### **Section 4 - Evaluation Framework (Hatecheckhin):**

The paper introduces Hatecheckhin, a framework for evaluating Hindi hate speech detection models. It outlines metrics such as accuracy, F1-score, and precision to compare various models. The framework aims to offer a more comprehensive assessment beyond simple accuracy measures.

#### **Section 5 - Models and Experiments:**

The authors evaluate multiple models while keeping mBERT as base model, including traditional machine learning models like SVM and other deep learning models and their Hindi variants. The



models are tested on the Hindi datasets, and the paper discusses how each model performs, detailing strengths and weaknesses.

#### Section 6 - Results and Discussion:

The authors present the results of the experiments, showing that BERT-based models outperform traditional machine learning models for Hindi hate speech detection. They discuss the challenges of detecting hate speech in Hindi, including cultural and linguistic factors that influence model performance.

#### Section 7 - Conclusion:

The paper concludes that there is a significant gap in hate speech detection for Hindi, but the Hatecheckhin framework can help assess and improve detection systems for this language. The authors suggest future work focused on improving the robustness of models and expanding datasets to enhance detection accuracy.

#### **Research Paper 4:**

**[4] Vijayaraghavan, P., Larochelle, H., & Roy, D. (2021). Interpretable multi-modal hate speech detection. *arXiv preprint arXiv:2103.01616*.**

This paper proposes a multi-modal approach to hate speech detection that incorporates both textual and socio-cultural context inputs. The authors introduce an interpretable framework that improves understanding of model predictions by highlighting which aspects of text and socio-cultural background contribute to the classification of hate speech.

#### Section 1 - Introduction

The paper proposes a novel multi-modal hate speech classification model that incorporates both textual and contextual features. It aims to improve the understanding of hate speech by considering social and cultural factors related to the user alongside the textual content of tweets. The work further investigates how socio-political contexts impact the classification and interpretation of hate speech.

#### Section 2 - Dataset

##### 2.1 Publicly Available Datasets:

The paper utilizes various publicly available datasets, primarily obtained through keyword-based filtering or hashtags from Twitter. These datasets are manually annotated and contain tweets related to different forms of hate speech.

## 2.2 Data Collection:

To enhance the model, additional data is collected using data augmentation and distant supervision methods. The data collection focuses on tweets containing swear words and phrases linked to extremist groups, which helps identify hate communities and expand the dataset.

## Section 3 - Model

The paper describes a deep learning-based model designed to handle the multi-modal nature of the data.

### 3.1 Extracting Semantic Features:

The model processes textual data and extracts relevant semantic features to understand the underlying meaning of the content.

### 3.2 Extracting Cultural Context Features:

It also accounts for cultural context by extracting features that relate to the user's cultural background, enhancing the model's ability to interpret hate speech from different perspectives.

### 3.3 Extracting Social Context Features:

The model incorporates social context by analyzing the social interactions and the user's online behavior, improving its sensitivity to different forms of hate speech.

## Section 4 - Evaluation

### 4.0.1 Identifying Categories of Hate Speech:

The authors categorize hate speech into various types, including Anti-Islam, Anti-Black, General Hate, Anti-Immigrant, and Anti-Semitic, to evaluate the model's ability to classify these forms of hate speech.

### 4.1 Interpretability:

The model's interpretability is evaluated by comparing clusters generated using embeddings from two models: a Text-Only Model and a Text+Social Context (SC) Model. The comparison highlights the importance of integrating social and cultural context features for improved classification.

## Section 5 - Conclusion

The paper concludes that a multi-modal classification method with late-fusion of textual, social, and cultural features significantly improves hate speech detection. The model can better classify and

understand hate speech code words and categorize them into distinct types, offering new insights into how contextual features influence hate speech classification.

### **Research Paper 5:**

**[5] Badjatiya, P., Gupta, S., Gupta, M., & Varma, V. (2017, April). Deep learning for hate speech detection in tweets. In Proceedings of the 26th international conference on World Wide Web companion (pp. 759 760).**

This paper addresses hate speech detection on Twitter, which is important for tasks like sentiment analysis and content recommendation. The study explores deep learning architectures, including CNNs, LSTMs, and FastText, to classify tweets as racist, sexist, or neutral. The complexity of natural language is tackled using semantic word embeddings.

#### **Section 1 - Introduction**

With the growing use of Twitter, hateful content targeting individuals and groups has become widespread. Effective detection is crucial for public sentiment analysis and preventing harmful activities. The paper evaluates multiple classifiers and deep learning models to address the challenges posed by informal language on Twitter.

#### **Section 2 - Proposed Approach**

- **Baseline Methods:**

- Char n-grams: Uses character n-grams for detection.
- TF-IDF: A traditional method based on word importance.
- BoWV: Averages word embeddings to represent tweets.

- **Proposed Methods:**

- CNN: Extracts hierarchical features from tweets.
- LSTM: Captures long-range dependencies in text.
- FastText: Fine-tunes word embeddings through back-propagation during training

Each model is trained with labeled data using back-propagation and classified into categories of racist, sexist, or neutral.

### Section 3 – Experiments:

- Dataset and Experimental Settings:

The dataset contains 16,000 annotated tweets, categorized into racist, sexist, and neutral. The paper experiments with different models and embeddings to compare performance.

- Results and Analysis:

Deep learning methods outperformed traditional ones. TF-IDF was better than character n-grams, while CNN performed best among proposed models. Surprisingly, random embeddings with GBDT outperformed GloVe embeddings. The top-performing model was LSTM + Random Embedding + GBDT.

### Section 4 - Conclusion

The study highlights that CNNs and LSTMs provide better results than traditional methods for hate speech detection. The integration of learned embeddings with gradient-boosted decision trees achieved the highest accuracy. Future work will explore how user network features can further improve detection accuracy.

### Research Paper 6:

**[6]: Zimmerman, S., Kruschwitz, U., & Fox, C. (2018, May). Improving hate speech detection with deep learning ensembles. In Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018).**

The paper presents an ensemble approach for detecting hate speech using deep learning techniques, improving accuracy by combining multiple models. The authors utilize various deep learning architectures to identify offensive content in social media text, with a focus on combining classifiers to enhance performance.

### Section 1 - Introduction

Hate speech detection is a critical challenge, especially in social media. Traditional methods often fail to capture the nuanced language used in such content. This paper explores combining multiple deep learning models—such as CNNs, LSTMs, and GRUs—into ensembles to improve the overall performance of hate speech detection systems.

### Section 2 - Proposed Approach

The authors apply ensemble learning, combining several models trained on different feature sets. These models include Convolutional Neural Networks (CNNs) for extracting local patterns, Long Short-Term Memory networks (LSTMs) to capture long-range dependencies, and Gated Recurrent

Units (GRUs) for handling sequential data. The models are ensembled to leverage the strengths of each, aiming to improve accuracy over individual models.

#### Section 4 - Experiments

The paper uses datasets such as Waseem and SemEval to train the models. The results show that ensemble approaches outperform individual models, providing better generalization and robustness. The ensemble method's ability to capture various aspects of the data results in significant performance improvements compared to single-model classifiers.

#### Section 5 - Conclusion

The paper concludes that deep learning ensembles are a promising approach for improving hate speech detection. The combination of different models allows for better performance in detecting offensive content, especially in noisy and ambiguous social media data. Future work involves expanding the dataset and refining the ensemble methods for broader applicability.

## **2.2 Integrated Summary of the Literature Studied:**

The reviewed literature highlights the challenges and advancements in hate speech detection, especially on social media platforms like Twitter. The studies reveal a consistent emphasis on leveraging deep learning techniques to overcome the complexities of detecting nuanced and context-dependent hate speech.

### **1. Role of Feature Representations:**

Early methods relied on traditional feature representations such as character n-grams, TF-IDF, and Bag of Words vectors, as explored by Badjatiya et al. (2017). While effective to a degree, these approaches struggled with informal and implicit language, leading to a shift toward word embeddings like GloVe and task-specific embeddings trained via neural networks. This evolution demonstrated that deep learning-based representations often outperform traditional techniques.

### **2. Deep Learning Architectures:**

Several studies tested diverse architectures, such as CNNs, LSTMs, and FastText (Badjatiya et al., 2017; Zimmerman et al., 2018). CNNs are noted for their ability to extract local text patterns, while LSTMs and GRUs capture long-range dependencies. Ensemble learning, particularly combining these architectures, has been proven to improve model robustness and accuracy (Zimmerman et al., 2018). The incorporation of these techniques enabled significant strides in capturing the semantic and syntactic subtleties of hate speech.

### **3. Incorporation of Contextual Features:**

Beyond textual content, Vijayaraghavan et al. (2021) introduced multi-modal approaches that include socio-cultural and user-network information. This approach underscores the importance of contextualizing hate speech within its social and cultural backdrop, significantly enhancing interpretability and classification accuracy.

### **4. Augmented Datasets and Challenges:**

Researchers have emphasized the role of data collection and augmentation to address the class imbalance inherent in hate speech datasets. For example, Zimmerman et al. (2018) and Badjatiya et al. (2017) utilized curated datasets like Waseem and SemEval but acknowledged the difficulty of distinguishing hate speech from offensive language. Augmentation techniques, such as keyword-based filtering and data synthesis, were proposed to expand datasets and capture diverse hate speech expressions.

### **5. Advanced Techniques for Robust Detection:**

Recent works (e.g., Das et al., 2022) employed ensemble strategies, distant supervision, and hierarchical deep learning models to improve classification performance. These approaches demonstrated enhanced generalization across different hate speech categories and reduced sensitivity to noisy or ambiguous inputs.

### **6. Future Directions:**

Across studies, there is a shared vision for incorporating user-network features, exploring multilingual hate speech, and refining interpretability. The importance of explainable AI methods to understand decision-making processes in hate speech detection models has been emphasized as a critical next step.

Collectively, the literature underscores the evolving sophistication of hate speech detection systems, transitioning from traditional machine learning methods to robust, multi-modal deep learning approaches. These advancements pave the way for comprehensive, context-aware systems capable of addressing the growing complexities of hate speech on social media platforms.

## Chapter-3 Requirement Analysis and Solution Approach

### 3.1 Overall Description of the Project:

This project analyzes the impact of hate speech by leveraging machine learning and advanced natural language processing (NLP) techniques. Various models were implemented to classify hate speech into three distinct categories: emotional distress, provoking violence, and individual harassment. These models were trained and evaluated using embeddings derived from BERT, and GPT-Neo, focusing on their efficiency and accuracy. The primary objective is to analyze the consequences of hate speech rather than merely detecting its presence.

#### I. Core Components of the Project:

**i. Hate-Speech Dataset from Hugging Face** consisting a total of 135,556 rows. It is designed to analyze hate speech by incorporating a **hate speech score** and a set of **ordinal labels**. The **primary outcome variable** is the continuous hate speech score, which indicates the degree of hate speech in a comment, with higher values denoting more hate speech and lower values suggesting counter or supportive speech. A score between -1 and +0.5 represents neutral or ambiguous content.

Key characteristics of the dataset include:

10 ordinal Labels	0	1	2	3	4
Sentiment (0-4)	No sentiments hurt		Sentiments might be hurt	Sentiments hurt (Religious/racial/cultural)	
Respect (0-4)	Respectful and supportive		Neutral	Disrespectful (Individual/group/both)	
Insult (0-4)	Supportive	Personal Level	Cultural/Racial Level	Personal / Group Level	
Dehumanize (0-4)	No dehumanization		Slightly Dehumanizing	Dehumanizing	
Humiliate (0-4)	Not humiliating		Ambiguous	Humiliating	

<b>Status (0-4)</b>	<b>Topic of status not initialized/Supportive</b>			<b>Slightly negative judgements</b>	<b>Strong negative judgements to demographics</b>
<b>Violence (0-4)</b>	<b>No violence</b>	<b>Explicit but not necessarily violent</b>		<b>Slightly more violent</b>	<b>Very violent</b>
<b>Genocide (0-4)</b>	<b>No Genocide</b>	<b>Ambiguous</b>		<b>Genocide provoked</b>	
<b>Attack_defend (0-4)</b>	<b>Defending/Supportive</b>		<b>Combination of defensive/attacking</b>	<b>Attacking on individual level and explicit</b>	<b>Attacking on level of genocide and violent</b>
<b>Hate_speech (0-2)</b>	<b>Neutral</b>	<b>offensive</b>	<b>Hate Speech</b>	<b>---</b>	<b>---</b>

## ii. Embedding models:

- **BERT:** Bidirectional Encoder Representations from Transformers, capturing semantic relationships.
- **GPT-Neo:** A generative transformer model, utilized for encoding nuanced language patterns.
- **TF-IDF:** Stands for term frequency-inverse document frequency, a statistical technique that measures how relevant words are to a document:
  - **Term frequency (TF):** The number of times a term appears in a document
  - **Document frequency (DF):** The number of documents that contain a word
  - **Inverse document frequency (IDF):** Weighs down frequent terms and increases the weight of rare terms
- **[11] Concatenated model of BERT and GPT-Neo :** may allow the model to better understand the relationships between the input and multiple labels it needs to predict.



```
# Concatenate GPT-Neo and BERT embeddings
concatenated_embeddings = torch.cat((projected_gpt_neo, projected_bert), dim=1)
print(concatenated_embeddings.shape) # Should be [54932, 768 + 768] = [54932, 1536]
```

```
torch.Size([54932, 1536])
```

emotionaldistress	provokingviolence	individualharrassment	embedding	emotionaldistress_binary	provokingviolence_binary	individualharrassment_binary
2	2	3	[0.004455165937542915, 0.33095115423202515, -0.0...	1	1	1
2	2	3	[0.8477875590324402, 1.9735788106918335, -1.19...	1	1	1
2	1	2	[0.011255793273448944, 1.4634126424789429, -0.0...	1	1	1
2	0	2	[-0.11528337001800537, 1.7641595602035522, -0.0...	1	0	1
2	2	3	[0.2188192456960678, 1.6596447229385376, -0.95...	1	1	1
...	...	...	...	...	...	...
2	1	3	[0.1726282835006714, 0.9378043677330017, -1.15...	1	1	1

### Labels converted to binarized labels

#### iii. Machine Learning Models for Multiclass Classification:

- **[10]BiLSTM:** BiLSTM extends LSTM by processing sequences in both forward and backward directions, enhancing context comprehension. This bidirectional approach improves classification accuracy, especially for context-sensitive tasks.
- **[9]XGBoost:** Uses gradient-boosted decision trees to efficiently handle structured data. Its scalability and regularization capabilities make it well-suited for predicting hate speech impact categories with high accuracy and robust performance.

#### iv. Machine learning models for Multilabel Classification:

- **[7][8] Hierarchical CNN-BiLSTM-Attention:** This model combines Convolutional Neural Networks (CNNs) for local feature extraction and BiLSTMs for capturing contextual dependencies in both forward and backward directions. The hierarchical structure processes data at multiple levels, and the attention mechanism helps the model focus on the most relevant parts of the input for each label. It excels in our multilabel classification by assigning multiple categories, improving its ability to detect various types of impact.
- **[12][13] Multi-Layer Perceptron (MLP):** This model uses a **Multi-Layer Perceptron (MLP)** for **multi-label classification** on embedded text data. The input text is first converted into numerical embeddings, and the model predicts multiple labels independently for each sample using a **sigmoid activation** for each output. The architecture consists of two fully connected layers, with a **dropout layer** in between to prevent overfitting. The model is trained

using **binary cross-entropy loss** and optimized with **AdamW**. It handles multi-label classification tasks, where each label is predicted separately, suitable for problems like emotional distress or harassment detection.

### 3.2 Development Process:

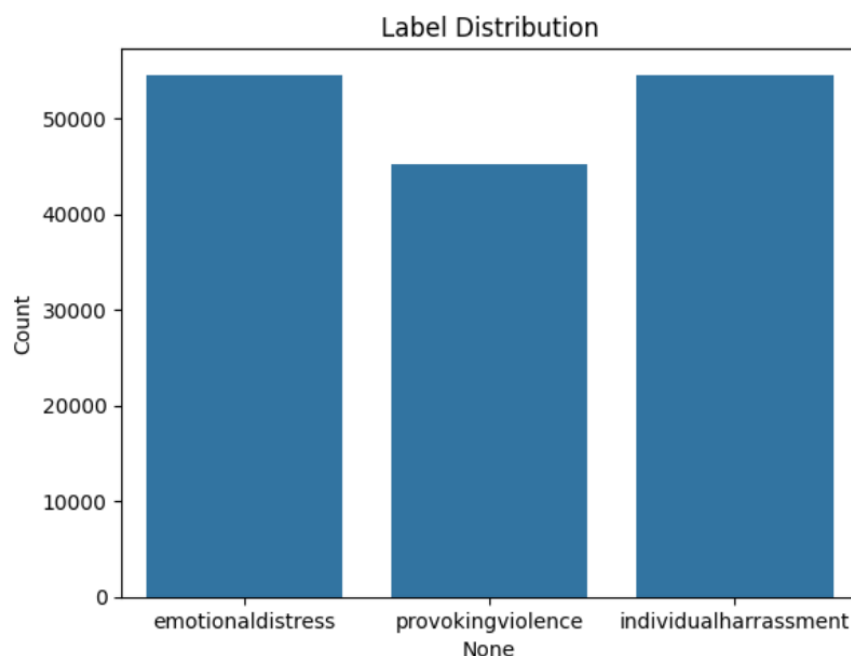
#### 3.2.1 Functional Requirements:

##### 1. Data Preprocessing:

- Dataset loading and Cleaning (handling null values) and mapping of new categories (feature engineering).
- **Tokenization:** Split text data into individual words or tokens for better analysis and **Stopword Removal:** Remove common, non-essential words (e.g., "the," "is," "and") to reduce noise in the dataset.
- **Dataset Filtering:** Identify and filter out hate speech instances based on pre-defined categories (e.g., emotional distress, provoking violence, individual harassment).

##### 2. Mapping with new categories (Average sum of upper and lower limits used):

1. Emotional Distress (0-2)
2. Provoking violence (0-3)
3. Individual Harassment (0-3)



**Label Distribution of for each label**

##### 3. Model Development:

- **Multiclass Classification Models:** Build models stated above that categorize hate speech into one of three classes—emotional distress, provoking violence, or individual harassment.
- **Multilabel Classification Models:** Implemented models stated above capable of predicting multiple categories for each instance, allowing for more comprehensive analysis.

#### 4. Performance Evaluation:

Use **precision, recall, weighted F1-score (for multiclass classification), and accuracy (for multilabel classification)** as key metrics to evaluate model performance. Also manually verify predictions of models for real world accuracy and performance.

#### 3.2.2 Non-Functional Requirements

- **Performance:** The models should process input text efficiently, with minimal latency during testing and real-time evaluations. Embeddings should be optimized to ensure fast and accurate predictions.
- **Scalability:** The system should support large-scale datasets and be adaptable to new or updated embeddings as they emerge, ensuring it can grow with increased data and technology advancements.
- **Usability:** The results from the models should be easy to interpret, providing clear classifications for users, allowing them to quickly understand the impact and context of hate speech instances.

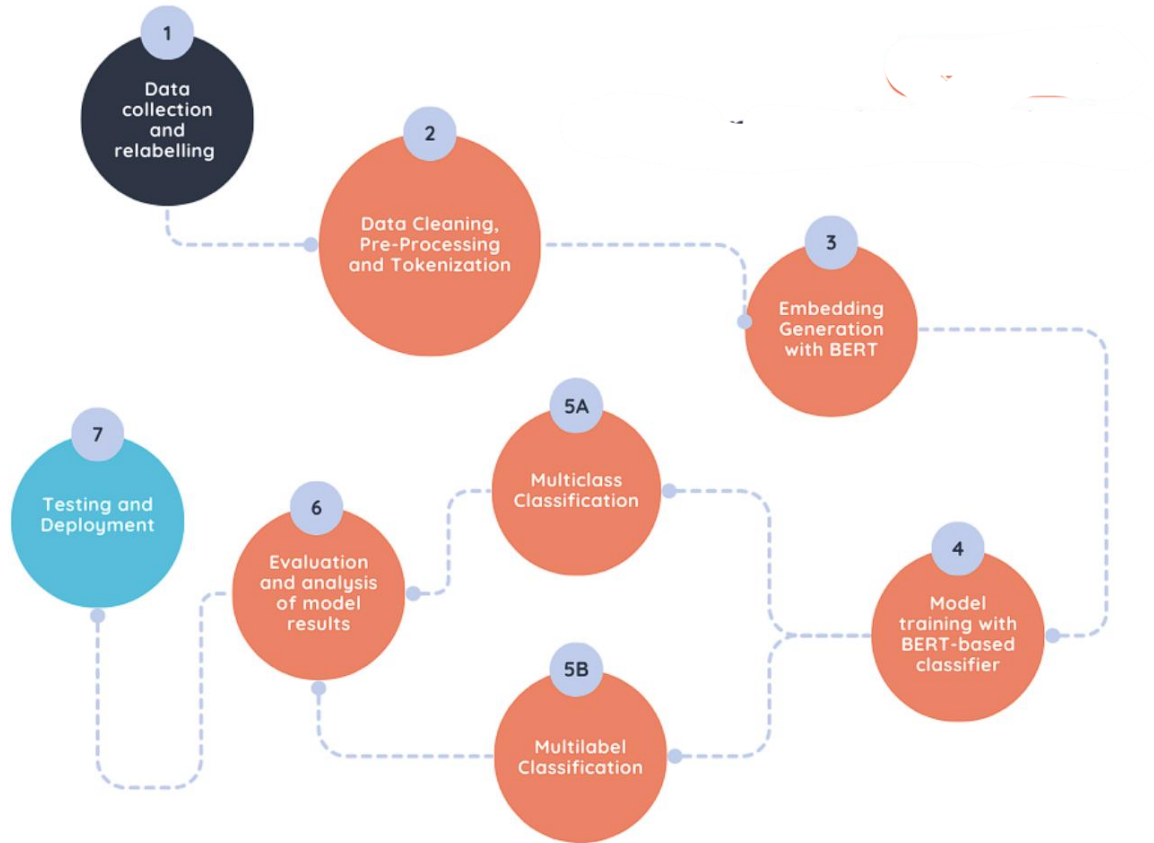
#### 3.3 Solution Approach:

- **Requirement Analysis:**
  - Dataset sourced from hate-speech dataset from hugging face, containing annotated instances categorized into normal, offensive, and hateful speech.
  - Labels refined to represent three impact categories: *emotional distress*, *provoking violence*, and *individual harassment*.
  - Preprocessed data stored in CSV or JSON format for ingestion by ML pipelines.
- **System Design:**

- Python-based frameworks employed for preprocessing, embedding generation, and model implementation.
- Machine learning models built and tested across embeddings (BERT, GPT-Neo) for performance comparisons.
- **Modules Built:**
  - Multiclass Classification Module: Implements models for exclusive classification of hate speech impacts.
  - Multilabel Classification Module: Designed for overlapping hate speech categories.
  - Evaluation Module: Calculates key metrics and visualizes model performance.

## Chapter-4 Modelling and Implementation Design

### 4.1 Design Diagram:



Embeddings were also done using GPT-Neo, TF-IDF and the concatenation of BERT and GPT-Neo,

### 4.2 Implementation Details:

#### 4.2.1 Multilabel Classification

##### 1. [12][13] Multilayer Perceptron (MLP)

###### a. Model Architecture

*Layers:*

- *Input Layer: Accept input features*
- *Hidden Layer: Fully connected, ReLU activation*
- *Dropout Layer: Prevent overfitting*
- *Output Layer: Fully connected, Sigmoid activation*

*Forward Function:*

- *Pass input through hidden layer, apply activation and dropout*
- *Pass result through output layer, apply sigmoid*

#### **b. BERT**

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	1.00	623
Provoking Violence	0.75	0.95	0.84	460
Individual Harassment	0.99	1.00	1.00	623
micro avg	0.91	0.99	0.95	1706
macro avg	0.91	0.98	0.94	1706
weighted avg	0.92	0.99	0.95	1706
samples avg	0.91	0.98	0.94	1706

Overall Accuracy: 0.7281

#### **c. GPT-neo**

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	1.00	10887
Provoking Violence	0.83	0.99	0.90	9012
Individual Harassment	0.99	1.00	1.00	10906
micro avg	0.94	1.00	0.97	30805
macro avg	0.94	1.00	0.96	30805
weighted avg	0.94	1.00	0.97	30805
samples avg	0.94	0.99	0.96	30805

Overall Accuracy: 0.8210

#### **d. TF-IDF**

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	1.00	10887
Provoking Violence	0.82	1.00	0.90	9012
Individual Harassment	0.99	1.00	1.00	10906
micro avg	0.93	1.00	0.97	30805
macro avg	0.93	1.00	0.96	30805
weighted avg	0.94	1.00	0.97	30805
samples avg	0.93	0.99	0.96	30805

Overall Accuracy: 0.8200

## **2. [7] Hierarchical CNN-BiLSTM Attention Model**

#### **a. Model Architecture**

- *CNN Layers:*
  - *Conv1: 1D Conv (input\_channels=1, output\_channels=128, kernel\_size=3)*

- *Conv2: 1D Conv (input\_channels=128, output\_channels=256, kernel\_size=3)*
  - *MaxPool: MaxPooling layer (kernel\_size=2)*
- *BiLSTM Layer:*
  - *BiLSTM (input\_size=256, hidden\_size=128, bidirectional=True)*
- *Fully Connected Layers:*
  - *FC\_shared: Linear (input\_size=128 \* 2 \* (input\_dim // 2), output\_size=128)*
- *Attention & Classification for each label:*
  - *Attention (Emotional Distress): Linear (input\_size=128)*
  - *FC (Emotional Distress): Linear (input\_size=128)*
  - *Attention (Provoking Violence): Linear (input\_size=128)*
  - *FC (Provoking Violence): Linear (input\_size=128)*
  - *Attention (Individual Harassment): Linear (input\_size=128)*
  - *FC (Individual Harassment): Linear (input\_size=128)*
- *Forward Pass:*
  - *CNN: Apply Conv1, Conv2, and MaxPool layers.*
  - *BiLSTM: Pass through Bidirectional LSTM.*
  - *Flatten: Reshape and pass through FC\_shared.*
  - *Attention & Classification:*
    - *For each label (Emotional Distress, Provoking Violence, Harassment):*
      - *Compute attention weights.*
      - *Multiply attention with shared features.*
      - *Pass through corresponding FC layer for classification.*

## b. BERT

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	0.99	10887
Provoking Violence	0.85	0.93	0.89	9012
Individual Harassment	0.99	1.00	1.00	10906
micro avg	0.95	0.98	0.96	30805
macro avg	0.94	0.97	0.96	30805
weighted avg	0.95	0.98	0.96	30805
samples avg	0.95	0.97	0.95	30805
Overall Accuracy: 0.8001				

## c. GPT-neo

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	1.00	10887
Provoking Violence	0.85	0.92	0.88	9012
Individual Harassment	0.99	1.00	1.00	10906
micro avg	0.95	0.98	0.96	30805
macro avg	0.94	0.97	0.96	30805
weighted avg	0.95	0.98	0.96	30805
samples avg	0.95	0.97	0.95	30805
Overall Accuracy: 0.7974				

## d. TF-IDF

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	1.00	5444
Provoking Violence	0.82	1.00	0.90	4505
Individual Harassment	0.99	1.00	1.00	5453
micro avg	0.93	1.00	0.97	15402
macro avg	0.93	1.00	0.96	15402
weighted avg	0.94	1.00	0.97	15402
samples avg	0.93	0.99	0.96	15402
Overall Accuracy: 0.8196				

## 4.2.2 Multiclass Classification

### 1. [10] BiLSTM

#### a. Model Architecture

- **Input Layer:**
  - Input size: (embedding\_size,)
- **Reshape Layer:**
  - Reshape input to (batch\_size, 1, embedding\_size) for LSTM compatibility.



- **BiLSTM Layer:**
  - Apply a bidirectional LSTM layer with desired number of units.
  - Add dropout for regularization.
- **Output Layers:**
  - Use a separate dense output layer for each label.
  - Softmax activation for multi-class classification.

## b. BERT

Classification Report for Provoking Violence:

	precision	recall	f1-score	support
0	0.53	0.06	0.11	3003
1	0.00	0.00	0.00	1429
2	0.60	0.93	0.73	8788
3	0.79	0.59	0.68	3260
accuracy			0.62	16480
macro avg	0.48	0.40	0.38	16480
weighted avg	0.57	0.62	0.54	16480

Classification Report for Individual Harassment:

	precision	recall	f1-score	support
0	0.00	0.00	0.00	125
1	0.57	0.09	0.15	3638
2	0.52	0.84	0.64	8117
3	0.56	0.34	0.42	4600
accuracy			0.53	16480
macro avg	0.41	0.32	0.30	16480
weighted avg	0.54	0.53	0.47	16480

Classification Report for Emotional Distress:

	precision	recall	f1-score	support
0	0.00	0.00	0.00	144
1	0.58	0.11	0.18	4736
2	0.72	0.97	0.83	11600
accuracy			0.71	16480
macro avg	0.43	0.36	0.34	16480
weighted avg	0.67	0.71	0.64	16480

### c. GPT-neo

Classification Report for Provoking Violence:				
	precision	recall	f1-score	support
0	0.54	0.02	0.04	3003
1	0.00	0.00	0.00	1429
2	0.58	0.95	0.72	8788
3	0.81	0.51	0.62	3260
accuracy			0.61	16480
macro avg	0.48	0.37	0.35	16480
weighted avg	0.57	0.61	0.52	16480

Classification Report for Individual Harassment:				
	precision	recall	f1-score	support
0	0.00	0.00	0.00	125
1	0.53	0.05	0.09	3638
2	0.51	0.93	0.66	8117
3	0.63	0.16	0.26	4600
accuracy			0.52	16480
macro avg	0.42	0.29	0.25	16480
weighted avg	0.54	0.52	0.42	16480

Classification Report for Emotional Distress:				
	precision	recall	f1-score	support
0	0.00	0.00	0.00	144
1	0.48	0.23	0.31	4736
2	0.74	0.90	0.81	11600
accuracy			0.70	16480
macro avg	0.40	0.38	0.37	16480
weighted avg	0.65	0.70	0.66	16480

### d. TF-IDF

Classification Report for Provoking Violence:				
	precision	recall	f1-score	support
0	0.55	0.08	0.14	3003
1	0.00	0.00	0.00	1429
2	0.60	0.90	0.72	8788
3	0.75	0.67	0.71	3260
accuracy			0.63	16480
macro avg	0.47	0.41	0.39	16480
weighted avg	0.57	0.63	0.55	16480

```

Classification Report for Individual Harassment:
      precision    recall  f1-score   support

     0         0.00      0.00      0.00        125
     1         0.56      0.14      0.22       3638
     2         0.52      0.81      0.63       8117
     3         0.56      0.36      0.44       4600

 accuracy          0.53       16480
 macro avg         0.41      0.33      0.32       16480
 weighted avg      0.54      0.53      0.48       16480

```

```

Classification Report for Emotional Distress:
      precision    recall  f1-score   support

     0         0.00      0.00      0.00        144
     1         0.65      0.25      0.36       4736
     2         0.75      0.95      0.84      11600

 accuracy          0.74       16480
 macro avg         0.47      0.40      0.40       16480
 weighted avg      0.72      0.74      0.69       16480

```

## 2. [9] XGBoost

### a. Model Architecture

- *XGBoost Classifier: For each target column, initialize an XGBoost model:*
  - *Set objective='multi:softmax' for multi-class classification.*
  - *Set num\_class based on the number of classes in the target.*
  - *Enable GPU support (tree\_method='gpu\_hist', gpu\_id=0).*
  - *Set max\_depth, learning\_rate, and n\_estimators for model parameters*

### b. BERT

```

Classification Report for 'provokingviolence':
      precision    recall  f1-score   support

     0         0.48      0.24      0.32       1975
     1         0.21      0.01      0.03        966
     2         0.63      0.85      0.72       5855
     3         0.75      0.70      0.73       2191

 accuracy          0.63      10987
 macro avg         0.52      0.45      0.45      10987
 weighted avg      0.59      0.63      0.59      10987

```

Overall Accuracy for 'provokingviolence': 0.6348

```

Classification Report for 'individualharrassment':
      precision    recall  f1-score   support

     0         0.29      0.02      0.05         81
     1         0.50      0.31      0.38        2386
     2         0.53      0.75      0.62        5430
     3         0.55      0.33      0.41        3090

 accuracy          0.53        10987
 macro avg         0.47      0.35      0.37        10987
 weighted avg      0.53      0.53      0.51        10987

```

Overall Accuracy for 'individualharrassment': 0.5316

```

Classification Report for 'emotionaldistress':
      precision    recall  f1-score   support

     0         0.50      0.03      0.06         100
     1         0.60      0.41      0.49        3151
     2         0.79      0.90      0.84        7736

 accuracy          0.75        10987
 macro avg         0.63      0.45      0.46        10987
 weighted avg      0.73      0.75      0.73        10987

```

Overall Accuracy for 'emotionaldistress': 0.7494

### c. GPT-neo

```

Classification Report for 'provokingviolence':
      precision    recall  f1-score   support

     0         0.50      0.22      0.31        1975
     1         0.15      0.01      0.01         966
     2         0.62      0.86      0.72        5855
     3         0.75      0.69      0.72        2191

 accuracy          0.63        10987
 macro avg         0.51      0.44      0.44        10987
 weighted avg      0.58      0.63      0.58        10987

```

Overall Accuracy for 'provokingviolence': 0.6347

```

Classification Report for 'individualharrassment':
      precision    recall  f1-score   support

     0         0.20      0.01      0.02         81
     1         0.51      0.29      0.37        2386
     2         0.53      0.77      0.63        5430
     3         0.56      0.32      0.40        3090

 accuracy          0.53        10987
 macro avg         0.45      0.35      0.36        10987
 weighted avg      0.53      0.53      0.51        10987

```

Overall Accuracy for 'individualharrassment': 0.5342

```

Classification Report for 'emotionaldistress':
      precision    recall  f1-score   support

     0       0.00      0.00      0.00        100
     1       0.59      0.40      0.48       3151
     2       0.78      0.90      0.84       7736

 accuracy          0.75      10987
 macro avg         0.46      0.43      0.44      10987
 weighted avg      0.72      0.75      0.73      10987

Overall Accuracy for 'emotionaldistress': 0.7477

```

#### d. TF-IDF

```

Classification Report for 'provokingviolence':
      precision    recall  f1-score   support

     0       0.58      0.15      0.24       1975
     1       0.00      0.00      0.00        966
     2       0.62      0.91      0.73       5855
     3       0.78      0.66      0.71       2191

 accuracy          0.64      10987
 macro avg         0.49      0.43      0.42      10987
 weighted avg      0.59      0.64      0.58      10987

```

Overall Accuracy for 'provokingviolence': 0.6418

```

Classification Report for 'individualharrassment':
      precision    recall  f1-score   support

     0       0.00      0.00      0.00         81
     1       0.56      0.18      0.27       2386
     2       0.52      0.83      0.64       5430
     3       0.57      0.30      0.39       3090

 accuracy          0.53      10987
 macro avg         0.41      0.33      0.33      10987
 weighted avg      0.54      0.53      0.49      10987

```

Overall Accuracy for 'individualharrassment': 0.5335

```

Classification Report for 'emotionaldistress':
      precision    recall  f1-score   support

     0       1.00      0.01      0.02        100
     1       0.61      0.30      0.41       3151
     2       0.76      0.93      0.84       7736

 accuracy          0.74      10987
 macro avg         0.79      0.41      0.42      10987
 weighted avg      0.72      0.74      0.71      10987

```

Overall Accuracy for 'emotionaldistress': 0.7417

### 4.2.3 Combined Embeddings – BERT and GPT-Neo

## 1. Multilabel Classification using Hierarchial CNN-BiLSTM-Attention classifier

	precision	recall	f1-score	support
Emotional Distress	0.99	1.00	1.00	10887
Provoking Violence	0.85	0.92	0.89	9012
Individual Harassment	0.99	1.00	1.00	10906
micro avg	0.95	0.98	0.96	30805
macro avg	0.95	0.97	0.96	30805
weighted avg	0.95	0.98	0.96	30805
samples avg	0.95	0.97	0.95	30805

Overall Accuracy: 0.8023

### 4.3 Risk Mitigation and Analysis

The "Hate Speech Impact Analysis" project addressed various risks to ensure the system's reliability, ethical integrity, and scalability. Key risks included model performance challenges, data privacy concerns, ethical issues, and scalability limitations.

#### 4.3.1 Model Performance Risks

Detecting nuanced impacts of hate speech required high model precision and recall. Misclassification could have compromised the system's effectiveness.

- **Strategies:** Models were trained using embeddings like BERT and GPT\_neo, evaluated through metrics like F1-score, and improved with hyperparameter tuning and ensemble methods (CNNs and Bi-LSTMs).
- **Outcome:** These measures ensured robust performance and accurate classifications.

#### 4.3.2 Data Privacy Risks

The use of real-world datasets posed risks of exposing sensitive information.

- **Strategies:** Data was anonymized by removing PII, and private systems with restricted access were used for secure storage.
- **Outcome:** Privacy was maintained, and ethical standards were upheld.

#### 4.3.3 Technical Scalability Risks

Growing data volumes required scalable solutions to avoid performance issues.

- **Strategies:** Scalable models like BERT and cloud infrastructure like Google Colab were used for efficient processing.
- **Outcome:** The system adapted to increased demand without delays or degradation.

## Chapter 5 - Testing

This chapter documents the testing process undertaken during the "Hate Speech Impact Analysis" project. The testing focused on ensuring the accuracy, scalability, and ethical compliance of the system. Key components were evaluated to validate performance and prepare for findings.

### 5.1 Testing Plan

The testing process aimed to:

- Validate the effectiveness of the classification models.
- Ensure data privacy and bias mitigation during training and evaluation.
- Assess the system's scalability to handle diverse datasets and increasing data volumes.

A structured, iterative methodology was adopted, emphasizing performance metrics, component validation, and ethical considerations.

### 5.2 Component Testing

The project's components were tested independently and as part of the complete workflow:

- **Classification Models:** Evaluated the CNN-BiLSTM-Attention hybrid architecture, MLP model, XgBoost, BiLSTM for reliability and accuracy.
- **Attention Mechanisms:** Tested the contributions of label attention and document self-attention layers to classification outcomes.

### 5.3 Dataset Validation

Efforts were made to ensure datasets met research and ethical standards:

- Data anonymization protected user privacy.
- Class balancing techniques addressed imbalanced labels, improving model performance.
- Dataset splitting avoided training-validation overlap to maintain unbiased evaluations.

### 5.4 Model Evaluation

Key metrics—accuracy, precision, recall, and F1-score—were used to evaluate the system:

- Cross-validation techniques ensured consistent results across data subsets.
- The CNN-LSTM hybrid model showed superior performance compared to other tested architectures.

## 5.5 Test Results

Example 1:  
True Labels (y): [1 0 1]  
Predicted Labels: [1 0 1]

Example 2:  
True Labels (y): [1 0 1]  
Predicted Labels: [1 1 1]

Example 3:  
True Labels (y): [1 0 1]  
Predicted Labels: [1 1 1]

Example 4:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 5:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 6:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 7:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 8:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 9:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 10:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 11:  
True Labels (y): [1 0 0]  
Predicted Labels: [1 1 1]

Example 12:  
True Labels (y): [1 1 1]  
Predicted Labels: [1 1 1]

Example 13:  
True Labels (y): [0 0 0]  
Predicted Labels: [1 1 1]

**Examples of the results as predicted by our models.**

0 represents absence and 1 represents presence of the labels in multilabel classification



## Chapter 6 - Findings, Conclusions, and Future Scope

This chapter summarizes the findings derived from the project and discusses conclusions drawn from the analysis. It also outlines potential future directions for advancing this work.

### 6.1 Findings

#### 1. Model Performance

- The CNN-BiLSTM-Attention hybrid architecture emerged as the most robust model, achieving reliable classification across labels such as emotional distress, provoking violence, and individual harassment.
- Performance evaluation using metrics like accuracy, precision, recall, and F1-score highlighted the system's ability to balance precision and generalizability.

#### 2. Testing and Validation

- Cross-validation and rigorous evaluation processes ensured unbiased performance across datasets.
- Dataset validation highlighted the system's scalability, confirming its capability to handle increasing data volumes effectively.

### 6.2 Conclusions

The "Hate Speech Impact Analysis" project successfully addressed its core research objectives by developing a robust, ethical, and scalable system. The classification models demonstrated the ability to accurately detect and analyze hate speech impacts, including emotional distress, violence provocation, and harassment.

Key achievements include:

- Creating a pipeline that combines advanced embeddings, attention mechanisms, and hybrid architectures for accurate classification.
- Ensuring ethical compliance and privacy protection, integral to research in sensitive areas.
- Validating scalability and adaptability to accommodate diverse datasets and potential real-world applications.

These limitations, however, present avenues for future exploration.

### **6.3 Future Scope**

The findings and limitations of this project open up several pathways for future research and development:

#### **1. Enhanced Dataset Diversity**

- Expanding the dataset to include more culturally and linguistically diverse examples can improve the system's generalizability.

#### **2. Advanced Model Architectures**

- Incorporating cutting-edge models, such as GPT-based architectures or multi-modal systems, can further enhance classification accuracy and adaptability.

#### **3. Real-World Deployment**

- Developing a user-facing platform or extension for real-time analysis would allow for practical application and iterative refinement based on user feedback.

#### **4. Multi-Modal Integration**

- Including non-textual data, such as images or videos, would enable a comprehensive analysis of hate speech across different media types.

#### **5. Automated Bias Mitigation**

- Building automated systems to detect and mitigate biases dynamically during training can further ensure fairness and inclusivity.

### **Final Remarks**

The "Hate Speech Impact Analysis" project represents a significant step toward understanding and mitigating the impacts of hate speech. By combining advanced machine learning techniques, ethical considerations, and scalable systems, this research establishes a solid foundation for future innovations in hate speech detection and analysis. With further development and deployment, the system can provide real-world solutions to address the growing challenges of online hate speech.

## References

- [1]: Zhou, Y., Yang, Y., Liu, H., Liu, X., & Savage, N. (2020). Deep learning based fusion approach for hate speech detection. *IEEE Access*, 8, 128923-128929.
- [2]: Aluru, S. S., Mathew, B., Saha, P., & Mukherjee, A. (2020). Deep learning models for multilingual hate speech detection. *arXiv preprint arXiv:2004.06465*.
- [3]: Das, M., Saha, P., Mathew, B., & Mukherjee, A. (2022). Hatecheckhin: Evaluating hindi hate speech detection models. *arXiv preprint arXiv:2205.00328*.
- [4]: Vijayaraghavan, P., Larochelle, H., & Roy, D. (2021). Interpretable multi-modal hate speech detection. *arXiv preprint arXiv:2103.01616*.
- [5]: Badjatiya, P., Gupta, S., Gupta, M., & Varma, V. (2017, April). Deep learning for hate speech detection in tweets. In *Proceedings of the 26th international conference on World Wide Web companion* (pp. 759-760).
- [6]: Zimmerman, S., Kruschwitz, U., & Fox, C. (2018, May). Improving hate speech detection with deep learning ensembles. In *Proceedings of the eleventh international conference on language resources and evaluation (LREC 2018)*.
- [7]: F. Sun and N. Chu, "Text Sentiment Analysis Based on CNN-BiLSTM-Attention Model," *2020 International Conference on Robots & Intelligent System (ICRIS)*, Sanya, China, 2020, pp. 749-752, doi: 10.1109/ICRIS52159.2020.00186.  
keywords: {Sentiment analysis;Analytical models;Recurrent neural networks;Text categorization;Semantics;Interference;Feature extraction;Word Vector;Convolution Neural Network (CNN);Text Sentiment Analysis;Bidirectional Long Short-Term Memory (BiLSTM)},
- [8]:  
[https://www.academia.edu/118613168/Sarcasm\\_Detection\\_with\\_A\\_New\\_CNN\\_BiLSTM\\_Hybrid\\_Neural\\_Network\\_and\\_BERT\\_Classification\\_Model](https://www.academia.edu/118613168/Sarcasm_Detection_with_A_New_CNN_BiLSTM_Hybrid_Neural_Network_and_BERT_Classification_Model) -International Journal of Advanced Networking and Applications
- [9]: Sharma, Surbhi & Joshi, Nisheeth(2024) A fusion approach to detect sarcasm using NLTK models BERT and XG Boost, *Journal of Information and Optimization Sciences*, 45:4, 981–990, DOI: [10.47974/JIOS-1621](https://doi.org/10.47974/JIOS-1621)
- [10]: Wibawa, Aji & Cahyani, Denis & Prasetya, Didik & Gumilar, Langlang & Nafalski, Andrew. (2023). Detecting emotions using a combination of bidirectional encoder representations from

transformers embedding and bidirectional long short-term memory. International Journal of Electrical and Computer Engineering (IJECE). 13. 7137. 10.11591/ijece.v13i6.pp7137-7146.

[11] <https://medium.com/mantisnlp/how-to-combine-several-embeddings-models-8e7bc9a00330>

[12] <https://medium.com/deep-learning-study-notes/multi-layer-perceptron-mlp-in-pytorch-21ea46d50e62>

[13] <https://austingwalters.com/classify-sentences-via-a-multilayer-perceptron-mlp/>