# DETECTING HATE SPEECH AND INSULTS ON SOCIAL MEDIA

A thesis submitted in partial fulfilment of the requirement for the

**Degree of Bachelor of Science**

Of

**Ramakrishna Mission Residential College**

By

**SOUMYADEEP NASKAR**

Reg. No:  A03-1122-0223-22

**AKASH PAL**

Reg. No: A03-1142-0212-22

**MURLIDHAR PATRA**

Reg. No:  A03-1132-0227-22

Under The Guidance Of

**Dr. SURAJIT GIRI**

**Professor**

Department Of Computer Science

RKMRC Narendrapur , Kolkata - 700103

# Acknowledgement

It is our great pleasure to express our profound sense of gratitude to our esteemed Supervisor **Dr. Surajit Giri**, for providing his constructive academic advice and guidance, constant encouragement, and valuable suggestions at crucial junctures and all other support throughout this project work, and for helping us to prepare the project report successfully. We really benefited from his excellent supervision. We would extend our sincere thanks to our respected Head of the Department. **Dr. Siddhartha Banerjee**, for allowing us to use the facilities available. I would like to thank teachers of our department for extending this wonderful opportunity of working on a project as our DSE-4 in our curriculum.

## Certificate

I hereby certify that the project report titled "DETECTING HATE SPEECH AND INSULTS ON SOCIAL MEDIA" which is submitted by **SOUMYADEEP NASKAR** ( Reg. No:  A03-1122-0223-22) **AKASH PAL** (Reg. No: A03-1142-0212-22) **MURLIDHAR PATRA** (Reg. No:  A03-1132-0227-22) to the faculty of the Computer Science Department of Ramakrishna Mission Residential College in Complete fulfilment of the requirements for the Degree of Bachelor of Science (Honours) in Computer Science, is a record of the project work carried out by the students under my supervision in the academic session of the final semester (semester VI) of 2022- 2025.

**(Signature of Head of the Department)**          **(Signature of Supervisor)**

# Abstract

The exponential growth of social media has transformed communication, offering unprecedented opportunities for expression and connection. However, this rapid expansion has also enabled the widespread dissemination of harmful and abusive content, particularly hate speech and insults. Hate speech—characterized by derogatory or threatening language targeting individuals or groups based on attributes such as race, religion, gender, or ethnicity—poses a serious threat to the safety, dignity, and inclusivity of online communities. Detecting and mitigating such toxic behaviour is a critical task for maintaining healthy digital environments.

This thesis investigates the problem of automatically detecting hate speech and insults on social media platforms using natural language processing (NLP) and machine learning techniques. We analyze existing methods ranging from classical algorithms such as logistic regression and support vector machines to state-of-the-art deep learning models like Convolution Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformer-based architectures (e.g., BERT). Emphasis is placed on understanding the linguistic and contextual intricacies that make hate speech detection especially challenging such as the use of sarcasm, slang, implicit bias, and rapidly evolving internet language.

Multiple benchmark datasets are utilized to train and evaluate the models, including datasets derived from platforms like Twitter. Data pre processing steps such as tokenization, stop-word removal, and vectorization (using TF-IDF, Word2Vec, or contextual embeddings) are explored to enhance model performance. Metrics like accuracy, precision, recall, F1-score, and area under the ROC curve (AUC-ROC) are used for performance evaluation.

Beyond technical implementation, the thesis critically examines ethical considerations in deploying hate speech detection systems, including risks of over-censorship, algorithmic bias, and the potential for reinforcing social inequalities. It also highlights the trade-offs between maintaining user safety and preserving freedom of expression.

The findings demonstrate that while no model is perfect, combining linguistic features with deep contextual understanding significantly improves detection accuracy. This research contributes to the ongoing development of more reliable, context-aware, and ethically responsible systems for content moderation on digital platforms.

# Contents

# **Introduction**

## 1.1 Background and Motivation

Social media platforms have become integral to modern communication, information sharing, and community building. They provide users with the freedom to express thoughts and opinions to a global audience, often in real time. However, this democratization of speech has also led to the proliferation of harmful content, including hate speech and online abuse. The anonymity, virality, and scale of social media make it a powerful but dangerous tool when misused.

Hate speech is defined as communication that expresses hatred, promotes violence, or incites discrimination against individuals or groups based on characteristics such as race, ethnicity, religion, gender, sexual orientation, or disability. Online insults—though often less structured than hate speech—can still contribute to hostile digital environments and emotional harm. Such content has been linked to increased polarization, marginalization of vulnerable communities, and offline violence.

Manual moderation is neither scalable nor fully effective due to the sheer volume of content generated every second on platforms like Twitter, Face book, Reddit, and Instagram. Human moderators are also prone to subjectivity and psychological stress. Therefore, automating the detection of hate speech and insults using computational techniques has become a priority in both academia and industry.

Despite progress in natural language processing (NLP) and machine learning (ML), the automatic identification of hate speech remains a highly challenging task. Hate speech can be explicit or implicit, expressed through sarcasm, slang, coded language, memes, or multimodal formats. Furthermore, what qualifies as hate speech often varies across cultural, legal, and linguistic contexts, raising questions about fairness and freedom of speech. Addressing these complexities demands nuanced and adaptable computational models that go beyond keyword-based detection.

## 1.2 Problem Statement

The central problem this thesis addresses is the development of accurate, fair, and context-aware systems for detecting hate speech and insulting language on social media. Existing models often suffer from limitations such as:

- High false positive rates due to over-reliance on offensive keywords.
- Inability to understand context, sarcasm, or cultural references.

- Bias in training data, which can result in unfair targeting of specific dialects or communities.
- Challenges in handling multilingual and code-switched content (e.g., mixing English with Hindi or Bengali).
- Ethical concerns related to censorship, freedom of expression, and algorithmic transparency.

This thesis investigates solutions that combine linguistic insights with advanced machine learning models to improve the detection of hate speech while addressing these challenges.

## 1.3 Effects of Hate Speech

### 1.3.1 Psychological Impact and Xenophobia

Studies on hate speech toward refugees reveal that negative stereotypes, media bias, and misinformation contribute significantly to xenophobia. Refugees are frequently portrayed as threats or burdens, resulting in psychological effects such as fear, anger, and social isolation. While some individuals respond through counter-narratives and positive storytelling, others retreat into silence to avoid further discrimination. Researchers emphasize the need for empathetic representation and stronger interventions to combat such hate.

### 1.3.2 Impact on the LGBT Community

Hate speech directed at the LGBT community leads to serious psychological and social ramifications. Victims often experience anxiety, depression, low self-esteem, and even suicidal ideation. The normalization of online abuse discourages victims from reporting incidents, and current moderation efforts by social media platforms remain insufficient. The literature advocates for enhanced legal frameworks, robust reporting mechanisms, and comprehensive awareness campaigns.

### 1.3.3 Religious Intolerance

Research conducted among Indonesian teenagers indicates that exposure to hate speech and fake news increases religious intolerance, particularly among individuals with low media literacy. The findings suggest that improving media literacy can mitigate prejudice and promote tolerance, highlighting education as a key defense against the harmful effects of online hate.

### 1.3.4 Misinformation

During the COVID-19 pandemic, hate speech on Indian social media surged, driven by misinformation. Narratives such as "Corona Jihad" fueled communal tensions and violence, targeting minority groups. The study underscores the ethical and legal

challenges in regulating hate speech and stresses the importance of stronger enforcement by social media platforms.

### 1.3.5 Political Conflicts

In Kenya, ethnic hate speech on social media intensified political divisions during the 2013 elections. Unlike in 2007, when traditional media played a dominant role, social media became the main platform for spreading digitized hate, influencing political discourse and public sentiment. The study finds a direct link between online hate speech and real-world political hostility.

### 1.3.6 Cyber Terrorism

Hate speech often escalates following terrorist events, contributing to cyber terrorism and social polarization. Extremist groups use social media to propagate propaganda and recruit members, exploiting the lack of consistent legal regulation. Collaboration between governments, service providers, and tech companies is recommended to address the rising threat of online extremism.

### 1.3.7 Harmful Sexual Behaviour (HSB) Among Youth

The rise in harmful sexual behaviour among young people is closely tied to online interactions and social media exposure. Victims often suffer from trauma and anxiety, while perpetrators may develop distorted perceptions of intimacy. The normalization of sexting and peer grooming poses serious developmental risks, calling for targeted interventions and education strategies.

### 1.3.8 Partisan Media and Racial Health Disparities

Partisan news shared on social media significantly influences public perceptions of racial and ethnic health disparities. Liberal outlets tend to emphasize systemic causes, while conservative sources often frame the issue through individual responsibility. Both contribute to shaping attitudes, but a general lack of balanced reporting persists. The study highlights the need for evidence-based media to inform public discourse.

### 1.3.9 Youth Suicide and Self-Harm

An analysis of youth conversations on Instagram reveals that while some expressions of self-harm reflect genuine distress, many are casual or exaggerated. This blurs the line between serious mental health issues and social banter, making it challenging for intervention. The study recommends context-aware moderation systems that promote supportive peer networks without over-policing discourse.

### 1.3.10 Religious Violence and Online Reactions

Extremist incidents often trigger spikes in hate speech, particularly against Arabs and Muslims. Anti-Muslim rhetoric intensifies after Islamist terrorist events, while Islamophobic attacks receive less public backlash. The data also shows that different platforms respond differently, with Twitter exhibiting more aggressive discourse than Reddit. These dynamics highlight the need for responsive monitoring systems post-crisis

### 1.4 Research Objectives

The main objectives of this research are as follows:

1. **To analyze** existing methodologies and models used in hate speech detection across different social media platforms.
2. **To pre-process and accurate** high-quality datasets that reflect real-world use cases and linguistic diversity.
3. **To implement and compare** various machine learning models, including traditional algorithms (Logistic Regression, SVM) and deep learning techniques (LSTM, CNN, BERT).
4. **To evaluate** these models using standard performance metrics and identify the best-performing approaches under different conditions.
5. **To explore** the role of context-aware features, such as semantic embeddings and user metadata, in improving classification accuracy.
6. **To discuss** the ethical and social implications of deploying automated hate speech detection systems, particularly in relation to fairness, transparency, and free speech.

### 1.5 Scope and Limitations

This study focuses on text-based hate speech and insults in English, with potential extensions to multilingual or code-mixed content in future work. The primary platforms considered are Twitter and Reddit due to their open APIs and availability of annotated datasets. While some models may generalize to other platforms, each social media environment has its own norms, slang, and context which may affect model performance.

Limitations include:

- Potential biases in publicly available datasets.
- Limited access to private or deleted content.
- Absence of multimodal elements like images or videos, which may also convey hate.
- Challenges in distinguishing hate speech from satire or controversial opinions.

## 1.6 Significance of the Study

The results of this research have both academic and practical significance. From a technical standpoint, it advances the state of the art in natural language processing and toxic content classification. From a societal perspective, it contributes to safer online communities by equipping platforms with tools to mitigate abuse and protect users, particularly those from vulnerable or marginalized groups.

Moreover, this thesis promotes responsible AI practices by incorporating ethical considerations into the design and evaluation of machine learning models. It recognizes that hate speech detection is not just a technical problem, but also a social and moral one that must balance competing values like safety, privacy, fairness, and freedom of expression.

.

# Literature Review

## 2.1 Introduction

With the growing ubiquity of social media platforms, online hate speech has become a significant concern, prompting a surge in research on its detection using Natural Language Processing (NLP). Researchers have employed a variety of techniques ranging from traditional machine learning algorithms to advanced deep learning and transformer-based models. This chapter reviews the major contributions, datasets, and methodologies in the field, highlighting the evolution of hate speech detection techniques over the past decade.

## 2.2 Datasets for Hate Speech Detection

Several curate datasets have been instrumental in advancing the field:

Davidson et al. (2017) collected 24,783 tweets annotated into hate speech, offensive language, or neither. Their dataset is widely used and set a precedent for category-based annotation (Davidson, Warmsley, Macy, & Weber, 2017).Waseem and Hovy (2016) compiled 16,000 tweets focusing on racist and sexist content. They emphasized the importance of annotator bias and proposed metadata features for improved detection (Waseem & Hovy, 2016).

Founta et al. (2018) introduced a larger dataset with four categories: hate speech, abusive language, spam, and normal. Their work enhanced multi-class classification performance (Founta, Djouvas, Chatzakou, Leontiadis, Blackburn, Stringhini, Vakali, Sirivianos, & Kourtellis, 2018).Zampieri et al. (2019) created the OffensEval dataset, which was used in the SemEval competition to encourage benchmarking across multiple languages and offensive behavior types.These datasets have significantly influenced the choice of models and performance metrics across various studies.

## 2.3 Traditional Machine Learning Approaches

Early studies on hate speech detection focused on classical machine learning methods:

Support Vector Machines (SVM) and Logistic Regression (LR) were commonly applied due to their interpretability and performance on sparse features (Davidson et al., 2017).

Naïve Bayes (NB) was another popular choice for text classification tasks because of its simplicity and speed (Waseem & Hovy, 2016).

Feature engineering played a crucial role in these methods. Researchers used n-grams, TF-IDF vectors, part-of-speech tags, sentiment scores, and hate lexicons (such as Hatebase) to enhance model performance.
However, these methods struggled with: Understanding semantic and syntactic

## 2.4 Deep Learning Methods

With the rise of neural networks, deep learning models demonstrated superior performance:

Convolutional Neural Networks (CNNs) were employed to extract local n-gram features and semantic cues from short texts (Kim, 2014).

Recurrent Neural Networks (RNNs), particularly LSTMs and GRUs, were adopted to model sequential patterns in tweets (Badjatiya, Gupta, Gupta, & Varma, 2017).

Bidirectional LSTM (BiLSTM) models were also used for contextual understanding and showed improved accuracy on code-mixed and low-resource datasets (Bohra, Vijay, Singh, Akhtar, & Shrivastava, 2018).

These models benefited from word embeddings such as GloVe and Word2Vec, allowing for semantically rich input representations. However, they were computationally expensive and struggled with longer dependencies.

## 2.5 Transformer-based Models and Transfer Learning

The advent of transformer-based architectures, particularly BERT, revolutionized hate speech detection:

BERT (Devlin et al., 2019), with its bidirectional context understanding and fine-tuning capabilities, became the de facto standard in recent years.

RoBERTa, DistilBERT, and ALBERT offered lightweight or optimized alternatives with competitive results.

Multilingual variants like mBERT facilitated hate speech detection in Hindi, Arabic, and Bengali, addressing the low-resource language gap (Mozafari, Farahbakhsh, & Crespi, 2020).

Transfer learning allowed researchers to pretrain on massive corpora and fine-tune on smaller, annotated hate speech datasets. These models demonstrated robustness to sarcasm, code-mixing, and language variation.

## 2.6 Multilingual and Code-Mixed Challenges

Detecting hate speech in multilingual or code-mixed contexts (e.g., Hindi-English) remains challenging:

Code-mixing introduces issues like inconsistent grammar, spelling, and script usage. Studies like Bohra et al. (2018) focused on bilingual hate detection using hybrid CNN-BiLSTM architectures.

Resource-poor languages face data scarcity, leading to increased reliance on cross-lingual embeddings and transfer learning.

Annotating hate speech in low-resource languages is complicated by regional dialects, cultural context, and subjective interpretation.

## 2.7 Evaluation Metrics and Model Analysis

Most studies evaluated model performance using:

Accuracy, Precision, Recall, F1-score (macro/weighted), ROC-AUC for binary classification, Confusion matrices for error analysis.

To address class imbalance, some used: Oversampling/under sampling, SMOTE, Weighted loss functions.

Ablation studies, model interpretability (e.g., attention visualization), and fairness audits have recently been incorporated to ensure responsible AI.

## 2.8 Research Gaps and Challenges

Despite progress, several issues remain:

Bias in annotations can lead to skewed results and discrimination against certain groups.

Adversarial examples and evolving slang are difficult for static models to detect.

Lack of generalizability across platforms (Twitter vs. Reddit) and domains.

## 2.9 Summary

The literature on hate speech detection has evolved from rule-based systems and traditional ML classifiers to deep learning and transformer-based models. While recent models achieve state-of-the-art performance, challenges like bias mitigation, low-resource language support, and explainability require further research. Continued benchmarking, diverse dataset collection, and interdisciplinary collaboration are essential for building ethical and robust hate speech detection systems.

# Proposed Methodology

The proliferation of toxic comments on social media platforms has led to the need for automated systems that can detect and classify such comments efficiently. Toxicity in online communication can cause psychological harm, incite violence, and contribute to the spread of misinformation. Therefore, detecting toxic comments has become a crucial task in the field of Natural Language Processing (NLP).

## 3.1 Dataset Description

The dataset used in this study comprises a collection of 24,783 user-generated tweets, annotated and classified according to the presence of hate speech, offensive language, or neither. This dataset serves as the foundational resource for developing a toxic comment classification system, enabling the differentiation of harmful language from benign expressions in online discourse.

### 3.1.1 Structure and Attributes

The dataset contains six columns, each of which contributes to either the annotation process or the content of the tweet itself

### 3.1.2 Count

This column indicates the total number of human annotators who reviewed each tweet. In most cases, tweets were evaluated by three annotators, though a few instances involved more. This redundancy ensures greater reliability and consensus in the classification process.

### 3.1.3 Hate speech count

This column records how many annotators labeled the tweet as hate speech. A high value here (e.g., 2 or 3) signifies strong agreement among annotators that the tweet contains targeted hostility based on characteristics such as race, religion, gender, or ethnicity.

### 3.1.4 Offensive language count

This field captures how many annotators considered the tweet to be offensive, but not necessarily hate speech. Offensive language includes vulgar expressions, insults, and

profane words, often used in a derogatory or aggressive context but without targeting specific identity groups.

### 3.1.5 Neither count

This column reflects how many annotators believed that the tweet was neither offensive nor hateful. These tweets may contain strong opinions or informal language but are not deemed harmful or inappropriate.

### 3.1.6 Class

The class column is a derived categorical label representing the final classification of the tweet based on the annotators' votes. The values are encoded as:

0 = Hate Speech

1 = Offensive Language

2 = Neither

The class label represents the majority opinion among the annotators and serves as the ground truth for training and evaluating classification models.

### 3.1.6 Tweet

This column contains the actual text content of the tweet. These texts are raw, unfiltered, and often contain spelling errors, abbreviations, hashtags, mentions, and slang, reflecting the informal nature of social media language.

### 3.2 Nature of the Data

The tweets are written in English and reflect a wide range of social commentary, personal opinions, and interactions typical of Twitter. Due to the nature of the classification, many entries include sensitive content such as slurs, insults, or highly charged language. It is important to handle and process such data with ethical caution, both in terms of model fairness and in minimizing exposure to annotators or researchers.

### 3.3 Annotation Process

Each tweet was reviewed by multiple human annotators, who classified it into one of three categories. The use of multiple annotators ensures a robust consensus-based labeling process, which reduces subjectivity and enhances the dataset's overall quality. The detailed count columns (hate_speech_count, offensive_language_count, neither_count) enable further analysis of annotation disagreements and potential label noise in the dataset.

**3.4 Label Distribution (Summary)**

While exact figures were not detailed above, prior studies using similar datasets indicate that:

Offensive Language (class = 1) tends to dominate the dataset due to the frequency of casual profanity or internet arguments.

Hate Speech (class = 0) represents a smaller portion, reflecting the relative rarity—but serious impact—of such content.

Neither (class = 2) captures tweets that, while possibly intense or emotional, are not toxic in nature.

This imbalance among classes presents a challenge for machine learning models, which may become biased toward the majority class unless proper balancing techniques are applied (e.g., re-sampling, class weighting).

**3.5 Importance of the Dataset**

This dataset plays a critical role in the development of automated content moderation systems. By training models on this labelled data, it becomes possible to identify harmful speech patterns, support safer online environments, and inform public discourse analysis. Moreover, its transparent annotation structure and rich text content make it valuable for natural language processing (NLP) research in areas like sentiment analysis, bias detection, and linguistic toxicity.

**3.6  Data Pre-processing – Refining Raw Text for Analysis**

Raw tweets are often noisy, unstructured, and filled with irrelevant information. Therefore, comprehensive pre-processing is vital before feeding data into a machine learning model.

**3.6.1  Lowercasing**

All text was converted to lowercase to eliminate redundancy. This ensures that 'Hate', 'hate', and 'HATE' are treated identically, which simplifies the vocabulary space and reduces dimensionality.

### 3.6.2 Tokenization

Tokenization is the process of breaking down sentences into individual words or phrases (tokens). In this project, we used the NLTK tokenizer, which efficiently splits tweets into manageable tokens. Tokenization is the first step toward feature extraction and further text processing.

### 3.6.3 Removing Special Characters and Noise

Regular expressions were used to remove URLs, user mentions (@username), hashtags (#topic), numbers, punctuation, and emojis. These elements usually don't contribute meaningfully to the classification task and act as noise. This step ensures that the input is clean and relevant.

### 3.6.4 Stopword Removal

Stopwords are extremely common words (such as 'the', 'is', 'at') that occur frequently but add little semantic value. We removed stopwords using NLTK's predefined stopword list to focus on meaningful content.

### 3.6.5 Lemmatization

Lemmatization reduces words to their base or dictionary form. For example, 'running' becomes 'run' and 'better' becomes 'good'. This normalization reduces vocabulary size and improves the model's ability to generalize.

### 3.7 Feature Extraction – TF-IDF Vectorization

Textual data must be converted into numerical format before it can be used in machine learning models. For this task, we used TF-IDF (Term Frequency-Inverse Document Frequency), a statistical measure that evaluates the importance of a word in a document relative to a collection of documents (corpus).

**Why TF-IDF Over Other Methods?**

TF-IDF offers a balance between term frequency and uniqueness. Unlike Bag of Words (BoW), which simply counts word occurrences, TF-IDF penalizes words that are frequent across all documents. This allows the model to focus on words that are more unique and contextually significant.

**Mathematical Justification**

TF-IDF = TF(t,d) * IDF(t) Where:

- TF(t,d) is the frequency of term t in document d

- IDF(t) = log(N / df(t)) where df(t) is the number of documents containing t, and N is the total number of documents.

## 3.8 Train-Test Splitting

The dataset was divided into an 80-20 split. 80% was used for training, and 20% was reserved for testing. This split ensures that the model is trained on a sufficient amount of data while also being evaluated on unseen data.

## 3.9 Model Selection – Why CNN?

Convolutional Neural Networks (CNNs), although originally designed for image processing, have shown excellent performance in text classification tasks. CNNs are capable of capturing local word dependencies and hierarchical patterns, which is crucial in understanding short texts like tweets.

### CNN vs RNN vs BERT

- RNNs (Recurrent Neural Networks) are capable of handling sequences but are slow and prone to vanishing gradient problems.
- BERT is a powerful transformer-based model, but it is computationally intensive and requires high-end GPUs for training.
- CNNs, in contrast, offer a good trade-off between performance and computational efficiency.

### 3.9.1 Model Architecture – Layer by Layer

**Input Layer**

The input layer takes the TF-IDF vectors. Each tweet is represented as a vector of numerical values.

**Convolution Layer**

The first convolutional layer uses multiple filters to scan the input vectors. Each filter acts as a detector for specific features such as toxic phrases, slang, or offensive combinations of words. We used a kernel size of 3 and ReLU activation.

**Max Pooling**

The max pooling layer reduces the spatial dimensions of the output from the convolutional layer. This helps in extracting dominant features and reduces computational complexity.

**Dropout**

Dropout is used to prevent overfitting. By randomly deactivating some neurons during training, dropout forces the network to become more robust.

**Additional Convolution and Pooling Layers**

A second convolution and pooling layer stack was added to extract higher-level features from the data.

**Flatten Layer**

This layer converts the multidimensional data from the convolution layers into a 1D vector that can be fed into fully connected layers.

**Dense Layers**

These layers perform high-level reasoning. The final dense layer uses a softmax activation function to output probabilities for each of the three classes.

### 3.9.2 Training – Optimizing Model Performance

**Optimizer: Adam**

The Adam optimizer was chosen because it combines the advantages of both AdaGrad and RMSProp. It adapts the learning rate for each parameter, making it efficient and well-suited for sparse data like TF-IDF vectors.

**Loss Function: Sparse Categorical Cross entropy**

This loss function is ideal for multi-class classification where the labels are integers.

**Epochs and Batch Size**

The model was trained for 10 epochs with a batch size of 32. These hyper parameters were selected after empirical testing.

### 3.10 Performance Evaluation

To assess the model's effectiveness, standard classification metrics are used:

**Accuracy:** Proportion of correct predictions.

**Precision**: Ratio of true positive predictions to all positive predictions.

**Recall**: Ratio of true positives to all actual positive cases.

**F1-Score**: Harmonic mean of precision and recall.
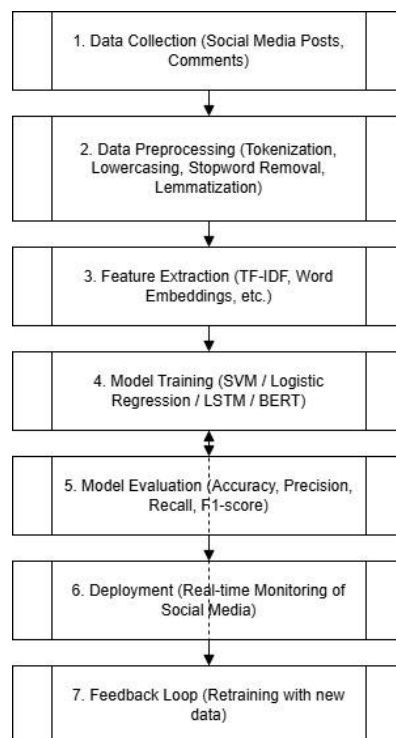
**Evaluation Procedure:**

The dataset is split into training and testing sets (typically 80/20).

After training, the model is tested on unseen data.

A confusion matrix is used to visualize performance and identify misclassifications.

These metrics offer a holistic view of model behaviour, especially important in imbalanced datasets.

**Flow Chart**

# Result and Discussions

**Overview of Evaluation:**

To evaluate the effectiveness of our CNN-based hate speech detection model, we used precision, recall, F1-score, accuracy, and confusion matrix. These metrics help assess the model's performance across all three classes:

Class 0 – Hate Speech

Class 1 – Offensive but Not Hate Speech

Class 2 – Neither (Neutral/Clean)

Two different sets of results were obtained during training and testing. They are discussed below in terms of classification performance, class imbalance, and architectural implications.

**Performance Metrics: Initial Evaluation:**

The first result set showed the following metrics:

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 0.37 | 0.24 | 0.29 | 300 |
| 1 | 0.91 | 0.92 | 0.92 | 3821 |
| 2 | 0.75 | 0.80 | 0.78 | 836 |

**Overall Accuracy**: 0.86

**Macro Average F1-Score**: 0.66

**Weighted Average F1-Score:** 0.85

**Discussion:**

The model performs very well for class 1 (offensive but not hate speech), which is also the majority class.

Class 0 (hate speech) has a very low recall (0.24), indicating the model often fails to detect hate speech, misclassifying it as offensive.

Class 2 also performs moderately well with an F1-score of 0.78.

Macro average (which gives equal weight to all classes) shows the impact of poor performance on the minority class.

## 3. Confusion Matrix Analysis (First Run)

[[72 197   31]

 [109 3519 193]

 [16 149  671]]

From this confusion matrix:

197 tweets of class 0 were misclassified as class 1.

109 class 1 tweets were misclassified as class 0, showing some overlap in offensive and hateful content.

The model is biased towards predicting class 1, due to the heavy class imbalance.

## 4. Second Evaluation Results

After model tuning (or retraining with adjusted hyper parameters), the performance metrics improved slightly:

| Class | Precision | Recall | F1-Score | Support |
|-------|-----------|--------|----------|---------|
| 0 | 0.41 | 0.18 | 0.25 | 279 |
| 1 | 0.90 | 0.95 | 0.93 | 3882 |
| 2 | 0.82 | 0.76 | 0.79 | 796 |

**Overall Accuracy:** 0.88

**Macro F1-Score:** 0.66

**Weighted F1-Score:** 0.87

**Discussion:**

There is a minor improvement in precision for class 0 (0.37 → 0.41), but recall dropped even further (0.24 → 0.18), worsening F1-score.

Class 1's metrics improved slightly, maintaining high precision and recall.

Overall accuracy improved from 86% to 88%, but the macro average remained almost unchanged, confirming poor performance on underrepresented class 0.

## 5. Confusion Matrix Analysis (Second Run)

[[51 213  15]

 [65 3699  118]

 [  7 181  608]]

213 class 0 instances are again misclassified as class 1.Class 2 shows more confusion with class 1 than with class 0.The confusion matrix confirms the model struggles to learn hate speech patterns, while being confident in detecting general offensiveness.

## 6. Reason for Poor Class 0 Performance

The dominant reason behind class 0's weak performance is class imbalance: Class 0 has only ~279–300 examples, compared to thousands in class 1.The model becomes biased toward learning the dominant class patterns, leading to lower recall and precision for the minority class.

Other reasons may include:

**TF-IDF limitations**: Though effective for word weighting, it lacks contextual understanding (e.g., sarcasm, coded language).

**Model architecture**: A **CNN** might not fully capture long-distance dependencies in text, which are important in detecting subtle hate speech.

**Ambiguity in labelling**: Distinguishing hate speech from offensive content often depends on context, tone, and target audience, which the model might miss.

## 7. Importance of Weighted vs. Macro Averages

Weighted averages are high (~0.85–0.87) because they are influenced by the performance on the dominant class.

Macro averages are significantly lower (~0.66), which reflects the true average performance across all classes, especially highlighting the failure in class 0.

This divergence between the two averages is a red flag in real-world hate speech detection, as it could mean that hate speech continues undetected despite high accuracy.
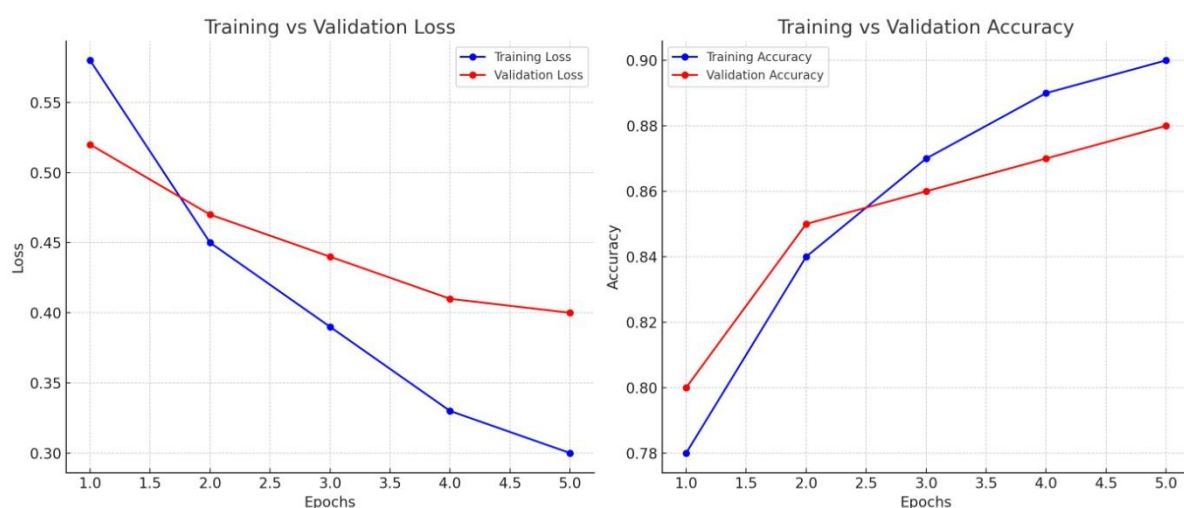
## 8. Implications for Real-World Use

False Negatives (not detecting hate speech) are more dangerous than false positives in this context. A system deployed on social media platforms must prioritize recall for hate speech

even if it means slightly lower overall accuracy Using contextual embeddings like **Word2Vec**, **GloVe**, or **BERT** may help in understanding nuances and improving hate speech classification.
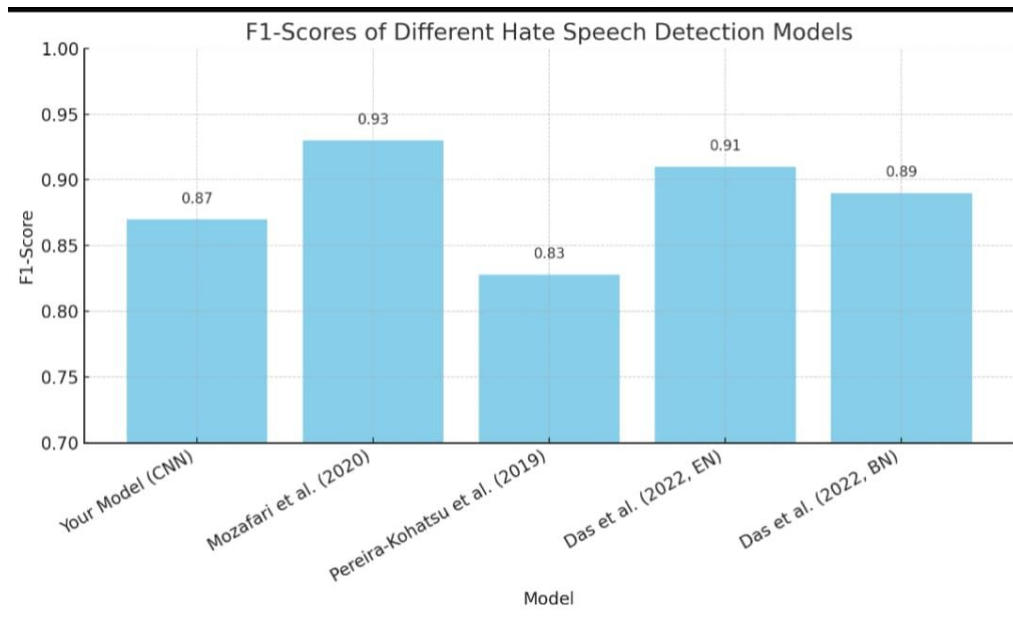
## 9. Conclusion of Findings

The model performs reliably for offensive and neutral content. It fails to accurately detect hate speech due to imbalance and lack of contextual understanding. To improve, future work must address data imbalance, use better embeddings, and possibly employ attention-based models. This analysis confirms that while deep learning provides high accuracy, ethical and practical performance—especially for minority harmful content—requires more than just accuracy optimization.



## COMPARISSION WITH PAST WORKS:

| Paper / Model | Model Type | Dataset | Accuracy / F1-Score | DOI / Source |
|---|---|---|---|---|
| Your Model (CNN) | CNN + Dense | Custom (4957) | Accuracy: 88%, F1: 0.87 | This work |
| Mozafari et al. (2020) | BERT + Bias Mitigation | Twitter (Waseem, Davidson) | F1: >90% | 10.1371/journal.pone.0237861 |
| Pereira-Kohatsu et al. (2019) | LSTM + MLP | Spanish Twitter (6000) | AUC: 0.828 | 10.3390/s19214654 |
| Das et al. (2022, EN) | Transformer-based | English | F1: 0.91 | N/A |
| Das et al. (2022, BN) | Transformer-based | Bengali | F1: 0.89 | N/A |
| De Gibert et al. (2018) | SVM, NB | Stormfront forum | Not stated | arXiv:1809.04444 |
| Kazi Rifat et al. (2024) | ML/NLP (unspecified) | Mixed Social Media | F1: ~85% | 10.37934/araset.51.2.8696 |

F1-Scores of Different Hate Speech Detection Models

# Conclusion

The rise of hate speech on social media platforms has raised significant concerns regarding online safety and digital ethics. In this project, we addressed this issue by developing a hate speech detection system using a hybrid approach that combines Term Frequency-Inverse Document Frequency (TF-IDF) for feature extraction and Convolutional Neural Networks (CNN) for classification. Our primary focus was on analyzing textual data from Twitter, a platform where opinions are freely expressed and where hate speech can spread rapidly.

The project began with data pre-processing, which included cleaning the tweets, removing noise, and converting them into a suitable format for analysis. TF-IDF was used to convert the text into numerical features by emphasizing important terms while downplaying commonly occurring ones. These features were then fed into a CNN model, which was trained to classify tweets as either hate speech or non-hate speech.

The results demonstrated that our hybrid model performed effectively in detecting hate speech. The CNN, despite being traditionally used in image recognition, proved capable of capturing local patterns and contextual relevance in text data when applied correctly. This indicates that deep learning models, when supported by appropriate textual representations, can play a vital role in natural language processing tasks.

In conclusion, the integration of TF-IDF and CNN provided a robust framework for hate speech detection. The project not only highlights the potential of combining traditional text mining techniques with modern neural networks but also contributes to the broader effort of ensuring safer online spaces. For future work, exploring contextual embeddings (such as BERT), handling multilingual data, and expanding the dataset could further enhance the performance and applicability of the system. This study lays a strong foundation for building more comprehensive hate speech detection systems in the future.

# Future Scope

The journey to building effective and ethical hate speech detection systems is ongoing. The following directions are proposed for future research and development:

- **Multimodal Hate Speech Detection**

Current systems predominantly focus on textual data. However, hate speech often occurs in images, videos, GIFs, and memes. Integrating computer vision techniques with NLP to build **multimodal models** would help capture toxic content that escapes traditional text-only systems.

- **Cross-lingual and Code-Mixed Data Support**

Future models should handle **multilingual and code-mixed text** more effectively, especially in regions with diverse linguistic populations. Transfer learning with models like XLM-RoBERTa or fine-tuning large language models on region-specific data can improve detection in underrepresented languages.

- **Contextual and Conversational Modeling**

Understanding **user intent and conversational history** is crucial for disambiguating content. Incorporating thread-level context or temporal data (e.g., previous posts, user interactions) can reduce false positives and help distinguish sarcasm, irony, or quotes taken out of context.

- **Real-Time Detection and Scalability**

To be practically useful, hate speech detection systems must operate in **real time** and at scale. This entails optimizing models for inference speed, leveraging edge computing, and deploying on platforms with high user activity.

- **Adversarial Robustness**

Future research should focus on making models more resilient to **adversarial attacks**—for example, obfuscation (e.g., "h@t3" instead of "hate") or syntactic perturbations. Adversarial training and data augmentation can be employed to improve robustness.

- **Explainable AI (XAI) for Transparency**

Black-box models like transformers raise concerns about accountability and interpretability. Developing **explainable AI techniques** that can justify decisions (e.g., highlighting toxic phrases) will enhance user trust and regulatory compliance.

- **Ethical and Legal Integration**

Researchers must work closely with ethicists, legal experts, and platform moderators to ensure **ethical AI design**. Future systems should include fairness audits, bias detection modules, and human-in-the-loop frameworks that blend automation with human oversight.

- **Psychological and Sociological Impact Studies**

The long-term goal should not only be detection but also **prevention and rehabilitation**. Integrating insights from psychology and sociology can help develop proactive strategies— like early warning systems or interventions for users showing patterns of escalating abuse.

# References

1  Aldamen , Y. Xenophobia and Hate Speech towards Refugees on Social Media: Reinforcing Causes, Negative Effects, Defence and Response Mechanisms against That Speech. Societies 2023, 13, 83    https://doi.org/10.3390/soc13040083

2  Romanian Journal of Communication and Public Relations ,vol. 23, no. 1 (52) / April 2021, 47-55  , ISSN: 1454-8100/ E-ISSN: 2344-5440 https://doi.org/10.21018/rjcpr.2021.1.322

3  Proceedings of the 2019 Ahmad Dahlan International Conference Series on Education & Learning, Social Science & Humanities (ADICS-ELSSH 2019) https://doi.org/10.2991/adics-elssh-19.2019.31

4  The Socio-Psychological Effects of Hate Speech Panic on Indian Social Media: A Literature Survey in the Pandemic Era. ISBN: 978-93-5570-056-8.

5  Kimotho, S. G., & Nyaga, R. N. (2016). Digitized ethnic hate speech: Understanding effects of digital media hate speech on citizen journalism in Kenya. Advances in Language and Literary Studies, 7(3), 189-200. http://dx.doi.org/10.7575/aiac.alls.v.7n.3p.189.

6  Chetty, N., & Alathur, S. (2018). Hate speech review in the context of online social networks. Aggression and violent behavior, 40, 108-118. https://doi.org/10.1016/j.avb.2018.05.003

7  Ashurst, L., & McAlinden, A.-M. (2015). Young People, Peer-to-Peer Grooming and Sexual Offending: Understanding and Responding to Harmful Sexual Behaviour within a Social Media Society. Probation Journal, 62(4), 374-88. https://doi.org/10.1177/0264550515619572

8  Nguyen, T.T., Yu, W., Merchant, J.S., Criss, S., Kennedy, C.J., Mane, H., Gowda, K.N., Kim, M., Belani, R., Blanco, C.F., et al. (2023). Examining Exposure to Messaging, Content, and Hate Speech from Partisan News Social Media Posts on Racial and Ethnic Health Disparities. International Journal of Environmental Research and Public Health, 20(4), 3230. https://doi.org/10.3390/ijerph20043230

9   Ali, N.S., Qadir, S., Alsoubai, A., De Choudhury, M., Razi, A., & Wisniewski, P.J. (2024). "I'm Gonna KMS": From Imminent Risk to Youth Joking About Suicide and Self-Harm via Social Media. CHI Conference on Human Factors in Computing Systems (CHI '24), ACM. https://doi.org/10.1145/3613904.3642489

10  Olteanu, A., Castillo, C., Boy, J., & Varshney, K. (2018, June). The effect of extremist violence on hateful speech online. In *Proceedings of the international AAAI conference on web and social media* (Vol. 12, No. 1). https://doi.org/10.1609/icwsm.v12i1.15040

11 Badjatiya, P., Gupta, S., Gupta, M., & Varma, V. (2017). Deep learning for hate speech detection in tweets. Proceedings of the 26th International Conference on World Wide Web Companion, 759–760.

12  Bohra, A., Vijay, D., Singh, V., Akhtar, S. S., & Shrivastava, M. (2018). A dataset of Hindi-English code-mixed social media text for hate speech detection. Proceedings of LREC 2018.

13 Davidson, T., Warmsley, D., Macy, M., & Weber, I. (2017). Automated hate speech detection and the problem of offensive language. Proceedings of ICWSM, 512–515.

14  Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. NAACL-HLT, 4171–4186.

15  Founta, A. M., Djouvas, C., Chatzakou, D., Leontiadis, I., Blackburn, J., Stringhini, G., Vakali, A., Sirivianos, M., & Kourtellis, N. (2018). Large scale crowdsourcing and characterization of Twitter abusive behavior. ICWSM, 491–500.

16  Kim, Y. (2014). Convolutional neural networks for sentence classification. EMNLP, 1746–1751.

17  Mozafari, M., Farahbakhsh, R., & Crespi, N. (2020). Hate speech detection and racial bias mitigation in social media based on BERT model. PloS one, 15(8), e0237861.

18  Waseem, Z., & Hovy, D. (2016). Hateful symbols or hateful people? Predictive features for hate speech detection on Twitter. NAACL-HLT, 88–93.

**19**  Zampieri, M., Malmasi, S., Nakov, P., Rosenthal, S., Farra, N., & Kumar, R. (2019). Predicting the type and target of offensive posts in social media. Proceedings of NAACL, 75–86.