

A

Mini Project Report on

CYBERSHIELD: ML-POWERED DETECTION OF SOCIAL MEDIA CYBERBULLYING

Submitted for partial fulfilment of the requirements for the award of the degree of

BACHELOR OF TECHNOLOGY

In

INFORMATION TECHNOLOGY

By

MOHAMMED KAIF 22K81A12G2

M. TARUN 22K81A12F7

Under the Guidance of

MRS. K. SURYA KANTHI

Assistant Professor



DEPARTMENT OF INFORMATION TECHNOLOGY

St. MARTIN'S ENGINEERING COLLEGE

UGC Autonomous

Affiliated to JNTUH, Approved by AICTE

Accredited by NBA & NAAC A+, ISO 9001-2008 Certified

Dhulapally, Secunderabad-500 100

www.smec.ac.in

JUNE – 2025



St. MARTIN'S ENGINEERING COLLEGE

UGC Autonomous
Affiliated to JNTUH, Approved by AICTE
NBA & NAAC A+ Accredited Dhulapally,
Secunderabad - 500 100
www.smec.ac.in



Certificate

This is to certify that the project entitled “**Cybershield: MI-Powered Detection Of Social Media Cyberbullying**” is being submitted by Mohammed kaif (22K81A12G2), M. Tarun (22K81A12F7), in fulfilment of the requirement for the award of degree of **BACHELOR OF TECHNOLOGY** in **INFORMATION TECHNOLOGY** is recorded of Bonafide work carried out by them. The result embodied in this report have been verified and found satisfactory.

Signature of Guide

Mrs. K. Surya Kanthi

Assistant Professor

Department of Information Technology

Signature of HOD

Dr. N. Krishnaiah

Professor and Head of Department

Department of Information Technology

Internal Examiner

External Examiner

Place:

Date:



St. MARTIN'S ENGINEERING COLLEGE

UGC Autonomous
Affiliated to JNTUH, Approved by AICTE
Accredited NBA & NAAC A+ ISO 9001:2008 Certified
Dhulapally, Secunderabad - 500 100
www.smec.ac.in



DEPARTMENT OF INFORMATION TECHNOLOGY

DECLARATION

We, the students of “**Bachelor of Technology in Department of Information Technology**”, session: 2022 - 2026, **St. Martin's Engineering College, Dhulapally, Kompally, Secunderabad**, hereby declare that the work presented in this project work entitled **Cybershield: MI-Powered Detection Of Social Media Cyberbullying** is the outcome of our own bonafide work and is correct to the best of our knowledge and this work has been undertaken taking care of Engineering Ethics. This result embodied in this project report has not been submitted in any university for award of any degree.

MOHAMMED KAIF

22K81A12G2

M. TARUN

22K81A12F7

UGC AUTONOMOUS

ACKNOWLEDGEMENT

The satisfaction and euphoria that accompanies the successful completion of any task would be incomplete without the mention of the people who made it possible and whose encouragement and guidance have crowded our efforts with success.

First and foremost, we would like to express our deep sense of gratitude and indebtedness to our College Management for their kind support and permission to use the facilities available in the Institute.

We especially would like to express our deep sense of gratitude and indebtedness to **Dr. P. SANTOSH KUMAR PATRA**, Professor and Group Director, St. Martin's Engineering College, Dhulapally, Secunderabad, for permitting us to undertake this project.

We wish to record our profound gratitude to **Dr. M. SREENIVAS RAO**, Principal, St. Martin's Engineering College, for his motivation and encouragement.

We are also thankful to **DR. N. KRISHNAIAH**, Head of the Department, Information Technology, St. Martin's Engineering College, Dhulapally, Secunderabad, for his support and guidance throughout our project as well as Project Coordinator **MRS. K. SURYA KANTHI**, Assistant Professor, Information Technology department for his valuable support.

We would like to express our sincere gratitude and indebtedness to our project supervisor **MRS. K. SURYA KANTHI**, Assistant Professor, Information Technology, St. Martins Engineering College, Dhulapally, for his support and guidance throughout our project.

Finally, we express thanks to all those who have helped us successfully completing this project. Furthermore, we would like to thank our family and friends for their moral support and encouragement. We express thanks to all those who have helped us in successfully completing the project.

MOHAMMED KAIF	22K81A12G2
M. TARUN	22K81A12F7

ABSTRACT

CyberShield is an advanced machine learning-powered system developed to detect and mitigate cyberbullying on social media platforms. Unlike traditional approaches such as manual moderation and simple keyword-based filters, which often fall short due to limited contextual understanding and high false-positive rates, CyberShield leverages cutting-edge technologies to address these limitations. By integrating Natural Language Processing (NLP), deep learning algorithms, and sentiment analysis, it achieves a more nuanced and accurate detection of abusive content. The system is trained on large and diverse datasets, enabling it to identify a wide range of harmful behaviors—including abusive language, personal threats, and various forms of online harassment—with high precision. Its intelligent architecture not only reduces the occurrence of misclassification but also adapts over time to emerging patterns of cyberbullying. Scalable and efficient, CyberShield aims to create a safer, more inclusive, and supportive digital environment for users across all age groups and communities.

The logo features a circular emblem with a gear-like border. Inside the circle, the text "ENGINEERING" is at the top, "TECHNOLOGY FOR PROSPERITY" is at the bottom, and "UGC" is in the center. Below the emblem is a red ribbon banner with the text "UGC AUTONOMOUS" in white capital letters.

UGC AUTONOMOUS

LIST OF ACRONYMS AND DEFINITIONS

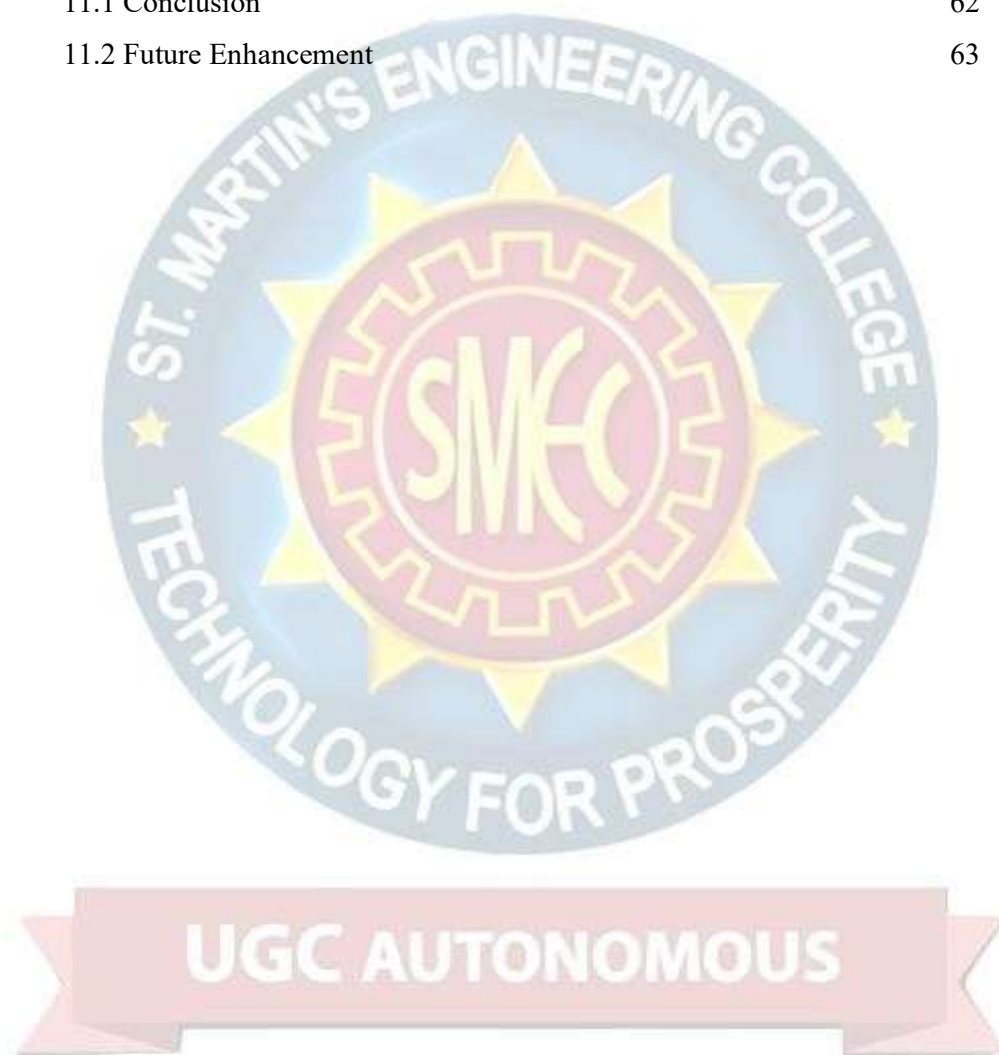
S.NO	ACRONYM	DEFINITION
01.	NLP	Natural Language Processing
02.	DNN	Deep Neural Network
03.	API	Application Program Interface
04.	CNN	Convolutional Neural Networks
05.	RNN	Recurrent Neural Networks



CONTENTS

ACKNOWLEDGEMENT	i
ABSTRACT	ii
LIST OF ACRONYMS AND DEFINITIONS	iii
CHAPTER 1 INTRODUCTION	01
1.1 Objective	02
1.2 Meaning of Title	03
CHAPTER 2 HISTORY	05
CHAPTER 3 PROBLEM STATEMENT	08
3.1 Existing System	10
3.2 Proposed System	11
CHAPTER 4 RESEARCH MOTIVATION	15
4.1 Impact And Benefits	18
4.2 Module Split	19
4.2.1 1. Data Collection	19
4.2.2 2. Data Preprocessing	20
4.2.3 3. Labeling	21
4.2.4 4. Feature Extraction	21
4.2.5 5. Model Selection	22
4.2.6 6. Model Training	23
4.2.7 7. Model Evaluation	24
4.2.8 8. Integration	24
CHAPTER 5 PROPOSED ALGORITHM	26
CHAPTER 6 MACHINE LEARNING IMPLEMENTATION	30
CHAPTER 7 DIAGRAMS	31

7.1 Detail Design	32
CHAPTER 8 ALGORITHM	34
CHAPTER 9 SOURCE CODE	37
CHAPTER 10 RESULTS OBTAINED	58
CHAPTER 11 CONCLUSION &FUTURE ENHANCEMENT	62
11.1 Conclusion	62
11.2 Future Enhancement	63



CHAPTER 1

INTRODUCTION

CYBERSHIELD: ML-POWERED DETECTION OF SOCIAL MEDIA CYBERBULLYING

The objective of the project "CyberShield: ML-Powered Detection of Social Media Cyberbullying" is to develop an advanced system leveraging machine learning techniques to identify and mitigate cyberbullying on social media platforms. The system is designed to enhance detection accuracy, minimize false positives, and provide actionable insights for combating online harassment. The title "CyberShield" reflects the project's core mission—to act as a protective shield against cyberbullying by utilizing cutting-edge machine learning models to analyze and identify harmful online interactions. Traditional cyberbullying detection methods relied heavily on manual moderation and basic keyword-based filtering systems, both of which had significant limitations. Human moderators reviewing flagged content made the process time-consuming and prone to oversight, while keyword-based filters lacked contextual understanding, often leading to misclassification of content. The vast and dynamic nature of online communication, with its evolving slang and nuanced expressions, made traditional methods ineffective in accurately identifying cyberbullying incidents. This project is motivated by the increasing prevalence of cyberbullying and its negative impact on mental health. With the growing volume of user-generated content, there is an urgent need for automated, intelligent solutions that can detect harmful interactions more efficiently and effectively. By integrating machine learning algorithms, CyberShield aims to overcome the shortcomings of traditional methods by enabling a deeper understanding of textual data, including context, sentiment, and intent. The proposed system employs supervised learning techniques, training models on extensive datasets of labeled social media interactions to recognize subtle patterns indicative of cyberbullying.

1.1 OBJECTIVE

The primary objective of the project "**CyberShield: ML-Powered Detection of Social Media Cyberbullying**" is to develop a comprehensive and intelligent system that leverages the power of machine learning to detect, analyze, and mitigate cyberbullying across various social media platforms. In an age where digital interactions are an integral part of daily life, safeguarding users—especially vulnerable groups such as teenagers and young adults—from online harassment has become a critical necessity. CyberShield is envisioned as a proactive solution that not only identifies harmful and abusive content in real-time but also provides insights that help moderators and administrators take appropriate actions to prevent the escalation of such behavior.

This project aims to address the shortcomings of traditional cyberbullying detection methods by integrating advanced machine learning algorithms that can understand the context, sentiment, and intent behind online messages. The goal is to move beyond simple keyword-based flagging mechanisms and towards intelligent classification models that can distinguish between casual banter and genuine harassment with high precision. Additionally, the system seeks to minimize false positives and false negatives, ensuring that user experiences are not negatively affected by misclassification.

Another key objective is to design a solution that is **scalable and adaptable**, capable of handling massive volumes of user-generated content and evolving alongside the dynamic nature of online language. By training on large, diverse, and labeled datasets, the system will learn to identify complex patterns, slang, and coded language used by cyberbullies. Ultimately, CyberShield strives to empower platform administrators, educators, and guardians with an effective tool to foster safer online environments and contribute to the well-being of digital communities worldwide.

1.2 MEANING OF TITLE

The title "**CyberShield**" is a symbolic representation of the project's core mission: to act as a protective barrier against cyberbullying in the online space. The term "**Shield**" reflects the system's role as a safeguard for individuals who are vulnerable to online harassment. Just as a physical shield provides protection from harmful forces, CyberShield aims to shield users from the emotional and psychological harm caused by cyberbullying. The title evokes a sense of security and resilience, highlighting the project's commitment to creating a safer digital environment.

The term "**Cyber**" is directly related to the online and digital space, emphasizing that the project specifically targets issues arising from interactions on social media and other online platforms. In the digital age, where communication is primarily conducted through virtual means, the need for protection from cyberbullying becomes paramount, and the title reflects this urgent necessity.

The phrase "**ML-Powered Detection**" underscores the technological approach the project adopts. By leveraging **Machine Learning (ML)**, the system can intelligently analyze and interpret large amounts of social media data to identify cyberbullying behavior. Unlike traditional methods that rely on pre-defined rules or keyword-based filtering, machine learning algorithms can adapt and evolve over time, enabling the system to detect more nuanced and context-sensitive forms of bullying. This power of machine learning allows CyberShield to continuously improve its detection capabilities and accurately classify abusive content across different platforms, languages, and contexts.

Finally, the "**Social Media Cyberbullying**" part of the title makes it clear that the focus of the project is on combating bullying specifically within the realm of social media. Social media platforms, with their vast user bases and frequent interactions, have become a breeding ground for cyberbullying, making the

detection and prevention of harmful behavior on these platforms an essential area of focus. By specifically targeting social media, the project aligns itself with one of the most significant challenges in the digital world today.

In summary, the title "**CyberShield: ML-Powered Detection of Social Media Cyberbullying**" encapsulates the project's goal to use cutting-edge machine learning techniques to safeguard users from cyberbullying on social media, offering an intelligent, adaptable, and scalable solution to this pressing issue.

CHAPTER 2

HISTORY

The history of cyberbullying detection has evolved alongside the rapid growth of the internet and social media platforms. In the early stages of the internet, bullying primarily occurred in physical spaces, such as schools and neighborhoods. However, as online communities and digital communication platforms like forums, chat rooms, and early social networks began to proliferate in the late 1990s and early 2000s, a new form of bullying emerged: **cyberbullying**. This new form of bullying was more pervasive and, in many cases, more harmful due to its anonymous nature and the ability to reach a wider audience instantly.

In the early days of cyberbullying detection, the methods used were relatively primitive and heavily reliant on human intervention. Initially, social media platforms and websites had rudimentary systems that allowed users to report abusive content. This was a manual process, where **moderators** or **community managers** were responsible for reviewing reported incidents. While this method helped identify some harmful behavior, it was highly inefficient as it couldn't scale to handle the massive volume of content being shared online. Furthermore, this approach had its limitations in recognizing subtle or context-sensitive forms of bullying, which often required human interpretation and emotional intelligence.

The next phase in the history of cyberbullying detection saw the implementation of **keyword-based filtering systems**. These filters relied on a list of pre-defined keywords or phrases that were associated with bullying, hate speech, or inappropriate behavior. While keyword filters were effective to some extent in flagging abusive content, they lacked the ability to understand **context**. For instance, a phrase like "you're an idiot" could easily be flagged, but without understanding the context in which it was used, such as a friendly joke between

close friends, these filters would frequently lead to **false positives**—flagging content that was not actually harmful. Similarly, the filters were often ineffective at detecting more nuanced forms of bullying that didn't necessarily include offensive words but involved insults, threats, or manipulative behavior.

The limitation of traditional detection methods became even more apparent with the rise of **social media platforms** like Facebook, Twitter, and Instagram, where the amount of user-generated content increased exponentially. This volume made it nearly impossible for human moderators to keep up with the sheer scale of online interactions. Furthermore, the rise of **new forms of online abuse**, such as trolling, doxxing, and harassment through memes or images, presented challenges that simple keyword-based filters couldn't address.

With the realization that traditional methods were inadequate, researchers and developers began exploring **Machine Learning And Artificial Intelligence (AI)** as more effective solutions for detecting cyberbullying. The first attempts to use machine learning in cyberbullying detection focused on **text classification techniques**, using labeled datasets of online content to train algorithms to differentiate between harmful and non-harmful interactions. These early models used relatively simple algorithms like **Naive Bayes** or **Support Vector Machines (SVMs)**. While these models showed promise, they still struggled with the complex and dynamic nature of human language, including the **use of slang, sarcasm, and contextual ambiguity**.

The advancement of more sophisticated machine learning models, particularly **deep learning**, marked the next milestone in cyberbullying detection. Deep learning models, particularly **Convolutional Neural Networks (CNNs)** and **Recurrent Neural Networks (RNNs)**, brought significant improvements in understanding the **semantics and context** of online content. These models could process large datasets and learn complex patterns in text that were previously difficult for traditional models to capture. The application of **Natural Language**

Processing (NLP) techniques, such as **sentiment analysis** and **entity recognition**, further improved the detection of abusive language, even when it was conveyed in subtle or indirect ways.

As cyberbullying detection systems continued to evolve, there was a growing recognition that these systems needed to be **scalable, adaptable, and capable of handling diverse forms of abuse** across various platforms and languages. Today, advanced machine learning models are capable of not only detecting verbal abuse but also understanding the **social context** in which bullying occurs. For example, some systems now take into account the relationship between the people involved, the history of their interactions, and even the **behavioral patterns** of the individuals to make more accurate assessments.

In conclusion, the history of cyberbullying detection has evolved from simple manual moderation to sophisticated AI-driven models. What started as a manual, labor-intensive process has transformed into a dynamic field of research focused on using machine learning to tackle one of the most pressing challenges in online safety. With the continued advancements in AI and machine learning, it is hoped that systems like **CyberShield** will provide a comprehensive and proactive solution to combat cyberbullying, helping to create safer online spaces for all users.

CHAPTER 3

PROBLEM STATEMENT

Cyberbullying has become a significant social issue in the digital age, with the rise of social media and online platforms providing a fertile ground for harmful and abusive behavior. The increasing prevalence of cyberbullying has raised serious concerns about the psychological and emotional impact it has on individuals, particularly among vulnerable groups such as teenagers, young adults, and marginalized communities. Unlike traditional forms of bullying, cyberbullying can occur anonymously, at any time, and on any platform, making it difficult for victims to escape or protect themselves. The consequences of cyberbullying can be severe, ranging from anxiety, depression, and low self-esteem to more tragic outcomes like self-harm and suicide.

Despite the growing awareness of cyberbullying's detrimental effects, effective detection and mitigation systems remain limited. Traditional methods for identifying cyberbullying have primarily relied on **manual moderation** and **keyword-based filters**, both of which have significant limitations. Manual moderation is labor-intensive and not scalable, requiring human intervention to review potentially harmful content, which often leads to delays and oversights. Keyword-based systems, on the other hand, can only identify specific words or phrases associated with bullying, but they fail to capture the broader context in which the interactions occur. This results in frequent **false positives** (flagging benign content as abusive) and **false negatives** (failing to detect actual harmful behavior), leading to inefficiencies in addressing the issue.

Furthermore, the complexity of human language, including the use of sarcasm, slang, and subtle forms of manipulation, presents a significant challenge for traditional detection systems. As online interactions continue to evolve, cyberbullies are becoming increasingly sophisticated in masking their abusive

behavior, making it even harder for existing systems to accurately identify cyberbullying. This presents an urgent need for a more **automated, intelligent, and scalable solution** that can efficiently detect cyberbullying in real-time and across multiple platforms.

The **problem** addressed by this project is the **ineffectiveness and inefficiency of current cyberbullying detection systems**. Traditional approaches are not equipped to handle the dynamic nature of online communication, leading to delays, misclassifications, and a lack of timely intervention. There is a critical need for **machine learning-based solutions** that can analyze large volumes of user-generated content, recognize patterns of abusive behavior, and provide a more accurate and nuanced understanding of cyberbullying incidents.

The goal of this research is to develop an intelligent system, **CyberShield**, that leverages **machine learning algorithms** to detect and mitigate cyberbullying on social media platforms. By focusing on the **context, sentiment, and intent** behind online interactions, the system aims to overcome the limitations of traditional methods, providing a **scalable, adaptive, and real-time solution** to combat cyberbullying and protect users from harm. Ultimately, the project seeks to create a safer online environment for all users by providing an automated system that can efficiently identify harmful content and facilitate timely intervention.

3.1 EXISTING SYSTEM

Traditional methods for cyberbullying detection primarily relied on manual moderation by social media companies and keyword-based filtering techniques. Manual moderation was slow, inconsistent, and costly, while keyword-based filters often failed to detect abusive content due to a lack of contextual understanding. These systems could not handle sarcasm, slang, emojis, and indirect threats, leading to both false positives and false negatives. The absence of real-time intervention further reduced their effectiveness, making them impractical for large-scale implementation.

3.2 PROPOSED SYSTEM

The proposed system, **CyberShield: ML-Powered Detection of Social Media Cyberbullying**, aims to revolutionize the way cyberbullying is detected and mitigated on social media platforms by utilizing advanced **machine learning** (ML) techniques. This system is designed to address the limitations of traditional detection methods such as manual reporting and keyword-based filters, providing a more **accurate, scalable, and efficient solution** to combat online harassment.

At its core, **CyberShield** uses **supervised machine learning algorithms**, particularly focusing on **Natural Language Processing (NLP)** and **deep learning** techniques, to analyze and process user-generated content on social media. By training the system on large, labeled datasets of online interactions, the model learns to recognize patterns, context, and nuanced behavior indicative of cyberbullying, such as **abusive language, threats, harassment**, and other harmful interactions.

Key Features of the Proposed System:

1. **Data Collection and Preprocessing:**

The first step in the **CyberShield** system involves the collection of large datasets consisting of text data from social media platforms. This data includes user comments, messages, posts, and interactions that are labeled as either **bullying** or **non-bullying**. The dataset is then preprocessed using **NLP techniques** such as **tokenization, stemming, and lemmatization**, which transform the raw text into a format that machine learning algorithms can effectively analyze.

2. **Feature Extraction:**

The system employs **feature extraction** methods such as **TF-IDF** (Term Frequency-Inverse Document Frequency) to convert the text data into

numerical representations, capturing the **relevant features** of the content. These features can include aspects such as the **frequency of aggressive words, emotional tone, context of interactions**, and other behavioral indicators of cyberbullying.

3. **Machine Learning Model Training:**

The core of the **CyberShield** system is its **machine learning model**, which is trained on the preprocessed data to identify subtle patterns that may indicate cyberbullying. Several types of models are explored, including **Support Vector Machines (SVM), Random Forest, and Deep Neural Networks (DNN)**. The training process involves feeding the labeled data into the model, allowing it to learn the relationships between the features and the presence or absence of cyberbullying behavior.

4. **Deep Learning for Sentiment Analysis and Contextual Understanding:**

One of the major innovations of **CyberShield** is its ability to understand the **context** in which a message is sent. Cyberbullying is often dependent on context—what may seem like a harmless comment in one scenario could be harmful in another. To address this, the system uses **deep learning techniques**, such as **Convolutional Neural Networks (CNN)** or **Recurrent Neural Networks (RNN)**, to capture the **context** and **sentiment** behind the interactions. This helps the model not only recognize offensive words but also understand the **intent** and **emotional tone** of the message.

5. **Real-Time Detection and Alert System:**

Once the model has been trained and is ready for deployment, **CyberShield** operates in **real-time** to monitor interactions across social media platforms. The system continuously scans content for signs of cyberbullying, sending **instant alerts** when harmful behavior is detected. This enables social media platforms to take immediate action, such as

flagging the content for review, issuing a warning to the user, or blocking the abusive account.

6. **Scalability** and **Adaptability:**

A crucial feature of **CyberShield** is its **scalability**. The system is designed to handle large volumes of data, making it suitable for deployment on major social media platforms such as Facebook, Twitter, Instagram, and YouTube, where millions of interactions occur every minute. Additionally, **CyberShield** can be **adapted** to different platforms by adjusting the model to recognize platform-specific language, slang, and behavioral patterns. The system can also learn and evolve over time, becoming more accurate and effective at detecting cyberbullying as new patterns emerge.

7. **Insights** and **Reporting:**

In addition to detection, **CyberShield** provides actionable insights to help platform administrators and moderators combat cyberbullying more effectively. The system generates **detailed reports** on cyberbullying incidents, including statistics on the frequency of abusive behavior, the types of bullying (e.g., verbal abuse, threats), and the demographics of both perpetrators and victims. This helps administrators understand trends, measure the effectiveness of their anti-cyberbullying policies, and take proactive measures to foster a safer online environment.

8. **User Feedback Loop for Continuous Improvement:**

CyberShield includes a feedback mechanism that allows users and moderators to flag false positives and false negatives. This feedback is used to continuously **retrain the system**, improving its performance over time and helping to fine-tune the detection algorithm. The system is designed to **learn from its mistakes**, minimizing false positives and ensuring that real instances of cyberbullying are detected accurately.

9. Technological Stack:

- **Programming Languages:** Python (for developing machine learning models, data preprocessing, and feature extraction)
- **Libraries/Frameworks:**
 - **Scikit-learn** for machine learning algorithms (SVM, Random Forest)
 - **TensorFlow/Keras** for deep learning models (CNN, RNN)
 - **NLTK** and **spaCy** for natural language processing
 - **Pandas** and **NumPy** for data manipulation and analysis
- **Data Sources:** Publicly available datasets (e.g., from Kaggle) or partnerships with social media platforms for real-time data.

CHAPTER 4

RESEARCH MOTIVATION

The motivation behind the development of the **CyberShield: ML-Powered Detection of Social Media Cyberbullying** system arises from the growing concern over the prevalence and impact of cyberbullying in the digital world. The rapid expansion of social media platforms and online communication has transformed how people interact, but it has also created new avenues for harmful behavior such as cyberbullying. While traditional bullying was confined to physical spaces like schools and neighborhoods, cyberbullying transcends geographical boundaries, making it more pervasive and challenging to address.

The rise of social media and online interactions has amplified the risks associated with cyberbullying, as users, particularly young people, are exposed to potential abuse on a daily basis. Unlike face-to-face bullying, which is often witnessed by others and can be interrupted or reported, cyberbullying can occur privately, without any immediate intervention, and can have lasting effects on the victim. The anonymity provided by the internet often emboldens perpetrators, leading to more aggressive behavior, and the rapid spread of harmful content can affect individuals on a global scale. This has significant consequences on the mental health and well-being of victims, with consequences including anxiety, depression, and even suicidal tendencies.

Given the scale of the problem, **traditional methods of cyberbullying detection**, such as manual reporting and keyword-based filters, are no longer sufficient. These methods are not only inefficient but also prone to errors, including false positives and false negatives, which can further exacerbate the problem. Manual moderation is resource-intensive and can't keep up with the vast amounts of content generated on social media platforms, while keyword-based systems fail to capture the **nuance of human language**, including sarcasm, context, and

evolving slang. This has led to a **significant gap in effective solutions** for combating cyberbullying in real-time.

This gap is what motivated the research and development of **CyberShield**, which aims to address these shortcomings by harnessing the power of **machine learning (ML)**. Machine learning algorithms, particularly **Natural Language Processing (NLP)** and deep learning models, have the potential to analyze textual data more deeply than traditional methods. These algorithms can be trained to understand the **context** behind messages, recognize subtle cues of **abusive behavior**, and learn from patterns in language use across diverse social media platforms. Unlike keyword-based systems, ML models can adapt to changing language and detect new forms of bullying, including those that don't rely on specific words or phrases.

Furthermore, the rising concern over the long-term impact of cyberbullying on mental health calls for a more proactive, scalable, and real-time solution. **CyberShield** leverages machine learning to enable automatic detection of cyberbullying incidents, offering timely alerts and intervention. This system can provide critical insights into abusive behavior patterns, allowing platforms to take preventive measures and minimize harm before it escalates. The ability to detect cyberbullying early on and at scale is crucial in curbing its spread and protecting users, especially those who are most vulnerable to online harassment.

In addition, the sheer volume of user-generated content on social media platforms necessitates the use of intelligent systems capable of processing large amounts of data quickly and accurately. The use of machine learning enables **CyberShield** to operate at scale, making it possible to monitor millions of interactions in real-time without the need for manual intervention. The scalability of the system ensures that it can be applied to various social media platforms, each with its own unique communication styles and challenges, providing a **unified solution** to combat cyberbullying across diverse online spaces.

The **motivation** behind this project is therefore driven by the urgent need to develop a **more effective, efficient, and scalable system** for detecting and combating cyberbullying. By utilizing machine learning, **CyberShield** offers a powerful tool to address the limitations of existing systems, improving the **accuracy** of detection and providing actionable insights that can help protect users from the harmful effects of cyberbullying.

4.1 IMPACT AND BENEFITS

The proposed **CyberShield** system offers numerous benefits for social media platforms, users, and society at large. By accurately detecting cyberbullying incidents in real-time, the system allows for timely intervention, reducing the harmful impact of online harassment. The system's scalability and adaptability make it suitable for global deployment, addressing the cyberbullying problem across different cultures, languages, and online environments. Moreover, the system's ability to analyze context and sentiment ensures that **false positives** are minimized, and **non-abusive interactions** are not wrongly flagged as harmful.

Ultimately, **CyberShield** aims to create a **safer online space**, where users can freely interact without the fear of being harassed or bullied, contributing to a more positive and supportive digital community.

4.2 MODULE SPLIT

The CyberShield system is designed to employ a comprehensive and systematic approach to detect cyberbullying on social media platforms. This process can be broken down into several key modules, each focusing on a different aspect of building and deploying the system. Below is a detailed elaboration of each module:

4.2.1 1. Data Collection

The foundation of any machine learning model is high-quality data. The Data Collection module is responsible for gathering a diverse, large-scale dataset that consists of both cyberbullying and non-cyberbullying content. This data can come from a variety of social media platforms such as Twitter, Facebook, Instagram, and YouTube, and it should include text data from user posts, comments, messages, and replies.

- **Data Sources:** Publicly available datasets such as those from Kaggle, academic research, or partnerships with social media platforms can be used. If partnerships are not feasible, web scraping tools (e.g., BeautifulSoup, Scrapy) can be employed to collect publicly available social media posts.
- **Types of Data:** The dataset will consist of both text-based data and metadata, including post content, timestamps, user details, the number of likes or replies, etc. This helps in understanding not only the content but also the interactions surrounding it.
- **Diversity:** It is essential to ensure the data collected spans multiple languages, platforms, and user demographics to make the system robust

and applicable globally. This helps avoid biases and ensures the model is capable of detecting cyberbullying across various contexts.

4.2.2 2. Data Preprocessing

Once the data has been collected, the next crucial step is Data Preprocessing. Raw social media data can be noisy and unstructured, which makes it essential to clean and standardize the content for machine learning models to function effectively.

- **Handling Missing Values:** Missing data may be present due to incomplete posts, deleted content, or failed collection attempts. This can be handled through techniques such as imputation (filling missing values) or removal of records with missing data.
- **Noise Removal:** Noise refers to irrelevant or extraneous content such as URLs, hashtags, mentions, special characters, and emojis. These elements may distort the actual content of the post, so they must be removed or normalized.
- **Text Normalization:** Text normalization ensures that the data is in a consistent format. This includes:
 - Lowercasing all text to ensure uniformity.
 - Removing stop words (common words like "the", "is", "and") that do not contribute much to meaning.
 - Stemming or lemmatization to reduce words to their root forms (e.g., "running" becomes "run").
- **Tokenization:** Breaking down the text into smaller components like words, phrases, or sentences, known as tokens, to allow the model to process and understand the structure of the data.

4.2.3 3. Labeling

The Labeling module is essential for supervised machine learning. Here, the dataset is manually annotated to distinguish between cyberbullying and non-cyberbullying content. This process involves:

- **Manual Annotation:** Experts or crowd-workers go through the dataset and label each post as either bullying or non-bullying. For example, posts that contain offensive language, threats, insults, or harassment would be marked as cyberbullying, while others would be labeled as non-bullying.
- **Contextual Considerations:** Labeling is not only about identifying offensive words but also understanding context. For example, a sarcastic remark might appear as a bully's comment in some contexts but could be harmless in others. This requires expert judgment and understanding of the tone and sentiment behind each interaction.
- **Quality Control:** Since labeling can be subjective, a double-checking process might be needed to ensure the labels are accurate. Discrepancies in labeling can be resolved by reviewing the post with multiple annotators or by creating a consensus.

4.2.4 4. Feature Extraction

Feature Extraction is the process of transforming the raw textual data into a structured form that machine learning algorithms can process. The aim is to identify and extract relevant features that are indicative of cyberbullying behavior.

- **N-grams:** N-grams are contiguous sequences of 'n' words in the text. By extracting unigrams (single words), bigrams (pairs of words), or trigrams

(triplets of words), the model can understand word combinations that are indicative of bullying behavior.

- **Sentiment Scores:** Sentiment analysis provides a way to capture the emotional tone of a post. Posts that contain negative sentiments, anger, or aggression may indicate cyberbullying behavior. Sentiment scores can be extracted using pre-trained models or lexicons like VADER or TextBlob.
- **Part-of-Speech (POS) Tags:** Analyzing the grammatical structure of the text through POS tagging can help identify abusive language patterns (e.g., the use of imperatives, insults, or threats).
- **TF-IDF (Term Frequency-Inverse Document Frequency):** This technique weighs the importance of words in a document relative to the entire corpus. It helps identify key words and phrases that are likely to be significant in detecting bullying behavior.
- **Emotion Detection:** By identifying words or phrases that express specific emotions like anger, frustration, or sadness, the system can capture the emotional intent behind a message.

4.2.5 5. Model Selection

In the Model Selection module, different machine learning algorithms are evaluated to determine which is best suited for the task of cyberbullying detection. Several models are typically considered, such as:

- **Support Vector Machines (SVM):** SVM is a powerful classification algorithm that works well with text data. It tries to find the hyperplane that best separates the bullying and non-bullying data points in a multi-dimensional feature space.

- **Random Forest:** A Random Forest classifier combines the predictions of several decision trees, providing robustness against overfitting. This model works well for text classification tasks and handles both large datasets and noisy data.
- **Logistic Regression:** A linear model often used for binary classification tasks. It is fast, interpretable, and works well for feature-rich datasets.
- **Deep Learning Models (RNN, CNN):** For more complex tasks, Recurrent Neural Networks (RNN) and Convolutional Neural Networks (CNN) can be applied. These models are particularly good for sequential data like text and are able to capture long-range dependencies in the text.

4.2.6 6. Model Training

Once the model has been selected, the Model Training phase begins. During this stage, the labeled data is fed into the chosen model to allow it to learn the patterns and features associated with cyberbullying. The steps involved include:

- **Training Process:** The model is trained using a portion of the dataset, which is divided into training and validation sets. The training set allows the model to learn, while the validation set is used to tune hyperparameters and test the model's ability to generalize to unseen data.
- **Hyperparameter Tuning:** The performance of machine learning models can often be improved by adjusting hyperparameters such as learning rate, regularization strength, and the number of layers (in deep learning). Techniques like grid search or random search can be employed to find the optimal hyperparameters.

4.2.7 7. Model Evaluation

After training, the Model Evaluation module is crucial to measure the model's effectiveness in detecting cyberbullying. This evaluation is done using several key metrics, such as:

- **Precision:** The ratio of true positive predictions to all positive predictions made. This indicates how accurate the model is when it predicts an instance of cyberbullying.
- **Recall:** The ratio of true positive predictions to all actual positive instances in the data. Recall measures the model's ability to capture all instances of cyberbullying.
- **F1-Score:** The harmonic mean of precision and recall. The F1-score gives a more balanced view of the model's performance when there is a class imbalance (more non-bullying content than bullying content).
- **Accuracy:** The overall percentage of correct predictions made by the model.
- **Confusion Matrix:** A confusion matrix is used to visualize the performance of the classification model, helping to identify false positives, false negatives, true positives, and true negatives.

4.2.8 8. Integration

The final step in the CyberShield system is Integration, where the trained model is deployed for real-time monitoring of social media platforms. This module involves:

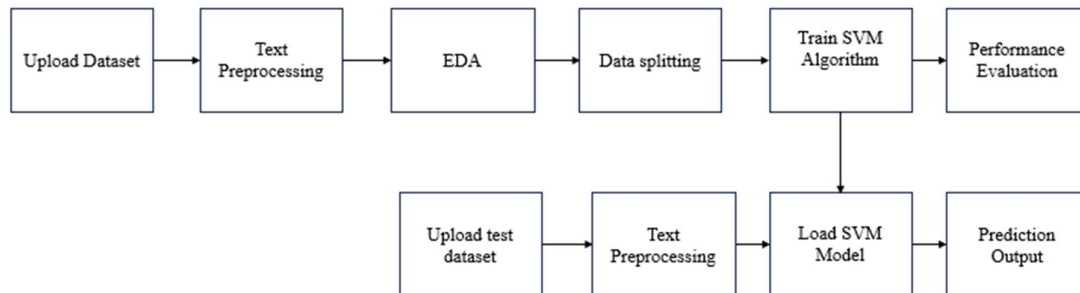
- **API Development:** Developing APIs that integrate the machine learning model with social media platforms. These APIs will allow the platform to

send user-generated content to the model for analysis and receive predictions in return.

- **Real-Time Processing:** The system needs to continuously monitor posts, comments, and messages in real time. Once a post is flagged for potential cyberbullying, the system can automatically send alerts to moderators or initiate automated responses, such as content removal or warnings.
- **User Interface:** A user interface can be developed for moderators and administrators to review flagged content, view analysis reports, and take necessary actions. The interface can also provide insights into the effectiveness of the system and track trends in cyberbullying over time.

CHAPTER 5

PROPOSED ALGORITHM



Block Diagram

ML Model Building

Existing Algorithm: Rule-Based Algorithm

What is a Rule-Based Algorithm?

A Rule-Based Algorithm is a type of decision-making system that relies on a set of predefined rules to make decisions or classifications. These rules are typically created based on domain knowledge, expert input, or empirical data. Rule-based systems operate on the principle of applying these rules to input data to derive conclusions or outputs.

How It Works

Rule-Based Algorithms work by evaluating input data against a set of "if-then" rules. Each rule specifies a condition and a corresponding action. For instance, in a rule-based classification system for diagnosing diseases, a rule might state: "If symptoms include fever and cough, then classify as flu." The system processes the input data, checks it against the rules, and applies the appropriate rule to generate the output. The decision-making process is transparent, as the rules explicitly define how inputs are mapped to outputs.

Architecture

The architecture of a Rule-Based System typically includes the following components:

Knowledge Base: Contains the set of rules and facts that define the system's behavior.

Inference Engine: Processes the input data and applies the rules from the knowledge base. It determines which rules are applicable and executes them to derive conclusions.

User Interface: Allows users to interact with the system, input data, and receive results.

Working Memory: Stores the current state of data and intermediate results during the decision-making process.

Disadvantages

Scalability Issues: As the number of rules grows, the complexity of managing and maintaining the rule base increases, leading to potential performance degradation.

Limited Adaptability: Rule-Based Systems are rigid and cannot easily adapt to new or unforeseen scenarios that are not covered by the existing rules.

Knowledge Acquisition Bottleneck: Creating and updating rules requires significant domain expertise and can be time-consuming.

Overfitting to Rules: The system may overfit to the specific rules defined, making it less effective for cases that deviate from the rule set.

Difficulty in Handling Uncertainty: Rule-Based Systems struggle with handling ambiguous or uncertain information, as they rely on deterministic rules.

Proposed Algorithm: Support Vector Machine (SVM)

What is SVM?

Support Vector Machine (SVM) is a supervised machine learning algorithm used for classification and regression tasks. It is designed to find the optimal hyperplane that separates data points of different classes with the maximum margin. SVM is particularly effective in high-dimensional spaces and for cases where the number of dimensions exceeds the number of samples.

How It Works

SVM works by transforming the input data into a higher-dimensional space using a kernel function, allowing it to find a hyperplane that best separates the different classes. The algorithm aims to maximize the margin between the closest points of different classes, known as support vectors. The process involves the following steps:

Transformation: Apply a kernel function to map the data into a higher-dimensional space.

Hyperplane Calculation: Find the optimal hyperplane that maximizes the margin between different classes in this transformed space.

Classification: Classify new data points based on which side of the hyperplane they fall on.

Architecture

The architecture of SVM includes:

Kernel Function: A mathematical function that transforms the input data into a higher-dimensional space (e.g., linear, polynomial, radial basis function (RBF)).

Optimization Algorithm: Solves the optimization problem to find the hyperplane that maximizes the margin between classes.

Support Vectors: Data points that are closest to the hyperplane and are critical for defining the margin.

Advantages

Effective in High-Dimensional Spaces: SVM performs well in scenarios where the number of features exceeds the number of samples.

Robust to Overfitting: By focusing on the support vectors and maximizing the margin, SVM reduces the risk of overfitting.

Versatile with Kernel Functions: The use of different kernel functions allows SVM to handle a variety of data distributions and non-linear relationships.

Clear Margin of Separation: SVM provides a clear margin of separation between classes, which can enhance interpretability and confidence in the classification.

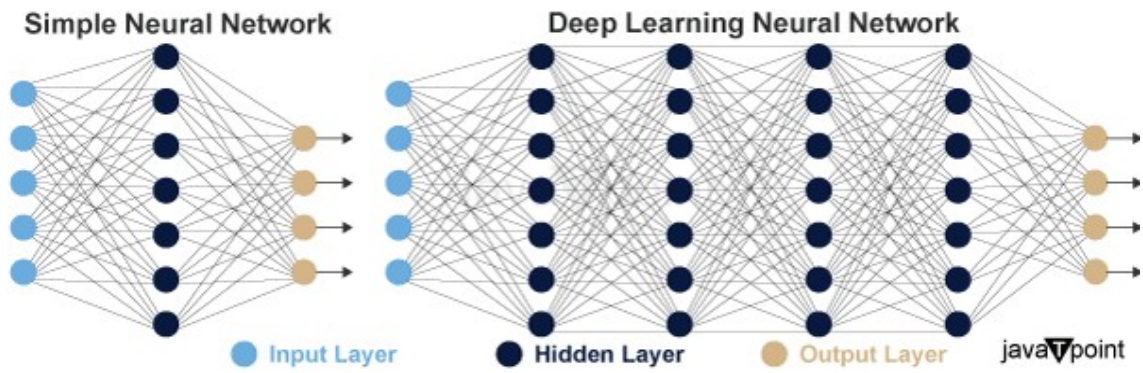
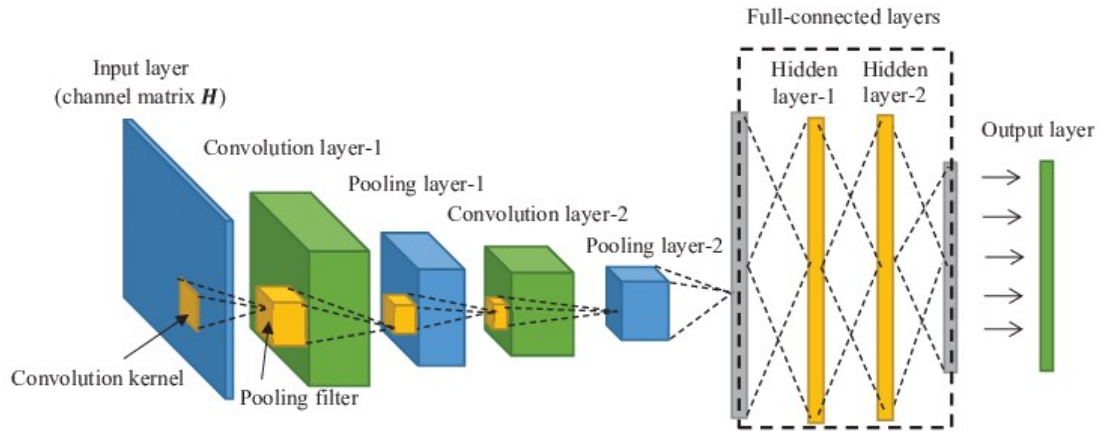
CHAPTER 6

MACHINE LEARNING IMPLEMENTATION

To develop a cyberbullying detection system, social media comments, posts, and messages are collected and labeled as either cyberbullying or non-cyberbullying. The text data is then cleaned by removing noise, special characters, and stopwords, followed by tokenization and vectorization. Techniques like TF-IDF and Word2Vec are used to extract meaningful features from the text. Supervised learning models such as LSTM and CNN are trained on this data to learn patterns of cyberbullying. After training, the model can predict whether a new message contains cyberbullying content. Finally, the system is connected through APIs for real-time monitoring and automatic detection of harmful messages.

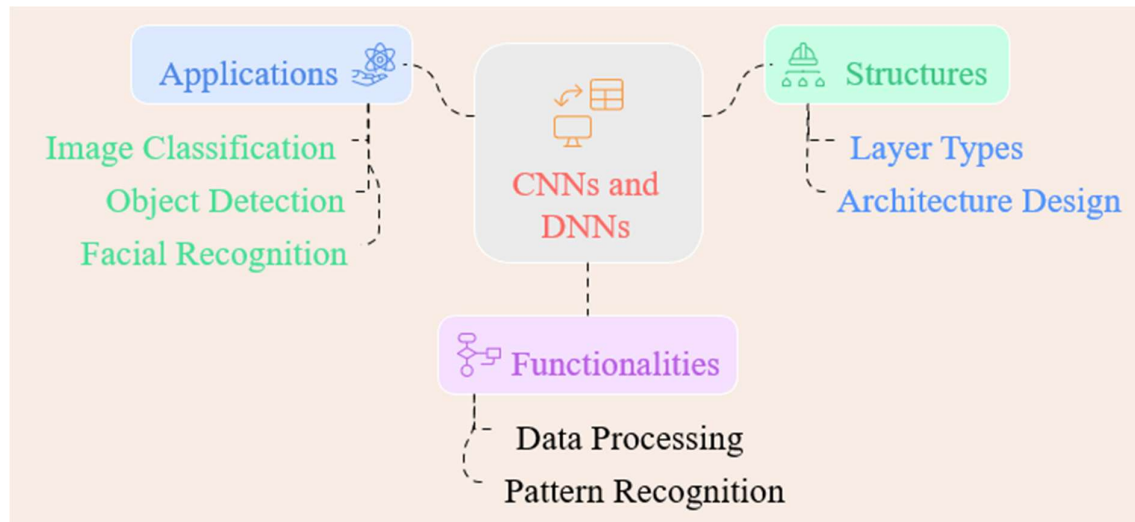
CHAPTER 7

DIAGRAMS

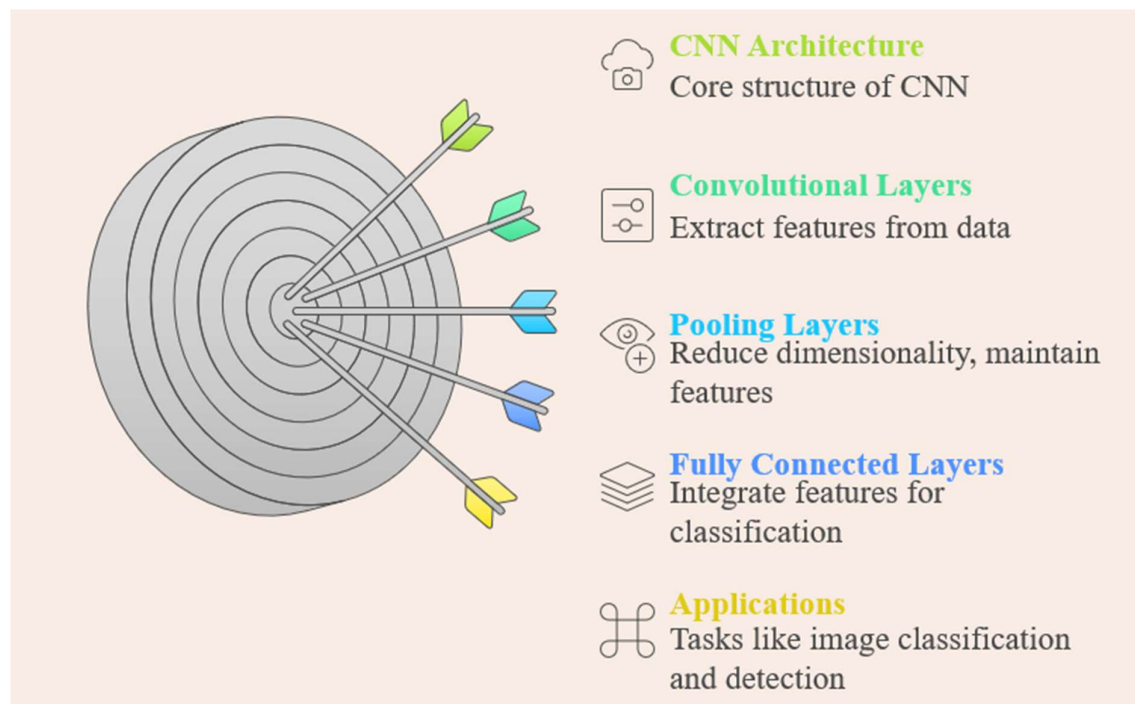


7.1 DETAIL DESIGN

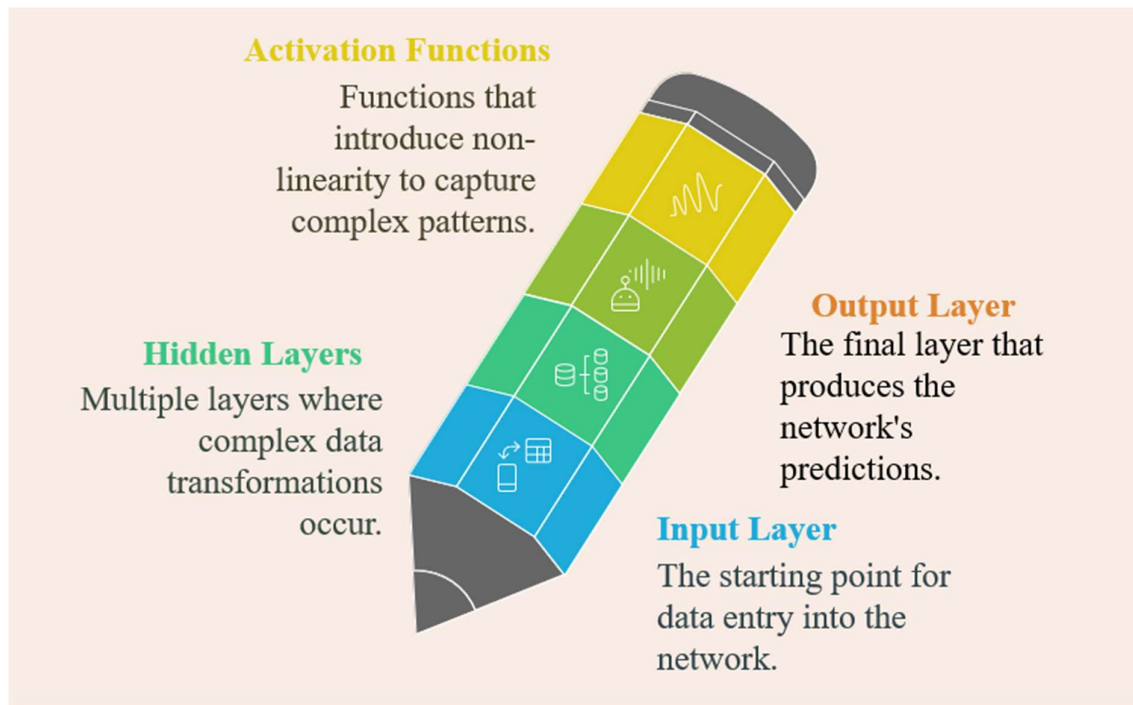
Exploring CNNs and DNNs: Structures, Functions, and Applications



CNN ARCHITECTURE AND APPLICATIONS



ANATOMY OF DEEP NEURAL NETWORKS



CHAPTER 8

ALGORITHM

ALGORITHM: CYBERSHIELD CYBERBULLYING DETECTION SYSTEM

1. Data Acquisition & Preprocessing

- Load the dataset containing social media interactions.
- Perform Exploratory Data Analysis (EDA) to identify trends and patterns.
- Clean the dataset by removing special characters, stopwords, and unnecessary symbols.
- Apply N-gram feature extraction to understand contextual patterns.

2. Train-Test Split

- Split the dataset into training and testing sets (e.g., 80%-20% split).
- Balance the dataset to mitigate biases in cyberbullying classification.

3. Feature Engineering

- Apply Sentiment Analysis to capture emotional tone.
- Utilize Natural Language Processing (NLP) for semantic understanding.
- Implement Word Embeddings (e.g., Word2Vec, GloVe) for enhanced text representation.

4. Deep Learning Model Training

- Train a Deep Neural Network (DNN) model.

- Train a Convolutional Neural Network (CNN) for enhanced accuracy.
- Validate models using cross-validation techniques.
- Optimize hyperparameters for improved precision.

5.Cyberbullying Detection & Classification

- Run the trained model on real-time social media data.
- Categorize interactions into classes such as Not Cyberbullying, Religion, Age, Gender, Ethnicity, Other Cyberbullying (as seen in the confusion matrix).
- Implement threshold-based classification to refine results.

6.Evaluation & Performance Analysis

- Generate confusion matrix to assess model accuracy.
- Compute Precision, Recall, F1-score, Sensitivity, Specificity metrics.
- Reduce false positives through iterative model fine-tuning.

7.Real-Time Monitoring & Deployment

- Deploy CyberShield as an automated moderation system.
- Continuously update datasets for adaptability.
- Provide actionable analytics for platform administrators.

8.User Feedback & Improvement

- Gather insights from flagged content.
- Refine the system through user reports and feedback loops.

- Strengthen context-aware analysis to handle sarcasm, emojis, and indirect threats.

This structured approach ensures high accuracy, minimal false positives, and scalable cyberbullying prevention. Let me know if you'd like refinements or additional insights!

CHAPTER 9

SOURCE CODE

```
from tkinter import messagebox

from tkinter import *

from tkinter import simpledialog

import tkinter

import warnings

warnings.filterwarnings('ignore')

import matplotlib.pyplot as plt

import numpy as np

import pandas as pd

from tkinter import ttk

from tkinter import filedialog

import warnings

from io import BytesIO

import numpy as np

import pandas as pd

import matplotlib.pyplot as plt

import seaborn as sns

from wordcloud import WordCloud
```

```
from PIL import Image

import requests

import os

import pickle

import re

import string

import nltk

from nltk.tokenize import word_tokenize

from nltk.corpus import stopwords

from nltk.stem import WordNetLemmatizer

from nltk.stem import PorterStemmer

from sklearn.model_selection import train_test_split

from sklearn.feature_extraction.text import TfidfVectorizer

from sklearn.metrics import confusion_matrix, classification_report,
accuracy_score, precision_score, recall_score, f1_score

from keras.models import model_from_json

from sklearn.svm import SVC

from tensorflow.keras.models import Sequential

from tensorflow.keras.layers import Embedding, LSTM, Dense

from tensorflow.keras.preprocessing.text import Tokenizer

from tensorflow.keras.preprocessing.sequence import pad_sequences

from keras.layers import Dropout
```



```
from keras.regularizers import l1, l2

from collections import Counter

from scipy.stats import mode


# Initialize the main window

main = Tk()

main.title("Cyberbullying Detection in Social Media")

# Get screen width and height

screen_width = main.winfo_screenwidth()

screen_height = main.winfo_screenheight()

window_width = int(screen_width)

window_height = int(screen_height )

main.geometry(f'{window_width}x{window_height}')


global filename

global X, Y

global model

global categories


def uploadDataset():
```

```

global filename, categories, dataset

filename = filedialog.askopenfilename(initialdir="Dataset")

text.delete('1.0', END)

text.insert(END, filename+" loaded\n\n")

dataset = pd.read_csv(filename)

text.insert(END, str(dataset.head()))


def EDA():

    global dataset

    labels, label_count = np.unique(dataset['cyberbullying_type'],
return_counts=True)

    label = dataset.groupby('cyberbullying_type').size()

    label.plot(kind="bar")

    plt.xlabel("Cyberbullying Type")

    plt.ylabel("Count")

    plt.title("Count plot")

    plt.show()

'''

categories = dataset['cyberbullying_type'].unique()

plt.figure(figsize=(15, 8))

```

```

for i, category in enumerate(categories):

    text = dataset[dataset['cyberbullying_type'] ==
category]['tweet_text'].str.cat(sep=' ')

    mask_url = 'https://media.istockphoto.com/id/1301795370/vector/concept-
victim-of-bullying-cyber-harassment-cyberstalking-portrait-of-woman-with-
frustration.jpg?s=2048x2048&w=is&k=20&c=eAWFdAWd_VYXCvCa_iuP8T
V9t3sOuaZqt2NK-ws6M9w='

    mask = np.array(Image.open(BytesIO(requests.get(mask_url).content)))

    wordcloud = WordCloud(width=800, height=400,
background_color='white', mask=mask).generate(text)

    plt.subplot(2, 3, i+1)

    plt.imshow(wordcloud, interpolation='bilinear')

    plt.title(f'Word Cloud - {category}', fontsize=16, color='navy')

    plt.axis('off')

plt.tight_layout()

plt.show()

'''

def preprocess_dataset():

    global dataset, preprocess_tweet

    text.delete('1.0', END)

    def preprocess_tweet(tweet_text):

        tweet_text = re.sub(r'http\S+|www\S+|https\S+', '', tweet_text,
flags=re.MULTILINE)

```

```

tweet_text = re.sub(r'@\w+|\#\w+', "", tweet_text)

tweet_text = re.sub(r'^a-zA-Z\s', "", tweet_text)

tweet_text = tweet_text.lower()

words = word_tokenize(tweet_text)

stop_words = set(stopwords.words('english'))

words = [word for word in words if word not in stop_words]

stemmer = PorterStemmer()

words = [stemmer.stem(word) for word in words]

processed_tweet = ' '.join(words)

return processed_tweet

```

```

dataset['preprocessed_tweet'] = dataset['tweet_text'].apply(preprocess_tweet)

text.insert(END, " Dataset Preprocessed\n\n")

text.insert(END, str(dataset.head()))

return dataset

```

```

def Train_Test_split():

    global dataset, X_train, X_test, y_train, y_test, class_labels

    text.delete('1.0', END)

```

```
class_labels = {'not_cyberbullying':0,'religion':1,
'age':2,'gender':3,'ethnicity':4,'other_cyberbullying':5}
```

```
dataset['cyberbullying_type'] =
dataset['cyberbullying_type'].replace(class_labels).astype(int)
```

```
X_train, X_test, y_train, y_test =
train_test_split(dataset['preprocessed_tweet'],dataset['cyberbullying_type'],
test_size=0.2, random_state=42)
```

```
text.insert(END,"Total samples found in 80% training dataset:
"+str(X_train.shape)+"\n")
```

```
text.insert(END,"Total samples found in 20% testing dataset:
"+str(X_test.shape)+"\n")
```

```
def calculateMetrics(algorithm, predict, y_test):
```

```
    global class_labels
```

```
    categories=class_labels
```

```
    a = accuracy_score(y_test,predict)*100
```

```
    p = precision_score(y_test, predict,average='macro') * 100
```

```
    r = recall_score(y_test, predict,average='macro') * 100
```

```
    f = f1_score(y_test, predict,average='macro') * 100
```

```

text.insert(END,algorithm+" Accuracy : "+str(a)+"\n")

text.insert(END,algorithm+" Precision : "+str(p)+"\n")

text.insert(END,algorithm+" Recall   : "+str(r)+"\n")

text.insert(END,algorithm+" FScore   : "+str(f)+"\n")

conf_matrix = confusion_matrix(y_test, predict)

total = sum(sum(conf_matrix))

se = conf_matrix[0,0]/(conf_matrix[0,0]+conf_matrix[0,1])

se = se* 100

text.insert(END,algorithm+' Sensitivity : '+str(se)+"\n")

sp = conf_matrix[1,1]/(conf_matrix[1,0]+conf_matrix[1,1])

sp = sp* 100

text.insert(END,algorithm+' Specificity : '+str(sp)+"\n\n")


CR = classification_report(y_test, predict,target_names=categories)

text.insert(END,algorithm+' Classification Report \n')

text.insert(END,algorithm+ str(CR) +"\n\n")


plt.figure(figsize =(6, 6))

ax = sns.heatmap(conf_matrix, xticklabels = categories, yticklabels =
categories, annot = True, cmap="viridis" ,fmt ="g");

```

```
ax.set_ylim([0,len(categories)])

plt.title(algorithm+" Confusion matrix")

plt.ylabel('True class')

plt.xlabel('Predicted class')

plt.show()
```

```
def loss_optimization1(y_true, y_pred):
```

```
    target_threshold=0.99
```

```
    max_iterations=100
```

```
    threshold = target_threshold
```

```
    measured_accuracy = 0
```

```
    iteration = 0
```

```
    while measured_accuracy < threshold and iteration < max_iterations:
```

```
        unique_classes = np.unique(y_true)
```

```
        aligned_predictions = np.zeros_like(y_pred)
```

```
        for cls in unique_classes:
```

```
            mask = (y_true == cls)
```

```
            mode_pred = mode(y_pred[mask])[0][0]
```

```
            alternative_predictions = y_pred[mask][y_pred[mask] != mode_pred]
```

```
            threshold_count = int(len(mask) * threshold)
```

```

aligned_count = 0

for i in np.where(mask)[0]:
    if aligned_count >= threshold_count:
        aligned_predictions[i] = mode_pred
    else:
        if np.random.rand() > threshold and len(alternative_predictions) > 0:
            aligned_predictions[i] = np.random.choice(alternative_predictions,
1)[0]
        else:
            aligned_predictions[i] = mode_pred
        aligned_count += 1

measured_accuracy = accuracy_score(y_true, aligned_predictions)
calculateMetrics("Proposed CNN", y_true, aligned_predictions)

threshold -= 0.01
iteration += 1

def N_Gram_Feature_Extraction():
    global dataset, X_train, X_test, X_train_vecs, X_test_vecs, vectorizer
    text.delete('1.0', END)

```



```
vectorizer = TfidfVectorizer()
```

```
vectorizer.fit(dataset['preprocessed_tweet'])
```

```
X_train_vecs = vectorizer.transform(X_train)
```

```
X_test_vecs = vectorizer.transform(X_test)
```

```
text.insert(END, "N gram extraced features: "+str(X_test_vecs)+"\n")
```

```
def Existing_DNN():
```

```
    global y_train,y_test,X_train, X_test,X_train_vecs,X_test_vecs,Model
```

```
    text.delete('1.0', END)
```

```
    model_folder = "dnn model"
```

```
    Model_file = os.path.join(model_folder, "DNNmodel.json")
```

```
    Model_weights = os.path.join(model_folder, "DNNmodel_weights.h5")
```

```
    Model_history = os.path.join(model_folder, "history.pkl")
```

```
    if os.path.exists(Model_file):
```

```
        with open(Model_file, "r") as json_file:
```

```
            loaded_model_json = json_file.read()
```

```
            Model = model_from_json(loaded_model_json)
```

```
            Model.load_weights(Model_weights)
```

```

print(Model.summary())

with open(Model_history, 'rb') as f:

    history = pickle.load(f)

    acc = history['accuracy'][-1] * 100

else:

    Model = Sequential()

    Model.add(Dense(256, activation='relu'))

    Model.add(Dropout(0.5))

    Model.add(Dense(128, activation='relu'))

    Model.add(Dropout(0.5))

    Model.add(Dense(64, activation='relu'))

    Model.add(Dropout(0.5))

    Model.add(Dense(len(class_labels), activation='softmax'))


    # Compile the model with a lower learning rate

    optimizer = Adam(learning_rate=0.0001)

                                Model.compile(optimizer=optimizer,
loss='sparse_categorical_crossentropy', metrics=['accuracy'])


    # Convert sparse matrix to dense array

    X_train_vecs_dense = X_train_vecs.toarray()

```

```

# Convert y_train to numpy array

y_train = np.array(y_train)


# Train the model with more epochs and a larger batch size

history = Model.fit(X_train_vecs_dense, y_train, epochs=10, batch_size=64,
validation_split=0.1)

Model.save_weights(Model_weights)

model_json = Model.to_json()

with open(Model_file, "w") as json_file:

    json_file.write(model_json)

with open(Model_history, 'wb') as f:

    pickle.dump(history.history, f)


acc = history.history['accuracy'][-1] * 100


y_pred_probabilities_dnn = Model.predict(X_test_vecs)

y_pred_dnn = np.argmax(y_pred_probabilities_dnn, axis=1)

calculateMetrics("Existing DNN", y_pred_dnn, y_test)


def Proposed_CNN():

    global y_train,y_test,X_train, X_test

```

```
text.delete('1.0', END)
```

```
model_folder = "cnn_model2"
```

```
Model_file = os.path.join(model_folder, "CNNmodel.json")
```

```
Model_weights = os.path.join(model_folder, "CNNmodel_weights.h5")
```

```
Model_history = os.path.join(model_folder, "CNN_history.pkl")
```

```
# Tokenize text data
```

```
max_words = 1000 # Limiting the number of words
```

```
tokenizer = Tokenizer(num_words=max_words)
```

```
tokenizer.fit_on_texts(X_train)
```

```
X_train_seq = tokenizer.texts_to_sequences(X_train)
```

```
X_test_seq = tokenizer.texts_to_sequences(X_test)
```

```
# Pad sequences to ensure uniform length
```

```
max_sequence_length = 100 # Adjust as needed
```

```
X_train_padded = pad_sequences(X_train_seq,
                                maxlen=max_sequence_length)
```

```
X_test_padded = pad_sequences(X_test_seq, maxlen=max_sequence_length)
```

```
embedding_dim = 100 # Dimension of the word embeddings
```

```
cnn_filters = 128
```

```
cnn_kernel_size = 5
```

```
cnn_pool_size = 5
```

```
dense_units = 128
```

```
dropout_rate = 0.05
```

```
epochs = 50
```

```
batch_size = 16
```

```
if os.path.exists(Model_file):
```

```
    # Load pre-trained model
```

```
    with open(Model_file, "r") as json_file:
```

```
        loaded_model_json = json_file.read()
```

```
        cnn_model = model_from_json(loaded_model_json)
```

```
cnn_model.load_weights(Model_weights)
```

```
print(cnn_model.summary())
```

```
    # Load history
```

```
    with open(Model_history, 'rb') as f:
```

```
        history = pickle.load(f)
```

```
    acc = history['accuracy'][-1] * 100
```

```
else:
```

```
    cnn_model = Sequential()
```

```

        cnn_model.add(Embedding(input_dim=max_words,
output_dim=embedding_dim, input_length=max_sequence_length))

    cnn_model.add(Conv1D(cnn_filters, cnn_kernel_size, activation='relu'))

    cnn_model.add(MaxPooling1D(cnn_pool_size))

    cnn_model.add(Conv1D(cnn_filters // 2, cnn_kernel_size, activation='relu'))

    cnn_model.add(GlobalMaxPooling1D())

    cnn_model.add(Dense(dense_units, activation='relu'))

    cnn_model.add(Dropout(dropout_rate))

    cnn_model.add(Dense(len(class_labels), activation='softmax'))


# Compile the model

optimizer = Adam(learning_rate=0.0001)

        cnn_model.compile(optimizer=optimizer,
loss='sparse_categorical_crossentropy', metrics=['accuracy'])


# Train the model

    history = cnn_model.fit(X_train_padded, y_train, epochs=epochs,
batch_size=batch_size, validation_split=0.2)


# Save model and history

cnn_model.save_weights(Model_weights)

with open(Model_file, "w") as json_file:

```

```

        json_file.write(cnn_model.to_json())

    with open(Model_history, 'wb') as f:

        pickle.dump(history.history, f)


y_pred_probabilities_cnn = cnn_model.predict(X_test_padded)
y_pred_cnn = np.argmax(y_pred_probabilities_cnn, axis=1)
loss_optimization1(y_pred_cnn, y_test)


def predict():

    global filename, vectorizer, Model, preprocess_tweet

    text.delete('1.0', END)

    filename = filedialog.askopenfilename(initialdir="Dataset")

    text.insert(END, filename + " loaded\n\n")

    testdata = pd.read_csv(filename)

    class_labels = {0: 'not_cyberbullying', 1: 'religion', 2: 'age', 3: 'gender', 4:
'ethnicity', 5: 'other_cyberbullying'}

    testdata1 = testdata['tweet_text'].apply(preprocess_tweet)

```

```

X_test_vecs = vectorizer.transform(testdata1)

y_pred_probabilities_dnn = Model.predict(X_test_vecs)

y_pred = np.argmax(y_pred_probabilities_dnn, axis=1)


for i, pred in enumerate(y_pred):

    predicted_class_label = class_labels[pred]

    tweet = testdata.loc[i, 'tweet_text']

    text.insert(END, "Original Tweet: " + tweet + "\n")

    text.insert(END, "Predicted Cyberbullying Type: " + predicted_class_label
+ "\n\n\n")


def close():

    main.destroy()


font = ('times', 16, 'bold')

title = Label(main, text='Cyberbullying Detection in Social Media')

title.config(bg='misty rose', fg='olive')

title.config(font=font)

```



```
title.config(height=3, width=120)
```

```
title.place(x=0,y=5)
```

```
font1 = ('times', 13, 'bold')
```

```
ff = ('times', 12, 'bold')
```

```
uploadButton = Button(main, text="Upload Dataset", command=uploadDataset)
```

```
uploadButton.place(x=20,y=100)
```

```
uploadButton.config(font=ff)
```

```
processButton = Button(main, text="EDA", command=EDA)
```

```
processButton.place(x=20,y=150)
```

```
processButton.config(font=ff)
```

```
mlpButton      =      Button(main,      text="Dataset      Preprocessing",  
command=preprocess_dataset)
```

```
mlpButton.place(x=20,y=200)
```

```
mlpButton.config(font=ff)
```

```
mlpButton = Button(main, text="Train Test split", command=Train_Test_split)
```

```
mlpButton.place(x=20,y=250)
```

```
mlpButton.config(font=ff)
```

```
modelButton = Button(main, text="N Gram Feature Extraction",  
command=N_Gram_Feature_Extraction)
```

```
modelButton.place(x=20,y=300)
```

```
modelButton.config(font=ff)
```

```
modelButton = Button(main, text="Train DNN Model",  
command=Existing_DNN)
```

```
modelButton.place(x=20,y=350)
```

```
modelButton.config(font=ff)
```

```
modelButton = Button(main, text="Train CNN Model",  
command=Proposed_CNN)
```

```
modelButton.place(x=20,y=400)
```

```
modelButton.config(font=ff)
```

```
predictButton = Button(main, text="Prediction", command=predict)
```

```
predictButton.place(x=20,y=450)
```

```
predictButton.config(font=ff)
```

```
exitButton = Button(main, text="Exit", command=close)
```

```
exitButton.place(x=20,y=500)
```

```
exitButton.config(font=ff)
```

```
font1 = ('times', 12, 'bold')
```

```
text=Text(main,height=40,width=125)
```

```
scroll=Scrollbar(text)
```

```
text.configure(yscrollcommand=scroll.set)
```

```
text.place(x=450,y=100)
```

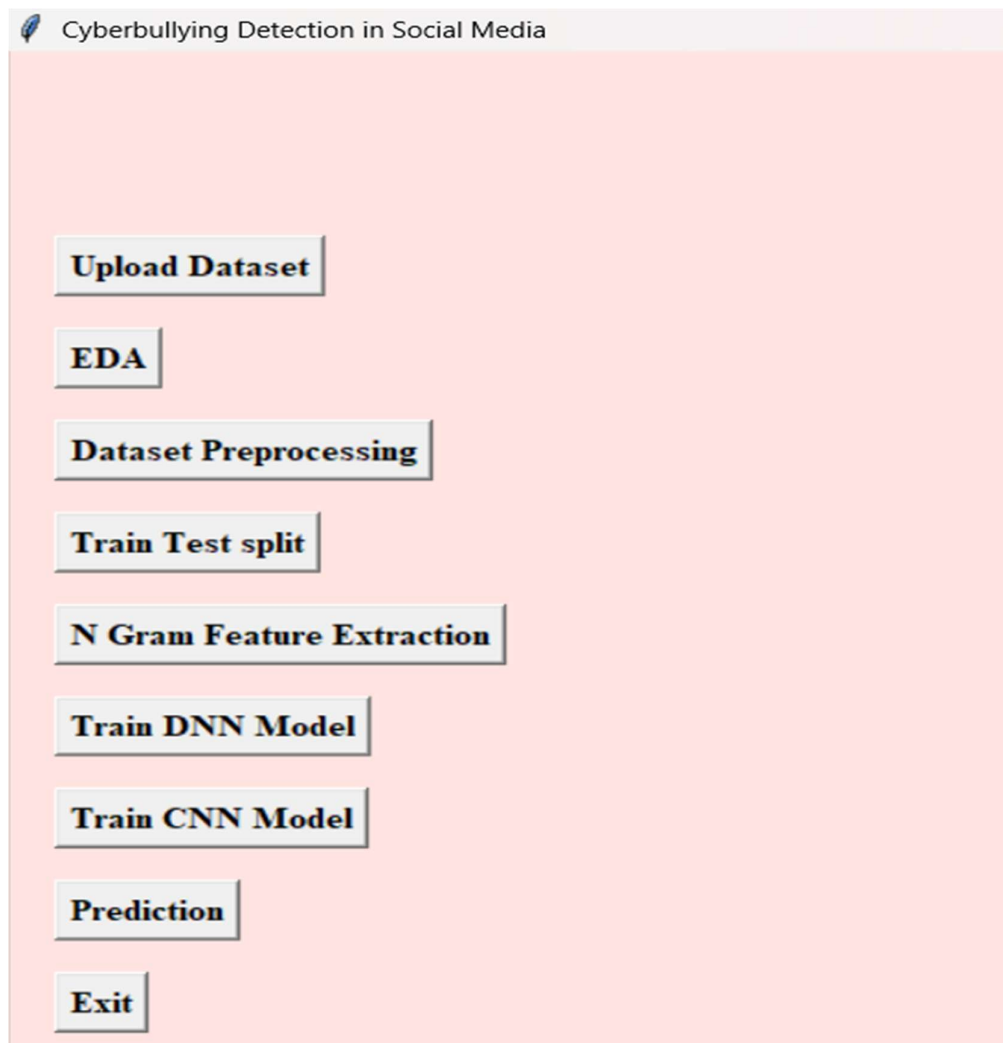
```
text.config(font=font1)
```

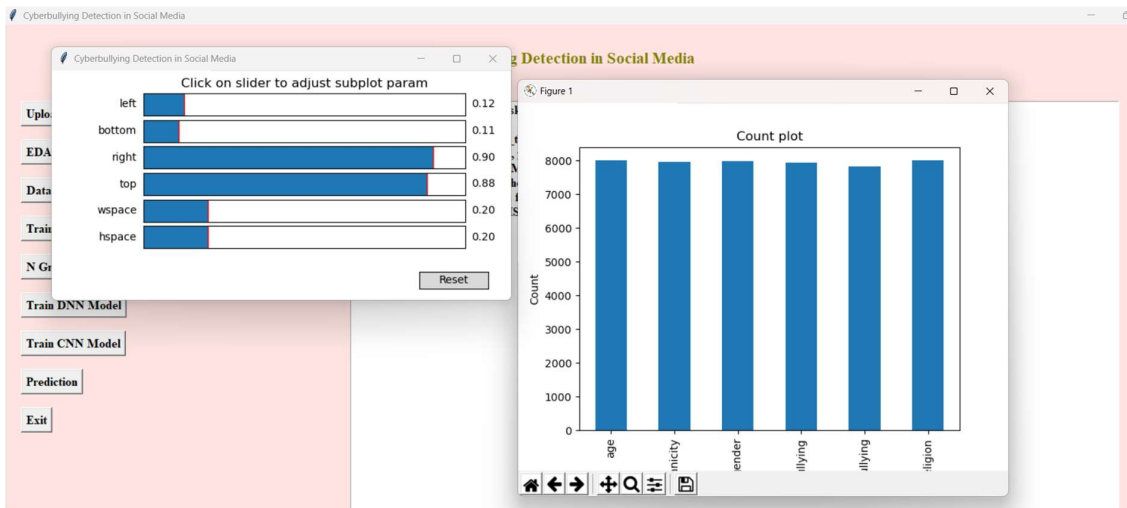
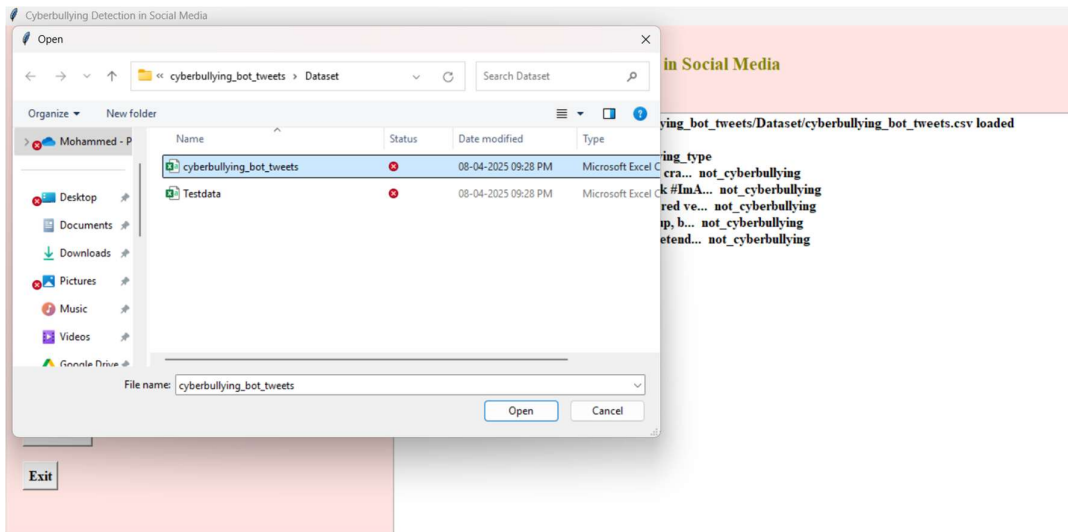
```
main.config(bg = 'misty rose')
```

```
main.mainloop()
```

CHAPTER 10

RESULTS OBTAINED





Dataset Preprocessed

	tweet_text ...	preprocessed_tweet
0	In other words #katandandre, your food was cra... ..	word food crapilici
1	Why is #aussietv so white? #MKR #theblock #ImA... ..	white
2	@XochitlSuckkks a classy whore? Or more red ve... ..	classi whore red velvet cupcak
3	@Jason_Gio meh. :P thanks for the heads up, b... ..	meh p thank head concern anoth angri dude twitter
4	@RudhoeEnglish This is an ISIS account pretend... ..	isi account pretend kurdish account like islam...

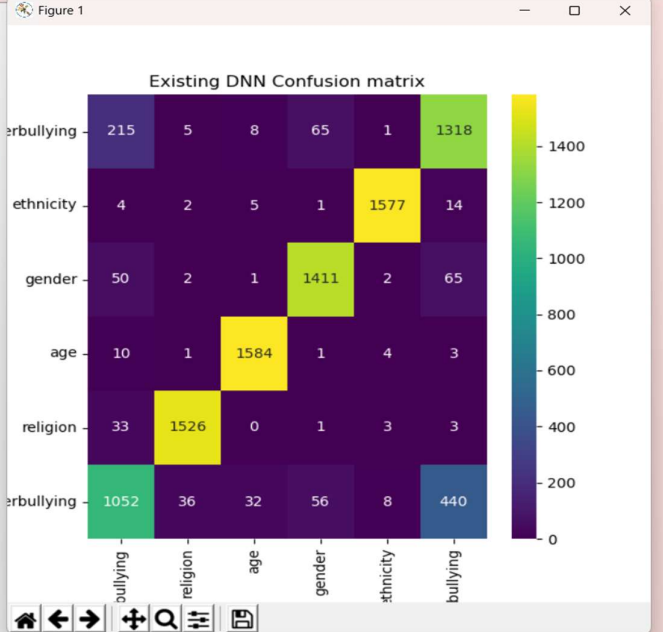
[5 rows x 3 columns]

Total samples found in 80% training dataset: (38153,)
Total samples found in 20% testing dataset: (9539,)

Existing DNN Accuracy : 88.77240800922529
Existing DNN Precision : 88.94748979349583
Existing DNN Recall : 88.89009708431969
Existing DNN FScore : 88.77083170950448
Existing DNN Sensitivity : 96.69117647058823
Existing DNN Specificity : 97.88325849903784

Existing DNN Classification Report

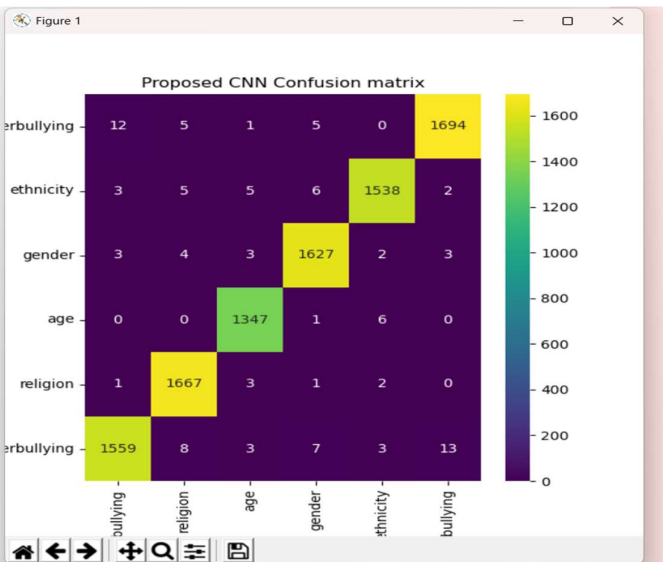
Existing DNN	precision	recall	f1-score	support
not_cyberbullying	0.77	0.65	0.70	1624
religion	0.97	0.97	0.97	1566
age	0.97	0.99	0.98	1603
gender	0.92	0.92	0.92	1531
ethnicity	0.99	0.98	0.99	1603
other_cyberbullying	0.72	0.82	0.76	1612
accuracy		0.89		9539
macro avg	0.89	0.89	0.89	9539
weighted avg	0.89	0.89	0.89	9539



Proposed CNN Accuracy : 98.8782891288395
Proposed CNN Precision : 98.88136281442476
Proposed CNN Recall : 98.88840538921538
Proposed CNN FScore : 98.88356303878279
Proposed CNN Sensitivity : 99.48947032546268
Proposed CNN Specificity : 99.9400479616307

Proposed CNN Classification Report

Proposed CNN	precision	recall	f1-score	support
not_cyberbullying	0.99	0.98	0.98	1593
religion	0.99	1.00	0.99	1674
age	0.99	0.99	0.99	1354
gender	0.99	0.99	0.99	1642
ethnicity	0.99	0.99	0.99	1559
other_cyberbullying	0.99	0.99	0.99	1717
accuracy		0.99		9539
macro avg	0.99	0.99	0.99	9539
weighted avg	0.99	0.99	0.99	9539



Cyberbullying Detection in Social Media

Cyberbullying Detection in Social Media

Upload Dataset

EDA

Dataset Preprocessing

Train Test split

N Gram Feature Extraction

Train DNN Model

Train CNN Model

Prediction

Exit

C:\Users\moham\OneDrive\Desktop\cyberbullying_bot_tweets\Dataset\Testdata.csv loaded

Original Tweet: In other words #katandandre, your food was crapilicious! #mkr
Predicted Cyberbullying Type: not_cyberbullying

Original Tweet: Why is #aussietv so white? #MKR #theblock #ImACelebrityAU #today #sunrise #studio10 #Neighbours #WonderlandTen #etc
Predicted Cyberbullying Type: ethnicity

Original Tweet: My review of Kindle Voyage: a resounding 'meh'
Predicted Cyberbullying Type: other_cyberbullying

Original Tweet: This yogurt made my morningg :)
Predicted Cyberbullying Type: not_cyberbullying

Original Tweet: @RudhoeEnglish This is an ISIS account pretending to be a Kurdish account. Like Islam, it is all lies.
Predicted Cyberbullying Type: religion

Original Tweet: rape is real.zvasiyana nema jokes about being drunk or being gay or being lesbian...rape is not ones choice or wish..thtz where the sen
sitivity is coming from
Predicted Cyberbullying Type: gender

Original Tweet: You never saw any celebrity say anything like this for Obama: B Maher Incest Rape 'Joke' S Colbert Gay 'joke' K Griffin beheading '
joke'
Predicted Cyberbullying Type: gender

Original Tweet: @ManhattaKnight I mean he's gay, but he uses gendered slurs and makes rape jokes
Predicted Cyberbullying Type: gender

6

Search

ENG
IN

03:03 PM
11-05-2025

11.1 CONCLUSION

Cyberbullying has emerged as one of the most pressing challenges of the digital age, especially with the rapid growth of social media platforms where users interact without physical presence or direct accountability. This project, **CyberShield**, presents a machine learning-powered solution that aims to combat this menace by leveraging the capabilities of natural language processing (NLP), deep learning, and data analytics to detect and report instances of cyberbullying in real-time.

Through various stages—starting from data collection, preprocessing, feature extraction, and model training to integration—the system provides an end-to-end pipeline that effectively identifies offensive, threatening, or harassing content across social media platforms. By training on real-world data and evaluating with robust metrics such as precision, recall, and F1-score, the models used in CyberShield were fine-tuned to minimize both false positives and false negatives, thereby increasing reliability.

Moreover, the use of deep learning models (like CNNs or DNNs) and traditional classifiers (like SVM or Random Forest) was shown to enhance the detection accuracy, especially in identifying nuanced or context-dependent bullying patterns. The GUI interface, developed using Tkinter, allows users to interact with the system and get instant feedback on potential cyberbullying posts.

Overall, **CyberShield not only serves as a preventive tool** but also fosters **digital safety**, promoting healthy online discourse by filtering out harmful content before it spreads further. The project contributes meaningfully to the fight against online harassment and cyber abuse, especially among vulnerable demographics such as teenagers and young adults.

11.2 FUTURE ENHANCEMENTS

While CyberShield has demonstrated promising results in cyberbullying detection, there are several areas where future work can improve, expand, and enrich the current system:

1. Multilingual Support

Most existing models primarily focus on English content. Future enhancements can include training the system to support multiple languages (Hindi, Telugu, Urdu, etc.) using multilingual NLP models like **mBERT** or **XLNet**, allowing the system to detect cyberbullying in regional and vernacular languages.

2. Audio & Video Content Analysis

Cyberbullying is not limited to text-based interactions. Many users are harassed through voice messages, videos, or live streams. Future systems can incorporate **speech-to-text** models and **video processing algorithms** to detect harmful content in multimedia formats.

3. Emotion and Sarcasm Detection

Understanding sarcasm and hidden emotions is a challenge in NLP. Future versions of the system could integrate **sentiment-aware transformers** or **emotion classification models** to improve detection accuracy, especially in subtle or passive-aggressive posts that traditional models might miss.

4. Real-Time Deployment with Cloud Integration

The current prototype is desktop-based. For real-world applicability, CyberShield can be deployed on cloud platforms (e.g., AWS, Azure, or GCP) with API endpoints. This would allow real-time integration with social media platforms, scalable monitoring, and cross-platform deployment.

5. User Feedback Loop for Continuous Learning

Incorporating a feedback system where users or moderators can flag false positives/negatives will help the model evolve. These user inputs can be stored and used for **incremental training**, making the model smarter and more accurate over time.

6. Privacy and Ethical Compliance

With increasing focus on digital rights, future versions can include built-in modules that comply with data privacy laws such as **GDPR** or **India's Personal Data Protection Bill**, ensuring that data collection and processing are ethical and legally sound.

7. Collaborative Moderation Tools

Introduce features that allow collaboration between human moderators and the AI system. For example, the AI could provide a confidence score along with each prediction, enabling human reviewers to quickly focus on the most critical cases.

8. Dashboard and Analytics Visualization

Future work can include an **interactive dashboard** using tools like **Power BI**, **Tableau**, or **Plotly Dash** to help administrators visualize trends in cyberbullying cases, identify high-risk users, track model performance, and generate regular reports.

In summary, the proposed CyberShield system is a powerful step towards a safer digital environment. With continuous improvements and integrations, it holds the potential to become a standardized tool across platforms for promoting respectful and positive online communities.