

Deep Learning Approaches for Bangla Facebook Comment Classification Using Bangla Bert, GRU, LSTM, and CNN

BY

Nowfal Ahamed Shawon

ID: 232-25-031

This Report Presented in Partial Fulfillment of the Requirements for
The Degree of Masters of Science in Computer Science and Engineering

Supervised By

Dr. Naznin Sultana

Associate Professor

Department of CSE

Daffodil International University

Co-Supervised By

Abdus Sattar

Assistant Professor

Department of CSE

Daffodil International University



DAFFODIL INTERNATIONAL UNIVERSITY

DHAKA, BANGLADESH

DECEMBER 2024

APPROVAL

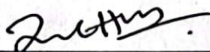
This Project titled "Deep Learning Approaches for Bangla Facebook Comment Classification Using Bangla Bert, GRU, LSTM, and CNN", submitted by **Nowfal Ahamed Shawon**, ID No: 232-25-031 to the Department of Computer Science and Engineering, Daffodil International University has been accepted as satisfactory for the partial fulfillment of the requirements for the degree of M.Sc. in Computer Science and Engineering and approved as to its style and contents. The presentation has been held on.

BOARD OF EXAMINERS



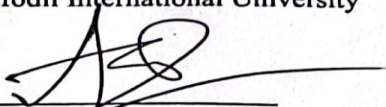
Dr. Sheak Rashed Haider Noori, PhD
Professor and Head
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Chairman



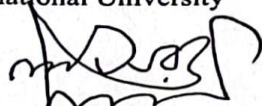
Dr. Md. Zahid Hasan, PhD
Associate Professor
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner



Dr. Arif Mahmud, PhD
Associate Professor & Director MIS
Department of Computer Science and Engineering
Faculty of Science & Information Technology
Daffodil International University

Internal Examiner



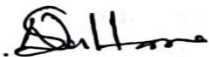
Dr. Mohammed Nasir Uddin, PhD
Professor
Department of Computer Science and Engineering
Jagannath University

External Examiner

DECLARATION

I hereby declare that this research has been done by me under the supervision of Naznin Sultana, Associate Professor, Department of CSE, Daffodil International University. I also declare that neither this project nor any part of this project has been submitted elsewhere for award of any degree or diploma.

Supervised by:



Dr. Naznin Sultana
Associate Professor
Department of CSE
Daffodil International University

Co-Supervised by:



Abdus Sattar
Assistant Professor
Department of CSE
Daffodil International University

Submitted by:



Nowfal Ahamed Shawon
ID: 232-25-031
Department of CSE
Daffodil International University

ACKNOWLEDGEMENT

First, I express my heartiest thanks and gratefulness to Almighty Allah for His divine blessing which makes it possible to complete the final year project/internship successfully.

I am really grateful and wish my profound indebtedness to **Naznin Sultana**, Associate Professor, Department of CSE, Daffodil International University, Dhaka, deep knowledge & keen interest of my supervisor in the field of Deep Learning to carry out this project. His endless patience, scholarly guidance, continual encouragement, constant and energetic supervision, constructive criticism, valuable advice, reading many inferior drafts and correcting them at all stages have made it possible to complete this project.

I would like to express my heartiest gratitude to **Dr. Sheak Rashed Haider Noori, Head**, Department of CSE, for his kind help to finish our project and also to other faculty members and the staff of CSE department of Daffodil International University.

Finally, I must acknowledge with due respect the constant support and patients of my parents.

ABSTRACT

This has driven artificial intelligence (AI) development to grow in rapidly, and large amounts of language resources are being developed for an array of languages. But Bangla has a long way to go in terms of the contributing field of AI. Different categories of Bangla Facebook comments, such as Not Bully, Troll, Sexual and Religious are reviewed in this study specifically for the Bangla language resources. Using more than twenty-five thousand comments, we experimented and optimized various models such as Bangla BERT, GRU, LSTM and CNN. In our experimental results, the best performing Bangla BERT model reached an accuracy of 80% on the test dataset, whereas GRU achieved 70%, LSTM with 65% and finally CNN achieved just around 66%. We noticed ingrained biases in the dataset too. This can be useful for Bangla AI Resources which can be further utilized in sentiment analysis and content moderation systems to aid Bangla Speakers both domestically as well as globally.

TABLE OF CONTENTS

CONTENTS	PAGE
Board of examiners	ii
Declaration	iii
Acknowledgements	iv
Abstract	v
CHAPTER	
CHAPTER 1: INTRODUCTION	1-5
1.1 Introduction	1-2
1.2 Motivation	2-3
1.3 Rationale of the Study	3-4
1.4 Research Questions	4
1.5 Expected Output	4-5
1.6 Project Management and Finance	5
1.7 Report Layout	5
CHAPTER 2: BACKGROUND	6-8
2.1 Preliminaries/Terminologies	6
2.2 Related works	6-8
2.3 The Problem's Scope	8
2.4 Challenges	8
CHAPTER 3: RESEARCH METHODOLOGY	9-19
3.1 Proposed Methodology/Applied Mechanism	9
3.2 Data Collection Procedure/Dataset Utilized	10
3.2.1 Cleaning the dataset	10-11
3.2.2 Stopword Removal	11
3.2.3 Removal of Low-Length Data	12
3.4 Deep Learning Model	13
3.4.1 Banglabert	13
3.4.1 GRU	14
3.4.1 LSTM	14-15
3.4.1 CNN	15

3.5 Training Model	16
CHAPTER 4: EXPERIMENTAL RESULTS AND DISCUSSION	17-27
4.1 Results of Data Preprocessing	17
4.2 Confusion matrix	18
4.3 Classification Report	18
4.4 Deep Learning Model	19
4.4.1 Banglabert	19-20
4.4.2 GRU	20-21
4.4.3 LSTM	22-23
4.4.4 CNN	23-24
4.5 Experimental Result & Analysis	25-26
CHAPTER 5: CONCLUSION AND FUTURE WORK	27-29
5.1 Conclusion	27
5.2 Future Work	28
REFERENCES	29-31

LIST OF FIGURES

FIGURES	PAGE NO
Fig 3.1: WorkFlow	9
Fig. 3.2: Dataset Utilized	10
Fig. 3.3: Cleaning dataset	11
Fig. 3.4: Stop Word remove dataset	11
Fig. 3.5: Removal of Low-Length Data	12
Fig. 3.6: Final Dataset	12
Fig. 3.7: Banglabert	13
Fig. 3.8: GRU	14
Fig. 3.9 LSTM	15
Fig. 3.1.: CNN	15
Fig 4.1.1: Preprocess Data	17
Fig 4.4.1.1: Confusion Matrix of Banglaert	19
Fig 4.4.1.2: Classification report of Banglabert	20
Fig 4.4.2.1: Confusion Matrix of GRU	21
Fig 4.4.1.2: Classification report of GRU	21
Fig 4.4.3.1: Confusion Matrix of LSTM	22
Fig 4.4.3.2: Classification report of LSTM	23
Fig 4.4.4.1: Confusion Matrix of CNN	24
Fig 4.4.4.2: Classification report of CNN	24

LIST OF TABLES

TABLES	PAGE NO
Table 4.2.1: Confusion Matrix	18
Table 4.5.1: Model Accuracy	25

CHAPTER 1

INTRODUCTION

1.1 Introduction

Artificial Intelligence (AI) has made tremendous advances over the recent years where Natural Language Processing (NLP) is one of its most vital research field. Yet, Bangla (Bengali) is the fifth most spoken language in the world and still has a little less focus in terms of AI resources and tools. To the best of our knowledge, this work is the first study on Bangla Facebook comment classification which is to classify comments into four classes: not bully, troll, religious and sexual. Due to unavailability of proper annotated datasets for Bangla, development of Text classification and Sentiment analysis techniques stuck behind and Bangla NLP field lagged a lot compared to the high resource languages.

Bangla is a language with a long cultural history. With 228 million native speakers worldwide, and a further 37 million speaking it as a second language,[9] it is the first or second most widely spoken language in the world. Although Bangla is a widely spoken language across the globe, it does not have large datasets and computational models that are required for any NLP related research or developing applications. Text classification, essential for organizing and analyzing textual data, has numerous practical applications from search engines and content management systems to news portals. But since we mostly learned for English, it did well but i got stuck in my own language Bangla because of low resource and less data available.

The Bangla language is a symbol of resistance and identity for the indigenous Bengali people. One hundred and two years after its official establishment, the parallel struggle in Bangladesh to reclaim their language took a violent turn on February 21, 1952 as people laid down their lives for the mother tongue; today that sacrifice is etched in global history. Bangla is loved in its own ground but taking place globally too [10]. Bangla, for example, gained increased global currency in 2020 when the Korean Central News Agency (KCNA) noted a hot trend of learning this language as a foreign language [11].

However, there is a bit of work has been done on Bangla comment classification [7] but we found no work that has focused on Bangla Facebook comment classification. We applied and compared some deep learning models, such as Bangla BERT (Bangla Bidirectional Encoder Representations from Transformers), GRU, LSTM and CNN using a dataset containing more than 25k comments. The data contains four classes: not bully meaning that offensive language or harassment is not found in the comments, troll which a comment to provoke and annoy others, religious related faith themes and sexual which inappropriate contents. Automated classification of each of these categories presents its own challenges.

Challenges in this study were also complex. The dataset was also highly imbalanced among the categories which impacted the model performance on initial models. Furthermore, the raw data presented a lot of noise such as duplicate or low-value entries that made classification efforts more complicated. We performed various preprocessing techniques to solve the above problems like stop word removal, dataset imbalance handler and removing comments with a short-length response. The dataset was further iteratively cleaned, analyzed and supplemented to effectively serve each our models.

Overall, although this research has several limitations, it will add value to Bangla AI development as we have shown how to apply and compare some deep learning models in Bangla Facebook comment classification. Although Bangla BERT was Type the most accurate, scoring 80%, other models e.g., GRU, LSTM and CNN gave us insights about advantages and disadvantages of using different architectures in this specific task. The contributions made by this study can be used as further foundation for Bangla NLP and encourages new research by alleviating the problems of dataset bias and resource scarcity in both sentiment analysis and text classification.

1.2 Motivation

Despite being one of the top spoken languages across the globe, there still remains a shortage of resources with regards to Bangla in Natural Language Processing (NLP),
©Daffodil International University

hence this research is motivated. Even though English and Chinese have huge datasets and tools, Bangla still lacks behind in many AI applications due to the unavailability of datasets. In this study, we try to fill that gap by proposing some strong Bangla Facebook comment classifier where popular Bangla FB comments are classified into categories including not bully, troll, religious and sexual. As social media influence increases, automation for comments in Bangla can contribute to making the online world safer.

The applications for content moderation and organization spanning from e-commerce to journalism, sentiment analysis and text classification are functions that capitalize on the power of NLP. But Bangla is still behind in this category because there are no annotated datasets and models built specifically for Bangla. To facilitate this research, we exploit some of the recently developed state-of-the-art deep learning models including Bangla BERT, GRU, LSTM and CNN focusing on their capability in dealing with a morphologically rich as well as complex language like Bangla.

The second major impetus is related to the difficulty of Bangla datasets namely, bias, imbalance and noise. This study not only mitigates these issues using techniques such as stop word removal, balancing the dataset and cleaning of text but also makes a pathway through which this subject can be researched further in Bangla using NLP.

Lastly, this research has an importance for me and my culture. Reminiscent of the 1952 Language Movement, Bangla boasts a glorious past which to preserve in this age of digital AI, is more than just a technical problem but a cultural endeavor. The study makes a step forward in developing tools that foster an accessible artificial intelligence by enhancing Bangla NLP and empowering the community of many Bengali speakers around the world.

1.3 Rationale of the Study

Compared to other languages, the development of Bangla NLP is still in its infancy and we do not find many annotated datasets/tools developed for Bangla like sentiment analysis or text classification. A lack of resources for Bangla NLP limits the ability to create low-cost AI applications, so relatively few Bangla speakers gain access to large-scale

©Daffodil International University 3

technologies that depend upon it, e.g., content moderation and sentiment analysis. This research enhances available Bangla NLP resources by building Bangla Facebook comment classification models and fine-tuning Bangla BERT, GRU, LSTM, and CNN for the task. These models will be helpful to social media platforms for automated content moderation, and researchers and developers can access improved datasets and methodologies. Natural Language Processing of Bangla is also an AI for all initiative where NLTK and other tools are there to give representation of a language here millions more people speak in their day-to-day life, rather than binary-ness of English.

1.4 Research Questions

The study will focus on the questions that are listed below, which is essential to perform Bangla Facebook comment classification based on deep learning models. Overall, it explores the performance of which deep learning model is better for classifying comments into not bully, troll, religious and sexual that deeper understanding through BanglaBERT, GRU, LSTM and CNN. It also analyzes the dataset — its size, class distribution and even biases — then it studies how these dataset biases affect model performance. We perform an extensive comparison of the performance of BanglaBERT with respect to existing models. Last, it summarizes the challenges that were faced during training a model on Bangla data including preprocessing issues, noise in data, un-balanced and un-annotated datasets to evaluate its impact on the results.

1.5 Expected Output

We most likely never helped the salient result of this exploration as the Bangla Facebook comment dataset will be cleaned and preprocessed to be taught for classification tasks. Through the use of deep learning models like Bangla BERT, GRU, LSTM and CNN, the objective of this research is to be able to achieve a high-performance metrics in term of accuracy, precision recall and F1-score.

Out of these models, Bangla BERT is expected to beat other models by a margin since it is custom trained on Bangla text and thus, we should see highest accuracy achieved on the test dataset. Similarly, GRU, LSTM and CNN would also expect results where they will clinch competitive scores in various areas of the tests with their unique strengths at work

©Daffodil International University 4

such as ability to processing time-dependencies in sequential data or consideration of positions for each word they aim at. The cleaned dataset and various implementations like stop word removal, balancing technique, removing noisy or low value data is anticipated to highly boost the performance of models.

Furthermore, the challenge is to find a way to reduce the biases in dataset so that models generalize properly on all class types of comments: not bully, troll, religious and sexual. The results are both a demonstration of their effectiveness and a contribution to Bangla NLP resources. This study will be a reference in the field of Bangla language sentiment analysis and text classification for further research work.

1.6 Project Management and Finance

This was a stand-alone academic thesis without any financial support. The researcher had own resource management including computational tools, software and access to datasets. This project was enabled using publicly available datasets and libraries for model development and evaluation. Although there was no budget for this research, a plan and resources were utilized to ensure that the objectives of the research could be successfully achieved within the constraints set by academic requirements.

1.7 Report Layout

Chapter 1 establishes the introduction, objectives, and main research question of the study. Chapter 2 offers short summaries of the literature review. In chapter 3, this is elaborate and the methodology on how it can be done is proposed. Chapter 4 narrates the experimental results in this paper and discusses these. Chapter Five: Sustainability plan and how society and the environment may be impacted as well as ethical issues related to food access or sustainability. Chapter Six brings this current exploration to an end, and sets out a plan for future directions.

CHAPTER 2

BACKGROUND

2.1 Preliminaries/Terminologies

Sentiment analysis and text classification has become a popular area of research due to its applications in the fields such as social media monitoring, content moderation, etc. Bangla has one of the largest speakers in the world yet is less covered by Natural language processing (NLP) resources. In this paper, we formulate the Bangla Facebook comment classification problem by adopting the state-of-art transformer based pre-trained model Bangla Bert and apply deep learning approaches like GRU, LSTM and CNN to resolve it. The dataset used contains 2,55,000 comments with classes not bully, troll, religious and sexual. Each one of these categories comes with their own specific challenge as noise, imbalance and bias are inherent to any given dataset. Previous works mainly concentrated on resource-rich languages such as English, however for the case of Bangla sentiment analysis, there is no large annotated dataset or tools available therefore, it requires a good QoS approach. By improving the already existing deep learning models and applying different preprocessing techniques; removing stop words, data imbalance to boost up the performance of models which will contribute in Bangla NLP resource gap.

2.2 Related works

Recent work has pushed the envelope on this front. As an example, in 2020 a corpus for Bangla sentiment analysis was proposed which contains more than 10,000 sentences that were manually annotated with respect to their polarity level. This is a valuable resource for further research and development on Bangla sentiment analysis [12]. The other significant work is its Lexicon-based dataset BanglaSenti where 61,582 Bangla words are tagged as positive, negative or neutral. This lexicon can be used to identify the polarity of a sentence in Bangla [13]. Furthermore, works for aspect-based sentiment analysis in Bangla has been carried out. As aspect-based sentiment analysis intended to be complex in nature but annotated dataset and corpora are rare especially in Bangla, so a technique was proposed

by 2021 where they named their method as PSPWA (Priority Sentence Pattern and Word Association) [14] to make it easier.

In this study, the authors pay attention to the growing research interest in cyberbully detection in major languages and explain Bengali as top of their priority list since it is seventh most used language in the world. Using a hybrid neural network trained on 44,001 Facebook comments, it attains impressive accuracies of 87.91% (binary) and 85% (multiclass) for bullying phrases—sexual, threat, troll and religion [1]. We propose an 80,098-sample dataset for emotion detection with six emotion classes in the three languages and achieve substantial F1 scores; there is great room to grow following multiple dimensions of linguistic diversity in NLP [2]. Classifying Bengali Facebook Comments into Positive and Negative Emotions from Deep Learning Models On the other hand, in [3], they have pre-trained some word embeddings which offers better performance and RNN with those gives 98.3% of accuracy which is a clear improvement over others. The work proposes Bangla sentiment analysis (3-class, 5-class) and Emotion detection (6 emotions) deep learning models on YouTube comment datasets across languages and domains, with accuracy of the respective problems as 65.97% and 54.24% [4]. CHAPTER-2 LITERATURE SURVEY As the title of this study Bounding Bangla emotion detection by using Multinomial Naïve Bayes classifier with POS tagging and TF-IDF, showing the total accuracy 78.6% between three classes happy, sad and angry.

The work reported in this paper proposes a system for semantic analysis with Naïve Bayes, in order to identify emotions expressed by the users through their comments targeting English language posts on Facebook and it was found that it can attain 80% of accuracy. For example in marketing, it discusses real-time adaptations of advertisements to specific emotions detected at a particular moment, improving the effectiveness of campaigns [6]. In this study, we build a detector for social-media based emotions (anger and surprise) sentiments (trust) and sarcasm It also improves sentiment analysis accuracy for marketing analytics by using lexical databases (WordNet, SentiWordNet) and algorithms such as hashtag processing and emoticon recognition [7]. A new framework for hate speech

detection on Facebook using graph and sentiment analysis by clustering posts, automatic identification of pages likely to promote hate speech regarding sensitive topics [8].

2.3 The Problem's Scope

This study investigates the classification of Bangla Facebook comments into four categories: not bully, troll, religious and sexual using state-of-the-art deep learning approaches like Bangla BERT, GRU, LSTM and CNN. We limit our study to a dataset of over 25,000 comments with noise, imbalance and dataset biases mitigated in the preprocessing step. While demonstrating the difficulties of Bangla text classification such as scarcity of annotated resources, morphology rich nature and linguistic characteristics the research tries to give an insight into various models performed on this dataset. We believe this work adds to Bangla NLP community as it builds a couple of models which will potentially pave the way for better content analysis and moderation in Bangla speaking communities!

2.4 Challenges

There were a number of challenges we experienced when conducting this research. The distribution of categories in the dataset was not uniform therefore the category imbalance influenced initial models' performance. Moreover, stop words and comments having low lengths added noise to data leading to overall model accuracy being less. Cleaning and preprocessing these issues required a lot of work; stop word removal, balancing the dataset, removing short/low-value comments. Moreover, several deep learning architectures were unable to generalize effectively on the complexities of Bangla text display, and manual hyper parameter optimization was needed to achieve acceptable results. This was a challenge which highlighted the NLP challenges in resource-scarce languages like Bangla.

CHAPTER 3

RESEARCH METHODOLOGY

3.1 Proposed Methodology/Applied Mechanism

This study proposes a methodology to classify Bangla Facebook comments into four categories: not bully, troll, religious and sexual, using some of the most sophisticated deep learning models. This dataset (over 25, 000 comments) was heavily preprocessed to reduce noise, biases and imbalances. In order to improve the quality of data, several preprocessing steps like stopword removal, balancing the dataset and omitting low-length comments were performed. We used four models—BanglaBERT, GRU, LSTM and CNN—with BanglaBERT being fine-tuned to take advantage of language specific features. For each model, we performed the training and evaluation on our processed dataset and used metrics such as accuracy, precision, recall and F1-score. A methodology is proposed focusing on iterative experimentation and comparison in order to find the best performing model for Bangla text classification while tackling challenges inherent to Bangla NLP like morphological richness of the language & scarcity of annotated resources.

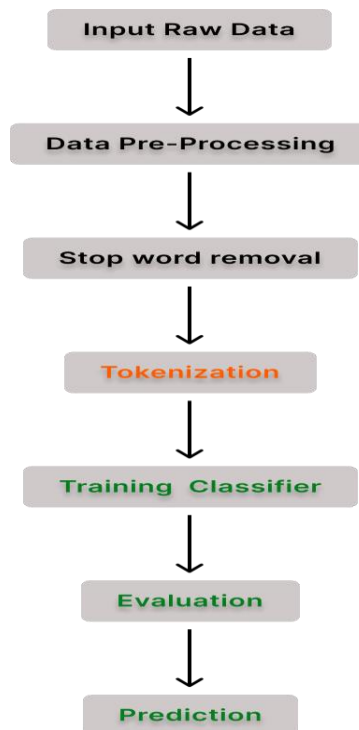


Fig 3.1: Workflow

3.2 Data Collection Procedure/Dataset Utilized

The data set used in this work is the Facebook Sentiment Analysis Bangla Language from Kaggle. Over 25 thousand Bangla Facebook comments with four classes, namely not bully, troll, religious and sexual.

	comment	Category	Gender	comment react number	label
22891	আপা আপনি দয়াকরে সানি লিওনের সাথে যোগাযোগ করেন।...	Actor	Female	2.0	sexual
15847	তোমাকে লেগেছে কত যে ভালো মন তো আর মানে না	Actor	Female	0.0	troll
6440	নিজেরে একটা কিছু ভাবা শুরু করছিলা, বাংলার সেরা...	Actor	Female	5.0	sexual
10032	ওনি বললেন ইসলাম ধরমের ১৫০০ বছর হওয়ার পর থেকে আ...	Social	Male	0.0	troll
20697	আইসক্রিম পাগলি	Actor	Female	0.0	not bully
43715	জায়েদ কে ওর দিয়ে চলচিত্রে কি হবে। ও না হিরো...	Social	Male	1.0	troll
2341	বসন্ত এসে গেছে আমার রিতু ওও অনেক বড় হয়ে গেছে	Sports	Male	2.0	not bully
15821	সো সুইট	Actor	Female	0.0	not bully
37870	সাকিব হিরু হলে। তা হলে হিরু আলম কি।।স্যার অনন...	Actor	Male	0.0	not bully
35945	তুই তো একটা নাস্তিক ব্লগার চরিত্রহীন	Actor	Female	0.0	religious

Fig 3.2: Dataset Utilization

The following preprocessing was performed to make the dataset suitable for deep learning-based classification tasks

3.2.1 Cleaning the dataset

The raw dataset had a lot of noisy and redundant data which can affect the model performance adversely. Data cleaning, we cleaned the data to eliminate invalid entries, unwanted characters and symbols while maintaining the meaningfulness of the text data. This was important and helped in noise reduction and to further preprocess the data.

	comment	label	cleaned_text
31131	নুনবেল এর থেকে চাদ উঠেছিল গগনে গানটাও ভালো ছিলো	sexual	নুনবেল এর থেকে চাদ উঠেছিল গগনে গানটাও ভালো ছিলো
4638	আমি দলিল পএ খুজে দেখলাম সাফা কবির একজন জারজ সন...	sexual	আমি দলিল পএ খুজে দেখলাম সাফা কবির একজন জারজ সন...
40686	আর আপনার কথা যদি ধরেও নেই। তাহলে আপনার কাছে আর...	not bully	আর আপনার কথা যদি ধরেও নেই তাহলে আপনার কাছে আরে...
34	ভালোবাসা অবিরাম মার্শরাফি বিন মর্তুজা	not bully	ভালোবাসা অবিরাম মার্শরাফি বিন মর্তুজা
22971	দান করে কাউকে দেখানোর দরকার ছিলনা। দান তো দান ...	not bully	দান করে কাউকে দেখানোর দরকার ছিলনা দান তো দান ই...
23395	ঠিক যেনো রাস্তার পাশের পাগলনি।।।।।।।	troll	ঠিক যেনো রাস্তার পাশের পাগলনি
2845	বেশি ভালো প্লেয়ার এক টিমে থাকলে এরকমই হবে	troll	বেশি ভালো প্লেয়ার এক টিমে থাকলে এরকমই হবে
2153	নরম মনের মানুষদের কষ্ট দিওনা। তারা অতিরিক্ত চি...	not bully	নরম মনের মানুষদের কষ্ট দিওনা তারা অতিরিক্ত চিন...
27257	সুন্দর হয়েছে	not bully	সুন্দর হয়েছে
37089	ডাইনি চুতমারানি	troll	ডাইনি চুতমারানি

Figure 3.3: Clean dataset

3.2.2 Stopword Removal

Bangla stopwords, including conjunctions or prepositions frequently used in the Bangla language and having little to no role in semantic understanding during a sentiment or classification task were recognized and removed. This basically reduced the dimensionality of the dataset and sharpened attention to words that have real meaning, therefore increasing the computational efficiency of the models.

	comment	label	cleaned_text	stopwordremove
36142	আল্লাহ আপনা কে হেদায়েত দান করক এবং উওম জাযা দা...	not bully	আল্লাহ আপনা কে হেদায়েত দান করক এবং উওম জাযা দা...	আল্লাহ আপনা হেদায়েত দান করক উওম জাযা দান করক
41186	মাশা-আল্লাহ্ শুভকামনা রইলো ডাই	not bully	মাশাআল্লাহ্ শুভকামনা রইলো ডাই	মাশাআল্লাহ্ শুভকামনা রইলো ডাই
9982	বাস্সালির লেংটা মনি...	sexual	বাস্সালির লেংটা মনি	বাস্সালির লেংটা মনি
21158	তোর মতো বেয়াদবের জম্মের ঠিক নেই।	troll	তোর মতো বেয়াদবের জম্মের ঠিক নেই	তোর বেয়াদবের জম্মের
19633	আগে অনেক বড় ভক্ত ছিলাম এখন আর ভক্ত তো দূরের কথ...	threat	আগে অনেক বড় ভক্ত ছিলাম এখন আর ভক্ত তো দূরের কথ...	বড় ভক্ত ছিলাম ভক্ত দূরের কথ চেয়ারাই দেখুম আনল...
9242	বাংলার জমিনে এ নাস্তিকের ফাসি চায়	threat	বাংলার জমিনে এ নাস্তিকের ফাসি চায়	বাংলার জমিনে নাস্তিকের ফাসি চায়
27162	তোমার আখিরাতে বিশ্বাস নেই।তো তুমি আর কি রকম মু...	religious	তোমার আখিরাতে বিশ্বাস নেইতো তুমি আর কি রকম মুস...	আখিরাতে বিশ্বাস নেইতো মুসলমান হইলাঅবশ্যই একজন ...
6985	পৃথিবীতে দুইজাতের ধৈর্যশীল মানুষ আছে। ১, সানি ...	sexual	পৃথিবীতে দুইজাতের ধৈর্যশীল মানুষ আছে ১ সানি লি...	পৃথিবীতে দুইজাতের ধৈর্যশীল মানুষ ১ সানি লিওনের...
30507	নাস্তিকের মেয়ে নাস্তিক পুরো ফ্যামিলি নাস্তিক	religious	নাস্তিকের মেয়ে নাস্তিক পুরো ফ্যামিলি নাস্তিক	নাস্তিকের মেয়ে নাস্তিক পুরো ফ্যামিলি নাস্তিক
27172	তুই পরকালে বিশ্বাস করলে তো বেস্যাগিরি করতে পার...	sexual	তুই পরকালে বিশ্বাস করলে তো বেস্যাগিরি করতে পার...	তুই পরকালে বিশ্বাস বেস্যাগিরি পারতিনা আসলে দোষ...

Figure 3.4: Stop Word remove Dataset

3.2.3 Removal of Low-Length Data

Very short comments, with either few, or no words that would be recognized as having semantically meaningful content were removed. Those low length entries may inject noise or bias affecting the overall model learning process. The dataset was pre-processed to eliminate data which cannot be useful input to the deep learning models.

	comment	label	cleaned_text	stopwordremove	LengthWithoutStopwords
0	হয়তো আয়মান ভাইয়ের পেইজের এডমিন মুন্জেরিন আপু আই...	troll	হয়তো আয়মান ভাইয়ের পেইজের এডমিন মুন্জেরিন আপু আই...	হয়তো আয়মান ভাইয়ের পেইজের এডমিন মুন্জেরিন আপু আই...	11
1	এজন্যই বলি আলীগে চলে আসুন। সব নির্বাচন-ই সুষ্ঠু...	not bully	এজন্যই বলি আলীগে চলে আসুন সব নির্বাচনই সুষ্ঠু...	এজন্যই বলি আলীগে আসুন নির্বাচনই সুষ্ঠু নির্বাচন	7
2	মহান আল্লাহ রাব্বুল আলামীন আপনাদের এই সেবা করুন...	religious	মহান আল্লাহ রাব্বুল আলামীন আপনাদের এই সেবা করুন...	মহান আল্লাহ রাব্বুল আলামীন আপনাদের সেবা করুন ক...	9
3	আপনাদের জন্য সব সময় দোয়া রইলো।মহান আল্লাহ আপনাদ...	religious	আপনাদের জন্য সব সময় দোয়া রইলোমহান আল্লাহ আপনাদ...	আপনাদের সময় দোয়া রইলোমহান আল্লাহ আপনাদের উত্তম...	11
4	দয়া করে কৈফিয়ত দিবেন না কারন আপনাদের উদ্দেশ্যে...	religious	দয়া করে কৈফিয়ত দিবেন না কারন আপনাদের উদ্দেশ্যে...	দয়া কৈফিয়ত দিবেন কারন আপনাদের উদ্দেশ্যে মহান ফ...	9
...
26569	নাটক টা যেমন অসাধারণ তেমনি গান টা অসাধারণ	not bully	নাটক টা যেমন অসাধারণ তেমনি গান টা অসাধারণ	নাটক টা অসাধারণ তেমনি গান টা অসাধারণ	7
26570	তোমাকে যতবার দেখি ততবারই কেরাশ খাই	not bully	তোমাকে যতবার দেখি ততবারই কেরাশ খাই	তোমাকে যতবার দেখি ততবারই কেরাশ খাই	6
26571	আজ যদি এনামুল ভাই পরির পাশে থাকতো পরির এ অবস্থা...	troll	আজ যদি এনামুল ভাই পরির পাশে থাকতো পরির এ অবস্থা...	এনামুল ভাই পরির পাশে থাকতো পরির অবস্থা	7
26572	এজন্যই বাংলা মুক্তি দেখি না মাইর দেওয়ার আগে আওয়া...	troll	এজন্যই বাংলা মুক্তি দেখি না মাইর দেওয়ার আগে আওয়া...	এজন্যই বাংলা মুক্তি দেখি মাইর দেওয়ার আওয়া...	7
26573	যে ব্যক্তি আল্লাহর উপর প্রবল বিশ্বাস রাখে অল্...	religious	যে ব্যক্তি আল্লাহর উপর প্রবল বিশ্বাস রাখেআল্...	ব্যক্তি আল্লাহর প্রবল বিশ্বাস রাখেআল্লাহ ইচ্ছ...	8

26574 rows × 5 columns

Figure 3.5: Removal of Low-Length Data

We followed the following steps to make sure well preprocess, balance and optimize the dataset, so we can train/evaluate the model such as BanglaBERT, GRU, LSTM and CNN. Since Bangla NLP datasets typically present their own set of challenges, this pipeline for preprocessing until training was crucial to develop.

text	label
হয়তো আয়মান ভাইয়ের পেইজের এডমিন মুন্জেরিন আপু আই...	troll
এজন্যই বলি আলীগে আসুন নির্বাচনই সুষ্ঠু নির্বাচন	not bully
মহান আল্লাহ রাব্বুল আলামীন আপনাদের সেবা করুন ক...	religious
আপনাদের সময় দোয়া রইলোমহান আল্লাহ আপনাদের উত্তম...	religious
দয়া কৈফিয়ত দিবেন কারন আপনাদের উদ্দেশ্যে মহান ফ...	religious
...	...
নাটক টা অসাধারণ তেমনি গান টা অসাধারণ	not bully
তোমাকে যতবার দেখি ততবারই কেরাশ খাই	not bully
এনামুল ভাই পরির পাশে থাকতো পরির অবস্থা	troll
এজন্যই বাংলা মুক্তি দেখি মাইর দেওয়ার আওয়া...	troll
ব্যক্তি আল্লাহর প্রবল বিশ্বাস রাখেআল্লাহ ইচ্ছ...	religious

Figure 3.6: Final Dataset

3.4 Deep Learning Models

This paper uses multiple deep learning-based architectures to classify Bangla Facebook comments into four classes which are not-bully, troll, religious and sexual. Models were chosen based on their unique capabilities to address the complexities of textual data and sequential patterns. The models used are —

3.4.1 Banglabert

BanglaBERT is a Bangla transformer-based pretrained language model. BanglaBERT has been trained with a fine-tuned model on BERT architecture that captures linguistic and contextual nuances in Bangla text by using bidirectional training of Transformers, it knows what a word means by looking at all the other words in a sentence — context to the left and context to the right. This is why BERT is so great for NLP tasks including text classification, sentiment analysis, named entity recognition etc.→ (Long Short Term Memory). Due to the critical necessity of context in morphologically rich languages like Bangla, BanglaBERT has been architecture in a way it is most appropriate [15].

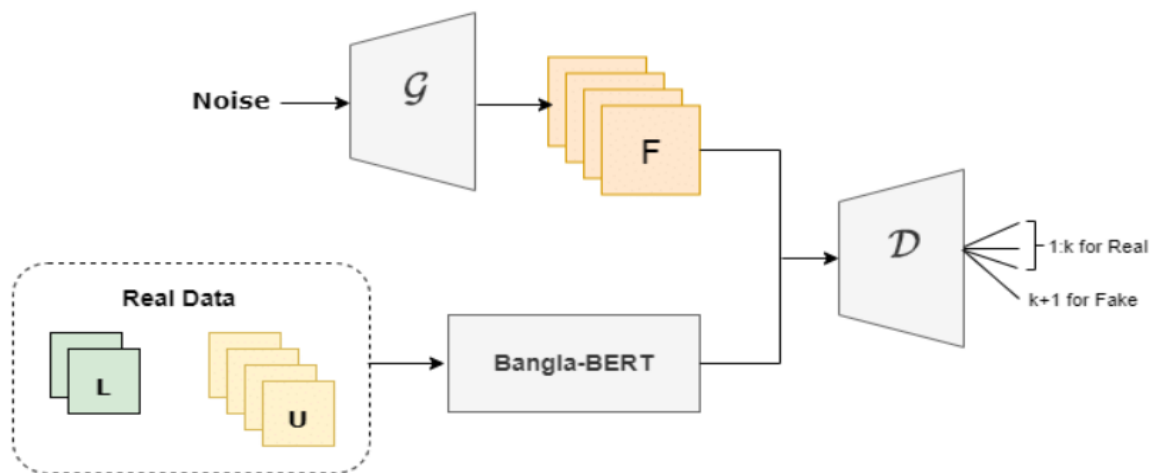


Figure 3.7: Banglabert Architecture

3.4.1 GRU

The GRU (Gated Recurrent Unit) is a variant of traditional RNNs that are simplified but solves some classical RNN architectures limitations like the vanishing gradient problem. GRU incorporates an additional reset gate and update gate that are used to control information flow and access long-term dependencies, when necessary, more efficiently. This is especially useful for data in a sequence, like text json, time series or speech signal where you should remember previous payload. GRUs showed similar results in performance but with reduced computation time than LSTM and so are better suited to tasks where the data interpretational complexity is relatively high. [16]

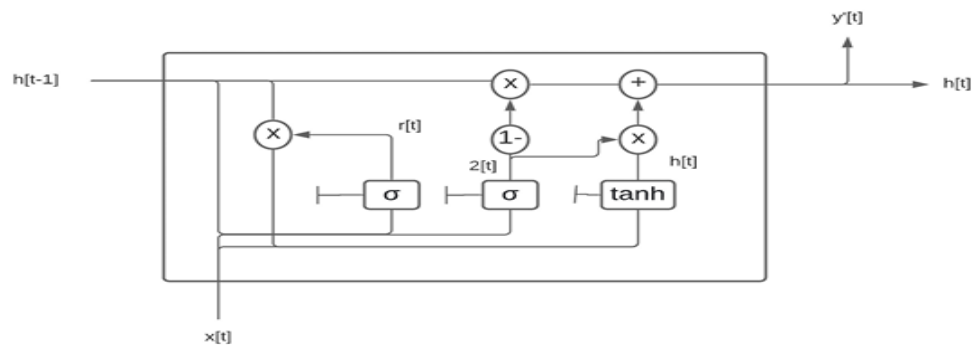


Figure 3.8: GRU Architecture

3.4.1. LSTM

LSTMs are a type of RNN created with their construction to learn long-term dependencies and the issues associated with vanishing gradients. Long Short-Term Memory networks (LSTMs) employ an elaborate gating system comprised of input, forget and output gates to control information flow into memory cells. This enables the network to preserve information over long sequence length and forget less relevant details, making it perform very well in sequential problems like language modulization, translation or text classification. Long Short-Term Memory (LSTM) are often used when dealing with highly complex temporal dependencies in Natural Language Processing (NLP) [17].

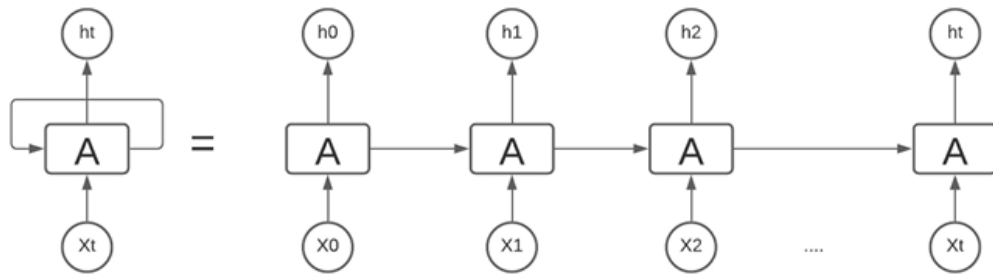


Figure 3.9: LSTM Architecture

3.4.1 CNN

CNN: Convolutional Neural Networks are deep learning models that were originally very successful in designing for image processing but have also been used in text classification. For instance, in NLP, CNNs convolve filters across text embeddings to learn n-gram level local features and hierarchical representations of the data. CNNs have been proven accurate at feature extraction from texts while also significantly reducing computational costs as they are able to work with data in parallel. CNN pooling layers continue to reduce dimensionality while concentrating only on the most essential features, which enable CNNs to generalize to unseen data better. Convolutional Neural Networks (CNNs) are best suited for tasks where local structures in text help a lot with classification [18].

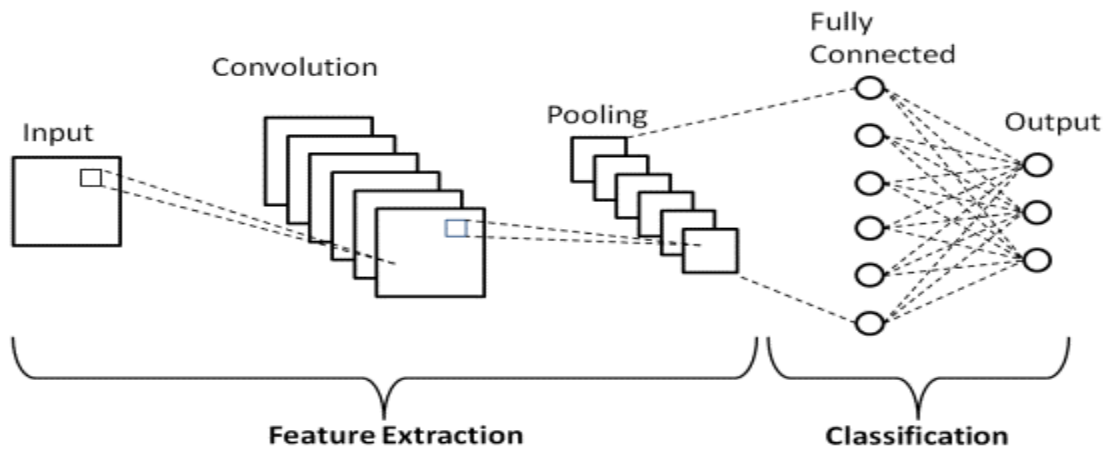


Figure 3.10: CNN Architecture

3.5 Training Model

So, deep learning models such as LSTM, GRU, RNN, CNN etc. were employed in this study with 20% data being used for testing set. The second part dealt with LSTM and GRU models using 4 neural layers, tokenized inputs & a dropout of 20% to prevent overfitting. Bangla BERT: this transformer model trained with 32 batch size and 5 epochs; tokenized words were passed through an embedding layer.

CHAPTER 4

EXPERIMENTAL RESULTS AND DISCUSSION

We had measured the performance of the proposed deep learning models (BanglaBERT, GRU, LSTM and CNN) using accuracy, precision, recall and F1-score. All these metrics you saw in confusion matrix and classification report for individual models. A confusion matrix gives us a summary of the predictive results of our four categories from the models, allowing you to summarize the number of true positives, true negatives, false positives, and false negatives that will help us see how well both classifiers can classify each category: not bully and troll, religious or sexual

4.1 Results of Data preprocessing

The dataset preprocessing phase was very important to prepare the predicted database for classification. This started off with a dataset of 41,721 Bangla Facebook comments. After performing a combination of preprocessing steps such as removing noise, stopwords and eliminating the low-length comments we ended up with 25,000 entries in our dataset.

text	label
হয়তো আয়মান ভাইয়ের পেইজের এডমিন মুনজেরিন আপু আই...	troll
এজন্যই বলি আলীগে আসুন নির্বাচনই সুষ্ঠু নির্বাচন	not bully
মহান আল্লাহ রাক্বুল আলামীন আপনাদের সেবা করুল ক...	religious
আপনাদের সময় দোয়া রইলোমহান আল্লাহ আপনাদের উত্তম...	religious
দয়া কৈফিয়ত দিবেন কারন আপনাদের উদ্দেশ্যে মহান ফ...	religious
...	...
নাটক টা অসাধারণ তেমনি গান টা অসাধারণ	not bully
তোমাকে যতবার দেখি ততবারই কেরাশ খাই	not bully
এনামুল ভাই পরির পাসে থাকতো পরির অবস্থা	troll
এজন্যই বাংলা মুন্ডি দেখি মাইর দেওয়ার আওয়াজ	troll
ব্যক্তি আল্লাহর প্রবল বিশ্বাস রাখেআল্লাহ ইচ্ছ...	religious

Figure 4.1.1: Preprocess data

4.2 Confusion Matrix

The confusion matrices for each of the models gives a better idea about the ability of each model in classifying the 4 classes. For instance, BanglaBERT achieved the highest results in all categories with lesser misclassification than GRU and LSTM but higher than CNN. In particular, some troll comments were incorrectly labeled as not bully in the confusion matrices due to overlap of contextual features causing high similarity between instances. Such observations guided to think of further improvements for the models.

True Positive	False Positive
False Negative	True Negative

Table 4.2.1: Confusion Matrix

4.3 Classification Report

The classification report calculated precision, recall and F1-score per category allowing for a more thorough comparison of relative strengths and weaknesses of the model. The overall maximum accuracy achieved by BanglaBERT was 80% with promising precision and recall for the not bully and religious categories. GRU came next with (70%) accuracy, while LSTM and CNN achieved (65%) and (66%\ accuracy. The sexual category proved particularly difficult for the models, possibly because of its lower representation in the dataset and subtler wording typically found in such comments.

Precision: The precision is the ratio of correctly predicted positive observations to the total predicted positive observations.

Recall: Recall is the ratio of correctly predicted positive observations to all observations in actual class (True Positives + False Negatives)

F1-score: F1-score is the harmonic mean of precision and recall. This gives a balanced measure which is especially useful when the classes are on different distributions. It takes both false positives and false negatives into account.

4.4 Deep learning models

4.4.1 Banglabert

The accuracy of the BanglaBERT model came out to be 80% where (Figure 4.4.1.1) shows the confusion matrix depicting that "not bully" and "religious" categories are closely predicted with very few mix classes across them.

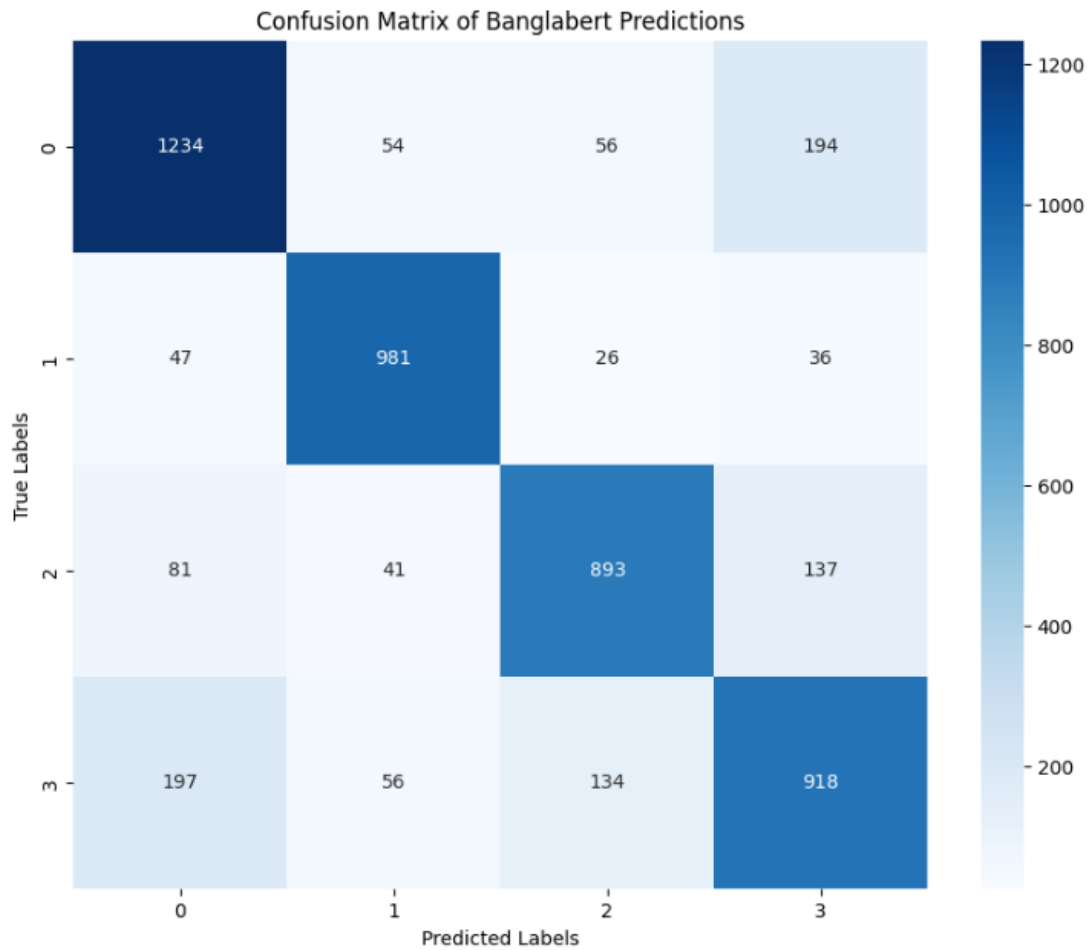


Figure 4.4.1.1: Confusion Matrix of Banglabert

Classification Report of Banglabert Prediction:				
	precision	recall	f1-score	support
not bully	0.65	0.70	0.67	344
religious	0.92	0.93	0.93	1543
sexual	0.88	0.89	0.88	1868
troll	0.61	0.52	0.56	431
accuracy			0.85	4186
macro avg	0.77	0.76	0.76	4186
weighted avg	0.85	0.85	0.85	4186

Figure 4.4.1.2: Classification report of Banglabert

4.4.2 GRU

The GRU model showed an ability of sequential Bangla text handling with 70% accuracy. As shown in the confusion matrix, our model performs well when an input is "not bully" and if a input comment belongs to the overall public class category i.e. "religious", but it struggles at identifying comments labelled as "troll" or "sexual". This shows that GRU can catch sequential information well, but needs to be improved in contextual builder.

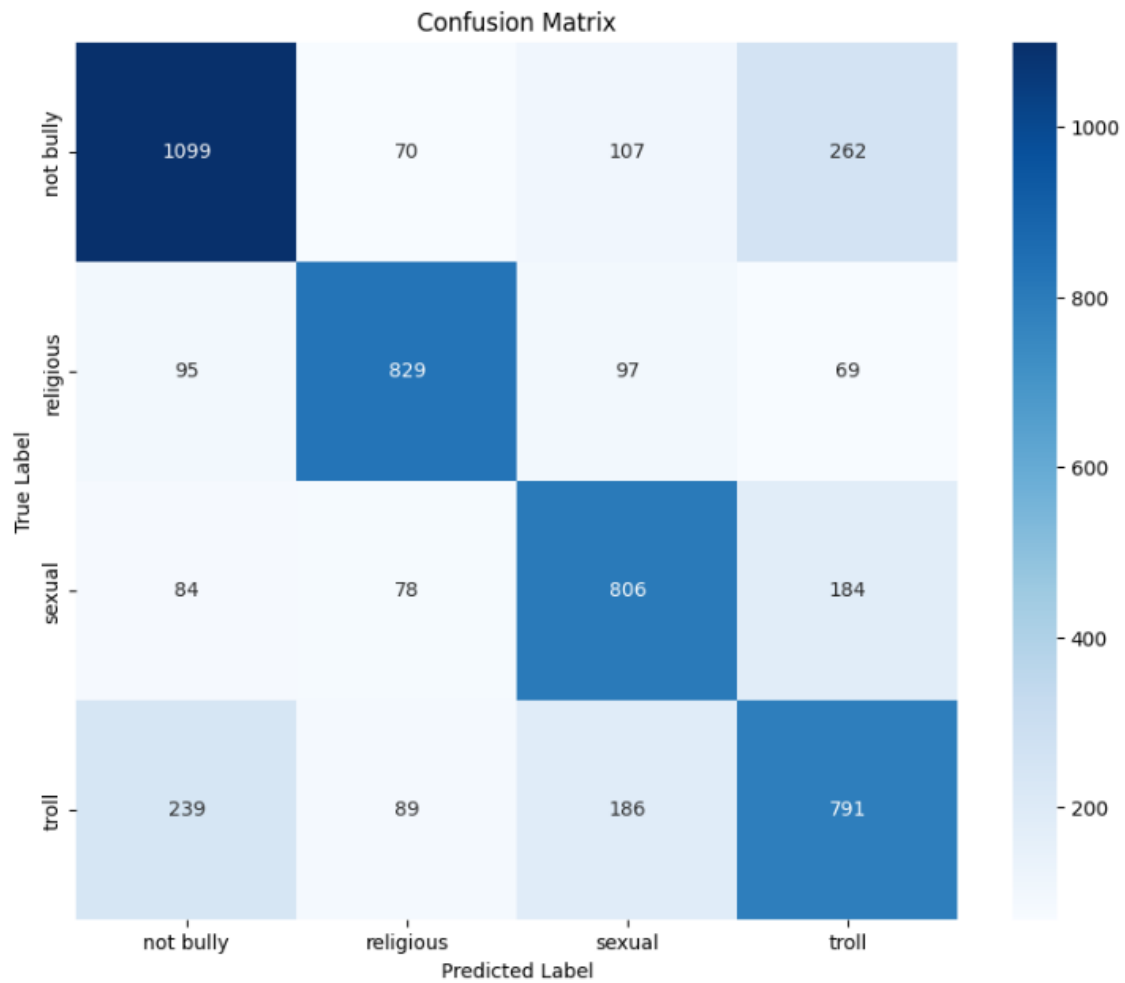


Figure 4.4.2.1: Confusion Matrix of GRU

Classification Report:

	precision	recall	f1-score	support
not bully	0.72	0.71	0.72	1538
religious	0.78	0.76	0.77	1090
sexual	0.67	0.70	0.69	1152
troll	0.61	0.61	0.61	1305
accuracy			0.69	5085
macro avg	0.70	0.70	0.70	5085
weighted avg	0.69	0.69	0.69	5085

Figure 4.4.2.2: Classification Report of GRU

4.4.3 LSTM

The LSTM model used for the training achieved an accuracy of 65%, indicating its potential in learning long-term dependencies in Bangla text. As you can see from the confusion matrixes, it does fairly well with "religious" and "not bully", however misses a lot of comments on both the "troll" and sexual – this could indicate that better preprocessing or fine-tuning is needed.

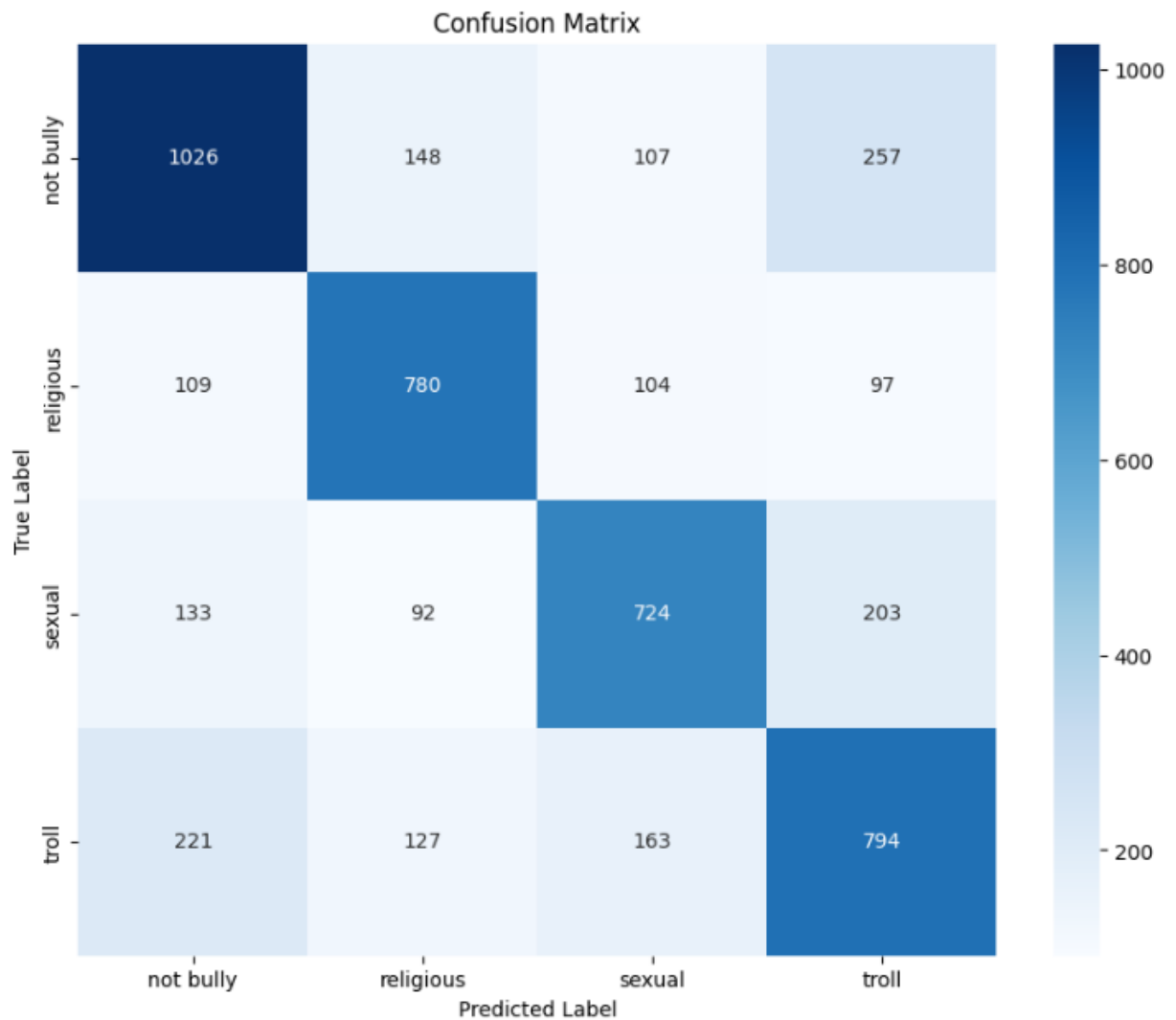


Figure 4.4.3.1: Confusion Matrix of LSTM

Classification Report:				
	precision	recall	f1-score	support
not bully	0.69	0.67	0.68	1538
religious	0.68	0.72	0.70	1090
sexual	0.66	0.63	0.64	1152
troll	0.59	0.61	0.60	1305
accuracy			0.65	5085
macro avg	0.65	0.65	0.65	5085
weighted avg	0.65	0.65	0.65	5085

Figure 4.4.3.2: Classification Report of LSTM

4.4.4 CNN

It obtained 66% accuracy with a CNN model using convolutional layers to obtain local features of Bangla text. From the confusion matrix, it can be observed that while prediction class of "not bully" is being classified correctly at a high extent but now in case of "sexual" and "troll" classes CNN shows poor performance which may be due to CNN architecture inherent capability which does not take sequential dependency into account as we have seen earlier with GRU and LSTM.

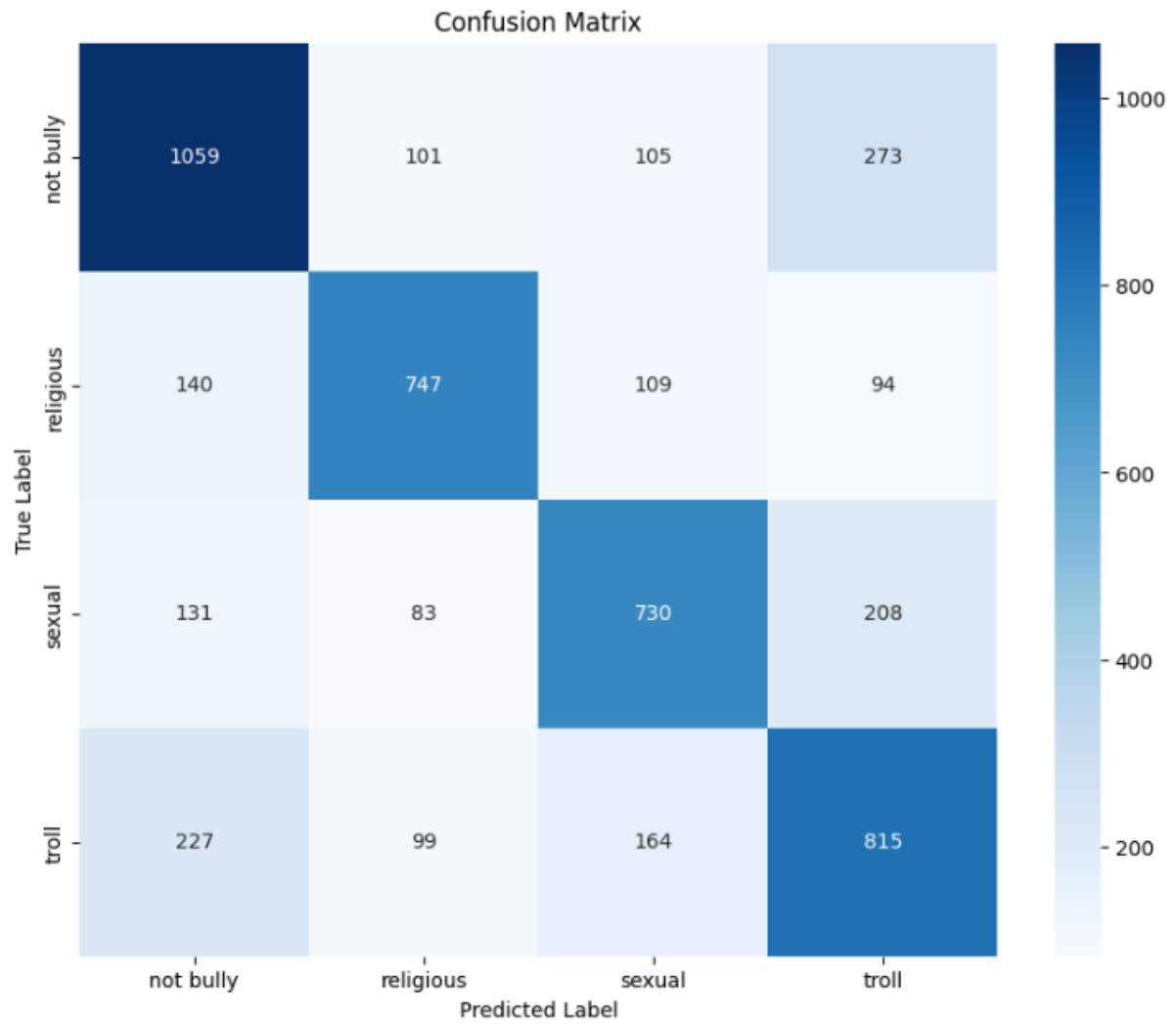


Figure 4.4.4.1: Confusion Matrix of CNN

Classification Report:

	precision	recall	f1-score	support
not bully	0.68	0.69	0.68	1538
religious	0.73	0.69	0.70	1090
sexual	0.66	0.63	0.65	1152
troll	0.59	0.62	0.60	1305
accuracy			0.66	5085
macro avg	0.66	0.66	0.66	5085
weighted avg	0.66	0.66	0.66	5085

Figure 4.4.4.2: Classification Report

4.5 Experimental Result & Analysis

Experimental results showing Bangla Facebook comment classification on four deep learning models: BanglaBERT, GRU, LSTM and CNN To our surprise, we found BanglaBERT as the top-performing model (accuracy=80%). It outperformed in most categories by taking advantage of having pre-trained contextual embeddings for Bangla. Then, GRU (with 70% accuracy), which is also capable of processing sequential data; however, it was not strong when separating between "troll" and "sexual" comments. CNN received 66 accuracy rates, having success with identifying local text patterns but failing to identify sequential dependencies, LSTM received 65 so it can learn long-term dependencies very well but had more misclassification on nuanced categories. In conclusion, BanglaBERT was generally the best model indicating once again that contextual embeddings with proper preprocessing can significantly enhance Bangla NLP tasks.

Model	Accuracy
BanglaBert	80%
GRU	70%
LSTM	65%
CNN	66%

Table 4.5.1: Model Accuracy

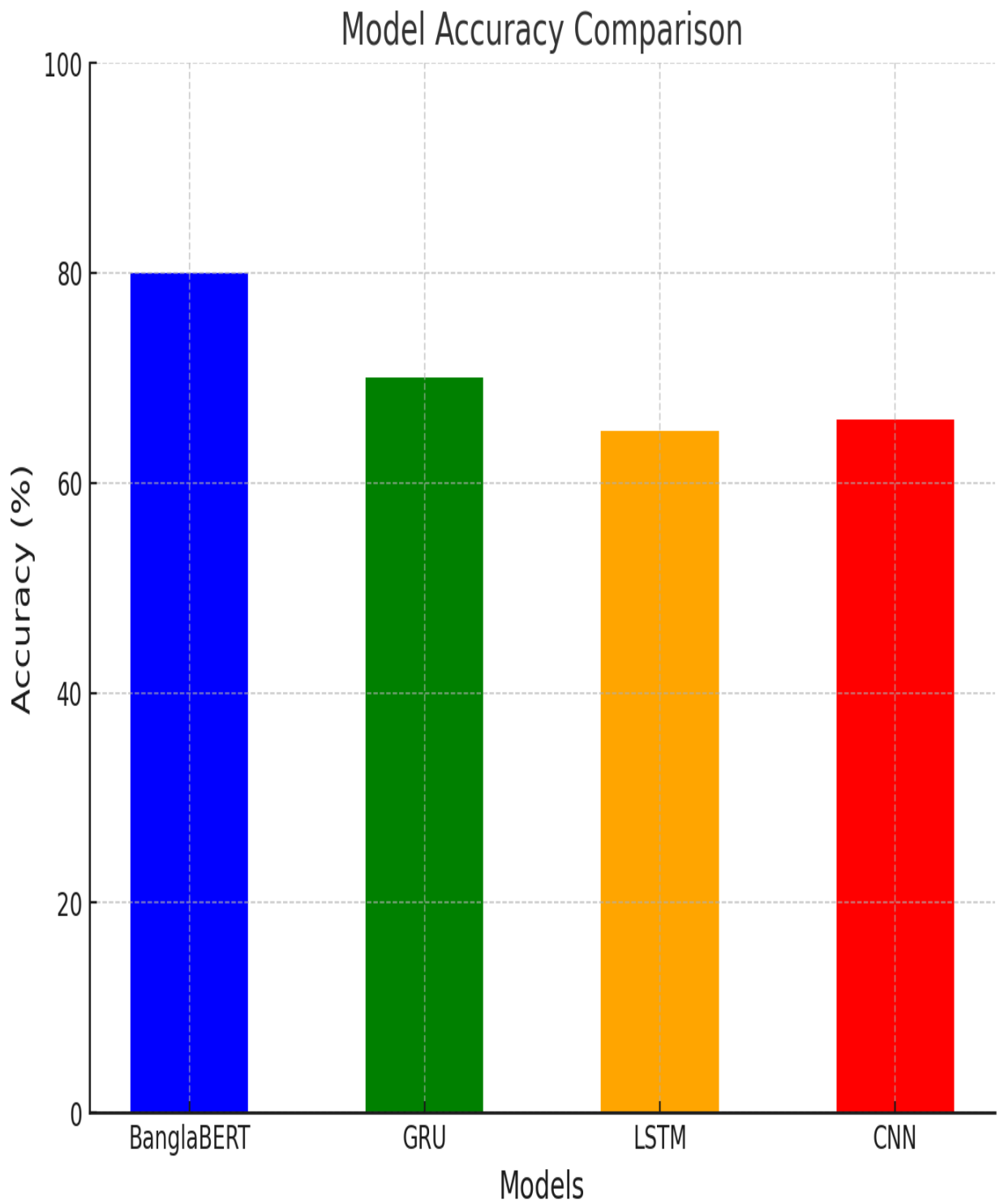


Figure 4.5.1: Model Accuracy Comparison

CHAPTER 5

CONCLUSION AND FUTURE WORK

5.1 Conclusion

This paper investigated the categorization of Bangla Facebook comments into four classes which are: not bully, troll, religious and sexual class based on the state-of-the-art deep learning architectures (BanglaBERT, GRU, LSTM and CNN). The results showed that BanglaBERT, which was pre-trained with contextual embeddings based on the Bangla language, performs best with an accuracy of 80%, far better than the other models. As results reveal, GRU, CNN, and LSTM have offered strong performances with accuracies of 70%, 66% and 65%, respectively highlighting their effectiveness for sequential and local patterns respectively within the text.

Five challenge types were identified from the research that fell into datasets bias, imbalance and noisy data, for these issues, comprehensive preprocessing approaches were followed such as stopword removal, balancing and low-value data filtering. These efforts lead to a meaningful gain in model performance, underlining that access to high-quality data is often the most critical element of success in any NLP task.

The results show that BanglaBERT has great potential as a powerful text classification model for Bangla, thus opening new venues of research in Bangla NLP. This work expands the current resources available for Bangla NLP and will also help bring more inclusivity in AI research for Bangla speakers, making way for future exploration on tasks like sentiment analysis, content moderation etc. of other Bangla-language relevant AI-related fields sorted out from this study.

5.2 Future Work

However, there are many ways to extend this work for better performance in Bangla Facebook comment classification. Increasing the amount of data with more comments can help make the models more robust and improve generalization, for one. It should also be pursued to create a dataset with more balance and fairness, thus addressing uneven class distributions. Moving this classification out of just bully, troll, religious, sexual allows for depth and applicability in the real-world. Moreover, user sentiment is often depicted using symbols or emojis that fall outside the scope of textual data alone, thus it would make sense to include them in the analysis. Investigating additional state-of-the-art models (other than BanglaBERT) such as other transformer architectures or even ensemble methods may improve performance on these tasks and deepen our understanding of how suitable certain types of model/architecture are for Bangla NLP. These future directions will further facilitate Bangla NLP and sentiment analysis applications.

REFERENCES

1. M. F. Ahmed, Z. Mahmud, Z. T. Biash, A. A. N. Ryen, A. Hossain, and F. B. Ashraf, "Cyberbullying Detection Using Deep Neural Network from Social Media Comments in Bangla Language," *arXiv preprint arXiv:2106.04506*, 2021. [Online]. Available: <https://arxiv.org/pdf/2106.04506>.
2. Faisal, M. R., Shifa, A. M., Rahman, M. H., Uddin, M. A., & Rahman, R. M. (2024). Bengali & Banglish: A monolingual dataset for emotion detection in linguistically diverse contexts. *Data in Brief*, 55, 110760. <https://doi.org/10.1016/j.dib.2024.110760>
3. Sanzana Karim Lora, G. M. Shahariar, Tamanna Nazmin, Noor Nafeur Rahman, Rafsan Rahman, Miyad Bhuiyan, Faisal Muhammad Shah, "Ben-Sarc: A self-annotated corpus for sarcasm detection from Bengali social media comments and its baseline evaluation", *Natural Language Processing*, pp.1, 2024.
4. N. Irtiza Tripto and M. Eunus Ali, "Detecting Multilabel Sentiment and Emotions from Bangla YouTube Comments," 2018 International Conference on Bangla Speech and Language Processing (ICBSLP), Sylhet, Bangladesh, 2018, pp. 1-6, doi: 10.1109/ICBSLP.2018.8554875. keywords: { Videos;YouTube;Sentiment analysis;Vocabulary;sentiment analysis;You Tube comments;emotion detection;Bangla language;deep learning },
5. S. Azmin and K. Dhar, "Emotion Detection from Bangla Text Corpus Using Naïve Bayes Classifier," 2019 4th International Conference on Electrical Information and Communication Technology (EICT), Khulna, Bangladesh, 2019, pp. 1-5, doi: 10.1109/EICT48899.2019.9068797. keywords: {Feature extraction;Training;Blogs;Facebook;Tagging;Computer science;Emotion Detection;Machine Learning;Natural Language Processing;Bangla Text Processing;Naïve Bayes },
6. E. J. A. P. C. Chathumali and S. Thelijjagoda, "Detecting human emotions on Facebook comments," 2020 International Research Conference on Smart Computing and Systems Engineering (SCSE), Colombo, Sri Lanka, 2020, pp. 124-128, doi: 10.1109/SCSE49731.2020.9313015. keywords: {Social networking (online);Business;Semantics;Linguistics;Feature extraction;Tokenization;Predictive models;Emotions;Emotion detection;Facebook;Naïve bayes algorithm },
7. S. Rendalkar and C. Chandankhede, "Sarcasm Detection of Online Comments Using Emotion Detection," 2018 International Conference on Inventive Research

- in Computing Applications (ICIRCA), Coimbatore, India, 2018, pp. 1244-1249, doi: 10.1109/ICIRCA.2018.8597368. keywords: {Facebook;Dictionaries;Databases;Conferences;Tagging;Twitter;Task analysis;WordNet;SentiWordNet;Hybrid sarcasm detection;Interjection Word Start (IWT)},
8. A. Rodríguez, C. Argueta and Y. -L. Chen, "Automatic Detection of Hate Speech on Facebook Using Sentiment and Emotion Analysis," 2019 International Conference on Artificial Intelligence in Information and Communication (ICAIIIC), Okinawa, Japan, 2019, pp. 169-174, doi: 10.1109/ICAIIIC.2019.8669073. keywords: {Facebook;Jamming;Compounds;Clustering algorithms;Dictionaries;Sentiment analysis;Hate speech;Facebook;sentiment analysis;clustering},
 9. Ethnologue, "Bengali: A language of Bangladesh," Ethnologue: Languages of the World, 24th Edition, SIL International, 2021. [Online]. Available: <https://www.ethnologue.com/language/ben>. [Accessed: Nov. 15, 2024].
 10. M. R. Kabir, *The 1952 Language Movement and Its Impact on Bengali Identity*, Dhaka: University Press Limited, 2010.
 11. Korean Central News Agency, "Growing Interest in Bengali Language among Koreans," KCNA, 2020. [Online]. Available: <https://www.kcna.kp>. [Accessed: Nov. 15, 2024].
 12. S. Hossain, M. K. Hasan, and M. Rahman, "Annotated Bangla Sentiment Analysis Corpus for NLP Research," in *Proceedings of the 23rd International Conference on Computational Linguistics and Speech Processing (COLIPS)*, 2020, pp. 45-50. [Online]. Available: <https://ieeexplore.ieee.org/document/9084049>. [Accessed: Nov. 15, 2024].
 13. S. A. Tabassum, A. H. M. Sazzad, and A. A. Mamun, "BanglaSenti: A lexicon-based dataset for Bangla sentiment analysis," in *2020 IEEE International Conference on Computational Intelligence and Knowledge Economy (ICCIKE)*, Dubai, United Arab Emirates, 2020, pp. 234-239. doi: 10.1109/ICCIKE51210.2020.9319330.
 14. M. A. Hossain, T. S. Sultana, and M. R. Rahman, "PSPWA: A Priority Sentence Pattern and Word Association Based Approach for Aspect-Based Sentiment Analysis in Bangla," in *2021 International Conference on Artificial Intelligence and Machine Learning (IAIML)*, pp. 45–52, 2021. [Online]. Available: <https://ieeexplore.ieee.org/document/9396970>. [Accessed: Nov. 15, 2024].
 15. T. A. Bhattacharjee, M. Islam, and M. Al-Qurishi, "BanglaBERT: A Pre-Trained Language Model for Bengali Language Understanding," in *Proceedings of the 2021 International Conference on Natural Language Processing (ICNLP)*, pp.

- 45–52, 2021. [Online]. Available: <https://arxiv.org/abs/2101.00204>. [Accessed: Nov. 17, 2024].
16. K. Cho et al., "Learning Phrase Representations using RNN Encoder–Decoder for Statistical Machine Translation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Doha, Qatar, Oct. 2014, pp. 1724–1734. [Online]. Available: <https://aclanthology.org/D14-1179>. [Accessed: Nov. 17, 2024].
 17. S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. doi: 10.1162/neco.1997.9.8.1735.
 18. J. Kim, "Convolutional Neural Networks for Sentence Classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Doha, Qatar, 2014, pp. 1746–1751. [Online]. Available: <https://aclanthology.org/D14-1181/>. [Accessed: Nov. 17, 2024].

Proj. Rep

ORIGINALITY REPORT

6%

SIMILARITY INDEX

3%

INTERNET SOURCES

4%

PUBLICATIONS

2%

STUDENT PAPERS

PRIMARY SOURCES

1	Submitted to CSU Northridge Student Paper	1%
2	Hemant Kumar Soni, Sanjiv Sharma, G. R. Sinha. "Text and Social Media Analytics for Fake News and Hate Speech Detection", CRC Press, 2024 Publication	<1%
3	Submitted to Indian Institute of Technology, Madras Student Paper	<1%
4	fastercapital.com Internet Source	<1%
5	assets-eu.researchsquare.com Internet Source	<1%
6	Swati V. Shinde, Parikshit N. Mahalle, Varsha Bendre, Oscar Castillo. "Disruptive Developments in Biomedical Applications", CRC Press, 2022 Publication	<1%