
Investigating the Predominance of Large Language Models in Low-Resource Bangla Language Over Transformer Models for Hate Speech Detection: A Comparative Analysis

[Fatema Tuj Johora Faria](#) , [Laith H. Baniata](#) ^{*} , [Sangwoo Kang](#) ^{*}

Posted Date: 17 October 2024

doi: 10.20944/preprints202410.1348.v1

Keywords:

Hate speech detection; Bengali language; Low resource language; Large language models; Few-shot learning; Zero-shot learning; Natural language processing






Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Article

Investigating the Predominance of Large Language Models in Low-Resource Bangla Language Over Transformer Models for Hate Speech Detection: A Comparative Analysis

Fatema Tuj Johora Faria ¹ , Laith H. Baniata ^{2,*} , Sangwoo Kang ^{2,*} 

¹ Department of Computer Science and Engineering, Ahsanullah University of Science and Technology, Dhaka, Bangladesh

² School of Computing, Gachon University, Seongnam 13120, Republic of Korea

* Correspondence: laith@gachon.ac.kr (L.H.B.); swkang@gachon.ac.kr (S.K.)

Abstract: The rise of abusive language on social media is a significant threat to mental health and social cohesion. For Bengali speakers, the need for effective detection is critical. However, current methods fall short in addressing the massive volume of content. Improved techniques are urgently needed to combat online hate speech in Bengali. Traditional machine learning techniques, while useful, often require large, linguistically diverse datasets to train models effectively. This paper addresses the urgent need for improved hate speech detection methods in Bengali, aiming to fill the existing research gap. Contextual understanding is crucial in differentiating between harmful speech and benign expressions. Large language models (LLMs) have shown state-of-the-art performance in various natural language tasks due to their extensive training on vast amounts of data. We explore the application of LLMs, specifically GPT-3.5 Turbo and Gemini 1.5 Pro, for Bengali hate speech detection using Zero-Shot and Few-Shot Learning approaches. Unlike conventional methods, Zero-Shot Learning identifies hate speech without task-specific training data, making it highly adaptable to new datasets and languages. Few-Shot Learning, on the other hand, requires minimal labeled examples, allowing for efficient model training with limited resources. Our experimental results show that LLMs outperform traditional approaches. In this study, we evaluated GPT-3.5 Turbo and Gemini 1.5 Pro on multiple datasets. To further enhance our study, we considered the distribution of comments in different datasets and the challenge of class imbalance, which can affect model performance. The BD-SHS dataset consists of 35,197 comments in the training set, 7,542 in the validation set, and 7,542 in the test set. The Bengali Hate Speech Dataset v1.0 & v2.0 includes comments distributed across various hate categories: personal hate (629), political hate (1,771), religious hate (502), geopolitical hate (1,179), and gender abusive hate (316). The Bengali Hate Dataset comprises 7,500 non-hate and 7,500 hate comments. GPT-3.5 Turbo achieved impressive results with 97.33%, 98.42%, and 98.53% accuracy. In contrast, Gemini 1.5 Pro showed lower performance across all datasets. Specifically, GPT-3.5 Turbo excelled with significantly higher accuracy compared to Gemini 1.5 Pro. These outcomes highlight a 6.28% increase in accuracy compared to traditional methods, which achieved 92.25%. Our research contributes to the growing body of literature on LLM applications in natural language processing, particularly in the context of low-resource languages.

Keywords: Hate speech detection, Bengali language, Low resource language, Large language models, Few-shot learning, Zero-shot learning, Natural language processing

1. Introduction

Bengali, spoken by approximately 260 million people, stands as the sixth most spoken language globally. It holds the distinction of being the second most spoken language in India and serves as the national language of Bangladesh. With nearly 205 million native speakers, Bengali ranks as the seventh most spoken native language worldwide, encompassing about 3.05% of the global population. Online social media platforms have evolved into crucial channels for communication, information sharing, opinion expression, and connection between individuals and businesses. However, they also contend with issues such as the proliferation of hateful or toxic content, bullying, and intimidation. Detecting such content manually on such large platforms is impractical, necessitating automated detection

systems. Yet, deploying effective automated detection remains challenging due to the dynamic nature of hate speech. In Bangladesh, approximately 81.7 million people use the internet, with 30 million actively engaging on social media, primarily via mobile phones. Notably, 42 million Bengali-speaking Facebook users participate in commenting, posting, and sharing content, constituting nearly 1.9% of all Facebook users. The use of Bengali on other social media platforms is also seeing a notable increase [1].

Despite extensive research on detecting abusive text in English on social networks, the Bengali language remains significantly underrepresented, even as its online presence grows. While the internet has promoted free speech, it has unfortunately also facilitated an increase in hate speech and vulgar language, particularly targeting Bangladeshi women. Social media has become an essential part of modern life, providing a popular and convenient platform for individuals to communicate and publicly express their thoughts. These platforms, along with online streaming services, have democratized information and amplified freedom of speech, often under the veil of anonymity. However, these same platforms also enable the spread of misinformation and hate speech, presenting significant challenges for regulatory authorities and law enforcement. Addressing hate speech is crucial for protecting human rights and preventing marginalization based on race, gender, ethnicity, or other affiliations. The ease of communication has unfortunately also led to harassment and attacks on individuals based on their expressions of sexism, racism, political opinions, or other concerns. Consequently, incidents of blackmail, cyberterrorism, and online harassment are proliferating rapidly across various social media platforms [2,3].

The vast amount of user-generated content necessitates the application of natural language processing (NLP) and machine learning (ML) models to effectively address the issue of online hate speech. Rapid advancements in artificial intelligence (AI) and machine learning technologies have led to numerous studies achieving promising results in this domain. However, state-of-the-art AI models predominantly rely on supervised learning techniques, which are generally limited to simple binary predictions of hate speech. A critical challenge in AI-based hate speech detection is the highly contextual nature of the problem. Existing supervised learning methods often fail to fully capture this context, resulting in inaccurate predictions. This underscores the need for detection methods capable of understanding and utilizing the full context of hate speech. Despite ongoing research and the development of diverse datasets for training classifiers, most efforts concentrate on English, neglecting low-resource languages like Bengali. The task is further complicated by challenges such as informal language syntax, spelling errors, and non-standard acronyms. Addressing these issues requires robust and context-aware AI models to ensure accurate and comprehensive hate speech detection across different languages and contexts [4,5].

Hate speech is a complex and evolving concept, highly dependent on the prevailing societal norms and the specific context in which it occurs. The deployment of advanced Large Language Models (LLMs) for content moderation is gaining popularity as a method to identify harmful and toxic content online. These models are trained to detect various forms of hate speech, both explicit and implicit. Recent studies have extensively explored the capabilities of models like GPT-3 and GPT-3.5 in hate speech detection. Notably, OpenAI's internal testing has highlighted the potential of GPT-4 as an effective tool for content moderation. Similarly, the latest open-source model, Llama 2, has shown promising results in identifying instances of hate speech [6,7]. These advancements underscore the growing role of LLMs in combating harmful online content, reflecting ongoing efforts to enhance digital safety and community well-being.

A crucial step in content moderation is the filtering of abusive content. A common method for achieving this is training language models on human-annotated content for classification. However, this approach presents several challenges, including the substantial resources required in terms of labor and expertise to annotate hateful content [8].

Additionally, this task exposes annotators to a wide array of hateful content, which is almost always psychologically taxing. Many studies have explored the potential of LLMs in detecting abusive

language, but none have examined the role of incorporating additional context as input to or output from such LLMs. This gap highlights the need for further research into how contextual information can enhance the effectiveness of LLMs in moderating abusive content [9–11].

This research aims to fill existing gaps by exploring the potential of LLM based approaches, focusing on prompt-based strategies, to enhance the detection of hate speech in Bengali. The objective is to develop a more effective methodology for identifying and mitigating hate speech in low-resource languages, thereby contributing to a safer and more inclusive online environment. For the first time, we introduce several variations of prompts and input instructions to probe two LLMs, specifically GPT-3.5 Turbo and Gemini 1.5 Pro across three datasets: BD-SHS, Bengali Hate Speech Dataset v1.0 & v2.0, and Bengali Hate Dataset. These datasets provide ground truth explanations, such as rationales or implied statements, that justify annotator decisions. Our approach involves designing prompts containing only the hate post as input and querying for the output label. We conduct a comprehensive evaluation of both zero-shot learning (ZSL) and few-shot learning (FSL) using prompting techniques to assess their capabilities in detecting hate speech. The selection of appropriate verbalizers significantly influences the effectiveness of these techniques, prompting us to systematically compare various verbalizers across multiple models in our study.

This research paper makes several key contributions to the field of hate speech detection. Our key contributions are summarized as follows:

- We have proposed to employ **class-balanced weights in the loss** function with pretrained language models. This approach adjusts the contribution of each class to the overall loss, ensuring that minority classes are given appropriate importance during training.
- We demonstrate that LLMs **surpass current benchmarks** in hate speech detection accuracy, achieving significant improvements over traditional methods.
- Our **Zero-Shot Learning** proves effective in identifying hate speech without task-specific training data, demonstrating adaptability to new datasets and languages.
- Our **Few-Shot Learning** enables efficient model training with minimal labeled examples, allowing for scalable real-time detection in low-resource language contexts.
- Our **detailed analysis** of prompting strategies reveals insights into optimizing LLMs for hate speech detection tasks, enhancing their applicability in diverse online communities.
- We conducted an **error analysis** within the framework of pretrained language models for hate speech detection.
- Our detailed **hallucination analysis** of Zero-Shot and Few-Shot learning strategies provides insights into optimizing LLMs for hate speech detection tasks, improving their suitability in a variety of internet groups.

The remainder of this paper is organized as follows: Section 2 offers an extensive review of related literature, establishing the foundation for our research. Section 3 explores background studies. Section 4 details the datasets used in the study. Section 5 describes the implementation details. Section 6 interprets the results. Section 7 examines the study's limitations, and Section 8 proposes directions for future research.

2. Literature Review

In this section, we provide a concise overview of earlier studies pertinent to our research on hate speech detection in the Bangla language. The overview is structured into four main categories based on the methodologies employed: Traditional Approaches, Deep Learning Approaches, Transformer-Based Approaches, and Large Language Model-Based Approaches. Below are summaries of the studies categorized in Table 1 and Table 2.

2.1. Traditional Based Approaches

This paper [12] addresses the pressing issue of hate speech and anti-social behavior on social media in Bangladesh, focusing on Bengali language comments. The authors collected 2,000 comments

from Facebook and YouTube and developed a dataset for analysis. Utilizing a Gated Recurrent Unit (GRU) neural network and various machine learning classifiers like Logistic Regression, Random Forest, Multinomial Naive Bayes (MNB), and Support Vector Machine (SVM), the study aimed to distinguish between social and anti-social comments. The GRU model achieved 78.89% accuracy, while MNB attained 80.51% accuracy. This research highlights the scarcity of Bengali language datasets and emphasizes the importance of context-specific feature extraction. The study's findings contribute significantly to the field by providing a benchmark dataset and demonstrating the effectiveness of GRU and MNB models in detecting anti-social comments. However, this paper did not explore the use of LLMs, and it has limitations in terms of scalability and the ability to handle nuanced context in comments. Similarly, another paper [13] highlights the increasing significance of social media in everyday life and the rising problem of negative comments on these platforms. While substantial work has been done on abusive text detection in other languages, similar research in Bengali is limited. This study utilizes a dataset comprising 5,000 comments collected from various social media platforms, such as Facebook and YouTube, with 2,698 labeled as abusive and 2,196 as non-abusive. Six machine learning algorithms—Logistic Regression, Multinomial Naive Bayes, Random Forest, Support Vector Machine (SVM), K-Nearest Neighbors (KNN), and Gradient Boosting—were employed, finding that SVM achieved the highest accuracy at 85.7%. The research emphasizes the challenges of using multi-class classification for Bangla and opts for binary classification to categorize comments as abusive or not. Additionally, the study underscores the importance of data preprocessing and feature extraction using TFIDF Transformer and Vectorizer. However, this paper also did not explore the use of LLMs, limiting its potential to capture complex linguistic patterns and context. Furthermore, another paper [14] discusses the rapid growth of the internet and social media, highlighting how these platforms have enabled free expression but also facilitated the spread of hate speech targeting individuals based on ethnicity, religion, gender, and other characteristics. This increase in online hate speech has led to disputes and cyberbullying, prompting organizations to seek effective solutions. The study focuses on detecting hate speech in Bangla videos using machine learning classification methods. Due to the lack of available datasets, the authors created a dataset from scratch, involving collecting, transcribing, and preprocessing videos from YouTube. They experimented with various machine learning and deep learning models, finding that logistic regression and the GRU model demonstrated the best accuracy. The GRU model achieved an impressive 98.89% accuracy, while the logistic regression model also showed high precision, recall, and F1 scores. Despite these successes, this paper did not explore the use of LLMs, which could have provided a more nuanced understanding of the transcribed video content.

2.2. Deep Learning Based Approaches

This research paper [15] addressed the exponential growth of social media, which, while empowering free expression, also facilitated online harassment and hate speech. It highlighted the lack of computational resources for under-resourced languages like Bengali. In response, it developed BengFastText, the largest Bengali word embedding model, based on 250 million articles. It created three extensive datasets for hate speech detection, document classification, and sentiment analysis. Experiments with a Multichannel Convolutional LSTM (MC-LSTM) network, incorporating BengFastText, demonstrated superior performance compared to baseline models, achieving high F1-scores. Similarly, another research paper [4] introduced the impact of social media platforms and online streaming services on the proliferation of hate speech. To address the need for linguistically diverse datasets, it introduced BD-SHS, a large manually labeled dataset including hate speech in various social contexts. BD-SHS contains over 50,200 offensive comments, making it significantly larger than previous Bangla hate speech datasets. The dataset was annotated using a hierarchical process that included hate speech identification, target identification, and categorization of hate speech types. For benchmarking, various models, including SVM and Bi-LSTM architectures, were experimented with. The Bi-LSTM model, trained with informal embeddings derived from 1.47 million social media

comments, achieved the highest F1-score of 91.0% in hate speech identification. This model consistently outperformed other pre-trained embeddings like BengFastText and multilingual fastText, which were trained on formal texts. Another study [1] highlighted the growing prevalence of hate speech on social media, particularly in Bengali, despite extensive research in other languages. It addressed this gap by proposing an encoder-decoder-based machine learning model to classify Bengali comments on Facebook. It collected a dataset of 7,425 comments across seven hate speech categories, using a combination of automated and manual methods due to limitations with the Facebook Graph API. The preprocessing involved tokenization, stemming, stopword removal, and extracting features using TF-IDF and word embedding. Three models—LSTM, GRU, and attention-based decoders—were evaluated, with the attention-based model achieving the highest accuracy at 77%. The study also incorporated a Bangla Emot Module to detect emotions from emojis and emoticons, enhancing the model's interpretability. Furthermore, another research paper [16] investigated the detection and classification of hateful speech in Bengali language social media, focusing on Facebook comments. The study employed both traditional machine learning algorithms and a GRU-based deep neural network model. It compiled and annotated a dataset of 5,126 comments, categorizing them into six classes: Hate Speech, Communal Attack, Inciteful, Religious Hatred, Political Comments, and Religious Comments. This dataset marks the first significant contribution to Bengali language hate speech detection in social media. Comparing various machine learning algorithms, it achieved 52.20% accuracy with Random Forest, which improved to 70.10% using the GRU model. The study highlighted the importance of linguistic and quantitative feature extraction tailored to the Bengali social context and underscored the superior performance of the GRU model in understanding context and semantics, which is critical for accurate hate speech detection in Bengali. Collectively, these studies emphasize the need for more advanced computational models and resources tailored to the Bengali language, highlighting the absence of LLMs in addressing hate speech detection comprehensively.

Table 1. Summary of Studies on Hate Speech Detection in Bengali Language

Types	Authors	Year	Models Employed	Performance Metrics	Key Findings
Traditional Approaches	Manash et al. [12]	2022	Gated Recurrent Unit (GRU), Logistic Regression, Random Forest, Multinomial Naive Bayes (MNB), Support Vector Machine (SVM)	GRU: 78.89% accuracy, MNB: 80.51% accuracy	Developed a dataset of 2,000 Bengali comments; highlighted scarcity of Bengali datasets and importance of context-specific feature extraction; MNB and GRU models effective in detecting anti-social comments.
	Sherin et al. [13]	2022	Logistic Regression, Multinomial Naive Bayes, Random Forest, Support Vector Machine (SVM), Gradient Boosting	SVM: 85.7% accuracy	Dataset of 5,000 comments; emphasized challenges in multi-class classification for Bangla; binary classification was used; highlighted importance of data preprocessing and TFIDF feature extraction.
	Istiaq et al. [14]	2021	Logistic Regression, Gated Recurrent Unit (GRU)	GRU: 98.89% accuracy	Created dataset from scratch with videos from YouTube; high accuracy with GRU model; logistic regression also showed high precision, recall, and F1 scores; focused on detecting hate speech in Bangla videos.
Deep Learning Approaches	Rezaul et al. [15]	2020	Multichannel Convolutional LSTM (MConv-LSTM), incorporating BengFastText	MConv-LSTM: F1-scores of 90.45%	Developed BengFastText, the largest Bengali word embedding model based on 250 million articles; created three extensive datasets; MC-LSTM with BengFastText outperformed baseline models.
	Nauros et al. [4]	2022	Bi-LSTM, Support Vector Machine (SVM)	Bi-LSTM: F1-score of 91.0%	Introduced BD-SHS, a large manually labeled dataset with over 50,200 offensive comments; Bi-LSTM trained with informal embeddings achieved highest F1-score; outperformed other pre-trained embeddings like BengFastText and MFT.

Table 1. Cont.

Types	Authors	Year	Models Employed	Performance Metrics	Key Findings
	Amit et al. [1]	2022	LSTM, GRU, Attention-based decoders	Attention-based model: 77% accuracy	Proposed an encoder-decoder-based model for classifying Bengali Facebook comments; collected 7,425 comments across seven hate speech categories; attention-based model achieved highest accuracy; included Bangla Emot Module
	Alvi et al. [16]	2019	GRU, Random Forest	GRU: 70.10% accuracy	Compiled and annotated a dataset of 5,126 comments into six classes; Random Forest achieved 52.20% accuracy, GRU model improved to 70.10%; emphasized importance of linguistic and quantitative feature extraction for Bengali.

2.3. Transformer Based Approaches

This research paper [3] emphasizes the critical need to address the proliferation of hate speech on social media in Bangladesh, particularly due to the lack of comprehensive Bangla datasets. It compiled a new dataset of 8600 user comments from Facebook and YouTube, categorized into sports, religion, politics, entertainment, and others. Various models were tested for hate speech detection, with BERT showing the highest accuracy at 80% on their dataset. When applied to an existing dataset of 30,000 records, BERT achieved an impressive accuracy of 97%, surpassing the performance of previously tested models like SVM, LSTM, and BiLSTM. However, there is a noticeable absence of Large Language Models specifically tailored for Bangla hate speech detection. Furthermore, another research paper [5] highlights that social media is a hotspot for hateful and offensive content, impacting race, gender, and religion in an unprejudiced society. Despite significant research on hate speech detection in English, there’s a notable gap in low-resource languages like Bengali, including its Romanized form used in social media interactions. To address this, it developed an annotated dataset of 10K Bengali posts (5K actual and 5K Romanized) and implemented several baseline models for classification. Experiments using m-BERT, XLM-Roberta, and IndicBERT were conducted, exploring interlingual transfer mechanisms. XLM-Roberta performed best when training datasets separately, while MuRIL outperformed in joint and few-shot training scenarios by better interpreting semantic expressions. Nonetheless, the absence of advanced Large Language Models for Bangla hate speech detection remains a significant gap. Moreover, this research paper [17] highlights the rapid expansion of social media and micro-blogging platforms, which have empowered freedom of expression but also facilitated the spread of antisocial behaviors such as online harassment, cyberbullying, and hate speech. The study evaluates four variants of the BERT model on a test set, reporting XML-RoBERTa as the top-performing model with an F1-score of 87%, outperforming other transformer models by 2% to 5%. Using WeightWatcher for ensemble prediction, they achieved the highest MCC score of 0.82, indicating a strong correlation between predictions and ground truths. Their ensemble approach improved overall accuracy by 1.8% across classes, effectively addressing misclassification rates. The effectiveness of BERT variants over traditional ML and DNN baselines was emphasized, particularly in minimizing classification errors across imbalanced datasets. The study provided class-specific classification reports,

highlighting nuances in detecting personal versus political hate speech, with specific challenges noted in identifying political hate due to overlapping terms with personal hate expressions. It also explored feature selection’s impact on ML baselines, noting significant improvements in models like GBT, which outperformed others with an MCC score of 0.571. However, SVM, LR, and NB models showed degraded performance due to feature selection’s impact on their assumptions of feature independence. In contrast, CNN and Bi-LSTM among DNN baselines performed reasonably well, although they fell short compared to transformer-based models like XML-RoBERTa. Despite these advances, the lack of Large Language Models specifically designed for Bangla hate speech detection persists as a major limitation.

Table 2. Summary of Studies on Hate Speech Detection in Bengali Language

Types	Authors	Year	Models employed	Em- ployed	Performance Metrics	Key Findings
Transformer Based Approaches	Jobair et al. [3]	2023	BERT, SVM, LSTM, BiLSTM		BERT: 80% ac- curacy on new dataset, 97% ac- curacy on exist- ing dataset	Compiled a dataset of 8600 com- ments; BERT showed highest ac- curacy at 80% on new dataset and 97% on existing dataset of 30,000 records; BERT outperformed SVM, LSTM, and BiLSTM.
	Mithun et al. [5]	2022	m-BERT, XLM- RoBERTa, IndicBERT, MuRIL		m-BERT: F1- score of 0.81	Developed an annotated dataset of 10K Bengali posts (5K actual, 5K Romanized); XLM-RoBERTa per- formed best in separate training; MuRIL outperformed in joint and few-shot training scenarios.
	Rezaul et al. [17]	2020	BERT vari- ants (includ- ing XLM- RoBERTa), traditional ML models, DNN models (CNN, Bi-LSTM)		XLM-RoBERTa: F1-score of 87%, MCC score of 0.82	Evaluated BERT variants; XLM- RoBERTa achieved highest F1- score of 87%; ensemble approach improved overall accuracy by 1.8%; highlighted challenges in detecting political hate speech; traditional ML models showed varied perfor- mance due to feature selection.
Large Lan- guage Models	Keyan et al. [7]	2024	GPT-3.5-turbo, Chain-of- Thought prompts		Accuracy: 0.85, Precision: 0.8, Recall: 0.95, F1 Score: 0.87	Chain-of-Thought reasoning prompts significantly outperform other strategies, capturing intricate contextual details for accurate hate speech detection.
	Sarthak et al. [8]	2023	Flan-T5-large, text-davinci- 003, GPT-3.5- turbo-0301		F1 Scores: Flan- T5-large: 0.59 (HateXplain), 0.63 (implicit hate), text- davinci-003: 0.45 (HateX- plain), 0.36 (implicit hate)	Flan-T5-large outperforms other models with vanilla prompts. In- corporating target community in- formation into prompts yields a 20- 30% performance boost. Precise prompt engineering is critical for optimizing LLMs in hate speech de- tection.
	Flor et al. [9]	2023	mT0, FLAN-T5, multilingual XLM-RoBERTa		Macro-F1 Scores: FLAN- T5: 65.34 (English), 62.61 (Spanish), 57.29 (Italian)	Zero-shot learning with prompting can match or surpass fine-tuned models’ performance, particularly with instruction fine-tuned models. Prompt and model selection signif- icantly impact accuracy.

2.4. Large Language Model Based Approaches

This research paper [7] addresses the critical issue of online hate speech detection, highlighting its contextual nature and the limitations of existing methods. It underscores the potential of LLMs

for context-aware detection due to their extensive training on diverse datasets. However, it notes the lack of effective prompting strategies for utilizing LLMs in this domain. Conducting a large-scale study using five established hate speech datasets, the researchers discovered that LLMs, especially with carefully crafted prompts, often surpass traditional models. They proposed four diverse prompting strategies, with the Chain-of-Thought reasoning prompt significantly outperforming others by capturing intricate contextual details. The Chain-of-Thought prompt achieved an accuracy of 0.85, precision of 0.8, recall of 0.95, and an F1 score of 0.87. This study emphasizes the importance of prompt engineering in optimizing LLMs for accurate hate speech detection, employing models such as GPT-3.5-turbo in their experiments. However, they only explored the English language for hate speech detection, leaving a gap for Bangla language applications. Furthermore, another research paper [8] provides a comprehensive analysis of various prompting strategies applied to LLMs for online hate speech detection. It demonstrates that Flan-T5-large outperforms other models with vanilla prompts, while text-davinci-003 shows superior results over GPT-3.5-turbo-0301. Incorporating target community information into prompts yields a 20-30% performance boost. Additionally, explanations and definitions as prompts enhance accuracy, though combining multiple strategies does not consistently yield further improvements. Detailed error analysis reveals frequent misclassifications, particularly for non-hate/non-toxic categories, underscoring the need for precise prompt engineering to optimize LLMs in hate speech detection. For the HateXplain and implicit hate datasets, Flan-T5-large achieved F1-scores of 0.59 and 0.63, respectively, outperforming gpt-3.5-turbo-0301 and text-davinci-003. However, this study also focused solely on the English language, indicating a gap in research for Bangla hate speech detection. Moreover, another research paper [9] identifies two key challenges in hate speech detection: the limited availability of labeled data and the high variability of hate speech across contexts and languages. Prompting offers a solution by enabling models to incorporate task-specific knowledge without labeled data. The study explores zero-shot learning (ZSL) with prompting for hate speech detection in three languages using eight benchmark datasets. Findings reveal that prompt selection significantly impacts results, with prompting—especially with recent large language models—often matching or surpassing fine-tuned models. This highlights the potential of prompting for under-resourced languages, showing that both prompt and model choice are crucial for accurate hate speech detection. Nonetheless, this research too only explored the English language, leaving a gap for Bangla language applications.

3. Background Study

3.1. Transformer Based Models

3.1.1. BERT-based Transformer Models

Bidirectional Encoder Representations from Transformers, commonly known as BERT [18], is built upon a sophisticated deep learning framework that establishes connections between input and output elements, adaptively determining their relationships. BERT's distinctive feature is its bidirectional training capability, which allows the model to understand the context of a word by considering both its preceding and succeeding words. This bidirectional approach contrasts with traditional models that only consider one direction, either forward or backward. BERT's architecture is designed to enable the model to capture intricate linguistic patterns and contextual dependencies in text, making it exceptionally powerful for a variety of natural language processing tasks. BanglaBERT Base [19], a variant tailored for the Bengali language, follows the same robust architecture as the original BERT model. This alignment ensures that BanglaBERT Base inherits the powerful contextual understanding capabilities of BERT, enabling it to effectively process and analyze Bengali text with high accuracy. mBERT [20], or Multilingual BERT, extends the BERT architecture to support multiple languages. It is pre-trained on a large corpus of text from 104 different languages, including Bengali, enabling it to understand and generate text across various languages. mBERT leverages the same bidirectional training approach as BERT, allowing it to capture the context of words from diverse linguistic backgrounds.

This makes mBERT a versatile tool for multilingual natural language processing tasks. XLM-RoBERTa [21], or Cross-lingual Language Model - RoBERTa, is an extension of the RoBERTa model designed to handle multiple languages. It is pre-trained on a vast dataset that includes text from 100 languages, which helps it learn cross-lingual representations. XLM-RoBERTa uses the same bidirectional training approach as BERT, allowing it to capture the context of words in a variety of languages. This makes it a powerful tool for multilingual tasks, such as cross-lingual understanding and translation.

3.1.2. ELECTRA-based Transformer Models

Efficiently Learning an Encoder that Classifies Token Replacements Accurately, known as ELECTRA [22], utilizes a unique pre-training task focused on identifying replaced tokens within an input sequence. This method involves two main components: a discriminator model and a generator model. The discriminator model is trained to recognize which tokens have been replaced in a corrupted sequence, while the generator model is simultaneously trained to predict the original tokens for the masked out ones. This setup is somewhat reminiscent of a generative adversarial network (GAN) training system, but without the adversarial component. In ELECTRA, the generator is not trained to deceive the discriminator, but rather to provide accurate token predictions. This collaborative process enhances the model's ability to understand and generate text. BanglaBERT [23] serves as the ELECTRA discriminator model, specifically designed for the Bengali language. This model effectively leverages the ELECTRA architecture to accurately classify and process Bengali text, improving the overall performance of text classification tasks in the Bengali language.

3.1.3. ALBERT-based Transformer Models

A Lite BERT, known as ALBERT [24], has shown that exceptional language models do not always require larger architectures. ALBERT achieves efficiency and high performance by utilizing the same encoder segment architecture as the original Transformer but introduces three crucial modifications: factorized embedding parameters, cross-layer parameter sharing, and employing Sentence-order Prediction (SOP) instead of Next Sentence Prediction (NSP). These modifications allow ALBERT to maintain a smaller model size while still delivering superior performance. Factorized embedding parameters reduce the number of parameters, cross-layer parameter sharing ensures consistency and reduces redundancy across layers, and SOP provides a more challenging pre-training objective than NSP, leading to better model understanding and contextual awareness. In the context of the Bengali language, sahajBERT¹ is a collaborative pre-trained ALBERT model that leverages these innovations. It utilizes masked language modeling (MLM) and Sentence Order Prediction (SOP) objectives to effectively learn from and understand Bengali text. This approach enables sahajBERT to perform various natural language processing tasks with high accuracy and efficiency, making it a powerful tool for Bengali language applications.

¹ <https://huggingface.co/neuropark/sahajBERT>

Table 3. Architectural and Training Objective Details of Various Pretrained Language Models

Architecture	Model	Layers	Attention Heads	Parameters	Objective Type During Training	Embedding Size
ELECTRA	BanglaBERT	12	12	110M	MLM with Replaced Token Detection (RTD)	768
BERT	BanglaBERT Base	12	12	110M	Masked Language Model (MLM)	768
	mBERT	12	12	110M	Multilingual Masked Language Model (MLM)	768
	XLM-RoBERTa	24	16	125M	Masked Language Model (MLM)	768
ALBERT	sahajBERT	24	16	18M	Multilingual Masked Language Model (MLM)	128

3.2. Large Language Models

3.2.1. GPT 3.5 Turbo

Generative Pre-trained Transformer 3.5 [25], or GPT-3.5, developed by OpenAI, is a sophisticated language model renowned for its ability to understand and generate human-like text. Built on the Transformer architecture, it leverages self-attention mechanisms to process and generate text efficiently. Key features of GPT-3.5 include its remarkable capability for few-shot learning, where it can perform tasks with minimal task-specific data. It also excels in zero-shot learning, requiring no examples, and one-shot learning, needing just a single example to generate responses. GPT-3.5 Turbo represents a significant advancement within the GPT series, particularly in comprehending and executing instructions accurately. This makes it ideal for tasks requiring specific formatting or outputs, such as creative content development. Developers can fine-tune the model to tailor its behavior to specific needs, enhancing performance across various applications. For instance, it can be adjusted to maintain consistent language use or to simplify prompts for desired responses. With an extensive context window capable of handling up to 16,385 tokens and improved precision in formatting, GPT-3.5 Turbo effectively addresses encoding challenges for non-English languages. It also offers rapid response times, restricted to generating outputs up to 4,096 tokens in length, making it a versatile and cost-effective solution for demanding text generation tasks.

3.2.2. Gemini 1.5 Pro

Gemini 1.5 [26] represents a significant leap forward, incorporating innovations in research and engineering across nearly every aspect of foundational model development and infrastructure. One of its standout features is the Mixture-of-Experts (MoE) architecture, which enhances efficiency in training and deployment. Unlike traditional Transformers that function as a single large neural network, MoE models are divided into smaller "expert" neural networks. These experts selectively activate based on the input type, greatly enhancing the model's efficiency. The first model in the Gemini 1.5 series is Gemini 1.5 Pro. This mid-size multimodal model is optimized for a wide range of tasks and performs at a level comparable to Gemini 1.0 Ultra, the largest model to date. A key innovation in Gemini 1.5 Pro is its breakthrough capability in long-context understanding. With a standard context window

of 128,000 tokens, Gemini 1.5 Pro is already impressive. Additionally, a select group of developers and enterprise customers can experiment with a context window of up to 1 million tokens through AI Studio and Vertex AI in a private preview. This extended context window allows the model to process vast amounts of information in a single prompt, enhancing the consistency, relevance, and usefulness of its outputs. Gemini 1.5 Pro's context window capacity significantly exceeds the original 32,000 tokens of Gemini 1.0. Now capable of handling up to 1 million tokens in production, Gemini 1.5 Pro can process extensive information in a single go, making it a powerful tool for large-scale data processing and analysis.

3.3. Evaluation Metrics

3.3.1. Accuracy

Accuracy [27] measures the proportion of true results (both true positives and true negatives) among the total number of cases examined. It is the most intuitive performance measure but can be misleading if the dataset is imbalanced, such as when hate speech instances are much fewer compared to non-hate speech instances. The formula for Accuracy is:

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \quad (1)$$

- TP = True Positives (correctly identified hate speech)
- TN = True Negatives (correctly identified non-hate speech)
- FP = False Positives (non-hate speech incorrectly identified as hate speech)
- FN = False Negatives (hate speech incorrectly identified as non-hate speech)

3.3.2. Precision

Precision [28], also known as Positive Predictive Value, measures the proportion of positive identifications (hate speech) that were actually correct. This metric is crucial in contexts where the cost of false positives is high, such as flagging non-hate speech content as hate speech. The formula for Precision is:

$$\text{Precision} = \frac{TP}{TP + FP} \quad (2)$$

3.3.3. Recall

Recall [28], also known as Sensitivity or True Positive Rate, measures the proportion of actual positives (hate speech instances) that were correctly identified. This metric is crucial when the cost of false negatives is high, such as failing to identify actual hate speech. The formula for Recall is:

$$\text{Recall} = \frac{TP}{TP + FN} \quad (3)$$

3.3.4. F1 Score

The F1 Score [27] is the harmonic mean of Precision and Recall. It provides a single metric that balances both Precision and Recall, especially useful when you need a balance between the two. The formula for F1 Score is:

$$\text{F1 Score} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (4)$$

4. Dataset Description

4.1. Dataset 1: BD-SHS

The BD-SHS [4] dataset is a comprehensive collection of Bangla hate speech (HS) comments sourced from various social media platforms. It consists of 50,281 comments meticulously annotated using a three-level hierarchical scheme by three annotators per comment. The majority decision among

annotators was adopted, resulting in a Fleiss Kappa score of 0.658, indicating moderate inter-annotator agreement. From this dataset, we have conducted Level 1 Hate Speech (HS) Identification, classifying comments as either Hate Speech (HS) or Non-Hate Speech (NH) based on the criteria outlined in Table 4. Figure 1 illustrating the overall distribution of comments and Figure 2 showcasing detailed examples of hate speech detection categories.

Table 4. Detailed Descriptions of Hate Speech (HS) and Non-Hate Speech (NH) Comment Categories

Category Name	Description
HS Comments	<ul style="list-style-type: none">• Attacks based on ethnicity, nationality, religion, sexual orientation, age, gender, disability, or disease.• Implicit support for hate speech acts without direct dehumanization.• Expressions advocating or supporting violence against individuals or communities.• Dehumanization through comparisons with animals, criminals, or historically vilified figures.• Expressions of disgust towards individuals or groups.
NH Comments	<ul style="list-style-type: none">• Lack of dehumanization or attacks based on the specified criteria.• Use of swear words not directed towards humans is not considered HS.

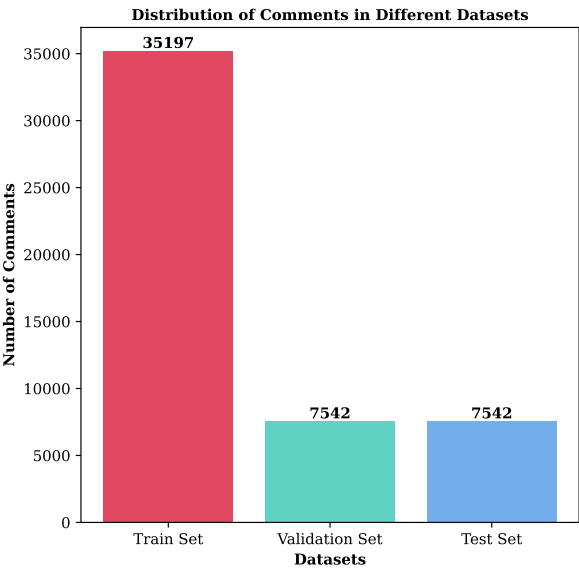


Figure 1. Visual Representation of Comment Distribution Across Datasets in Dataset 1

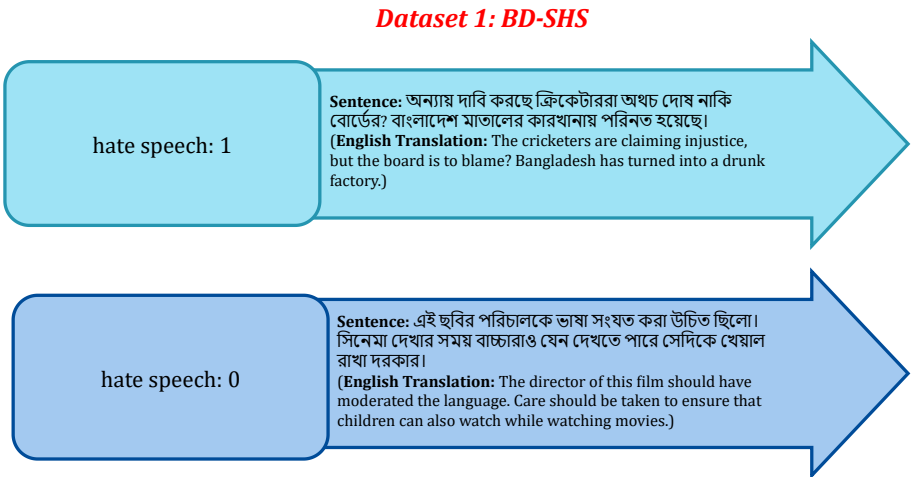


Figure 2. Detailed Visual Examples of Hate Speech Detection Categories in Dataset 1

4.2. Dataset 2: Bengali Hate Speech Dataset v1.0 & v2.0

The dataset [15,17] is an extensive compilation of Bengali articles from a variety of sources, including social media (Twitter, Facebook, LinkedIn), books, TV channels, news items from major newspapers, blogs, and sports websites. Using a bootstrapping approach, two linguists and three native Bengali speakers annotated hate speech from this dataset. At first, particular texts with common insults and signs of hate speech were found. The inclusion of 175 normalized abusive phrases that are frequently employed in hate speech in Bengali was the basis for a semi-automatic annotation process applied to a set of 10,000 statements, texts, or articles. If an annotation contained one or more of these specified terms, it was labeled as "hate". Notwithstanding difficulties in separating hate speech from unpleasant language and regional differences in hate speech categories, the procedure was centered on objective standards. Political, personal, gender-based abuse, geopolitical, and religious hate are among the categories of hate speech that have been recognized. The dataset contains 3,418 remarks (or roughly 3.5% of the annotated texts) that have been classified as hate speech. Three experts verified and edited the annotations, assuring their robustness and minimizing bias. An additional 3,000 labeled samples that fell into the categories of political, personal, geopolitical, religious, and gender abusive hate were added to the dataset. Semantic overlap made it difficult to differentiate between hate that is directed towards a person and hate that is directed towards a gender. Personal hate was defined as remarks that were antagonistic to specific people and primarily addressed to women in Bengali language. A bootstrap technique was used to collect data, with an emphasis on texts that contained particular sorts of slurs and phrases that were aimed towards individuals or groups. Texts were collected from newspapers, YouTube comments, and Facebook. The annotation procedure involved three annotators: an NLP researcher, a linguist, and a native Bengali speaker. Annotators were given objective materials to work with, and labels were decided upon by majority vote in order to reduce bias. In order to guarantee annotation quality and impartial criteria for decision-making, inter-annotator agreement was assessed using Cohen’s Kappa statistic. Figure 3 illustrating the overall distribution of comments and Figure 4 showcasing detailed examples of hate speech detection categories.

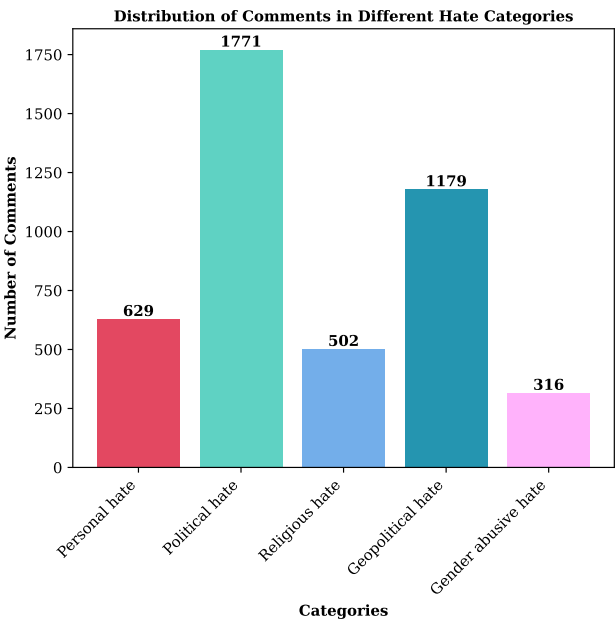


Figure 3. Visual Representation of Comment Distribution Across Datasets in Dataset 2

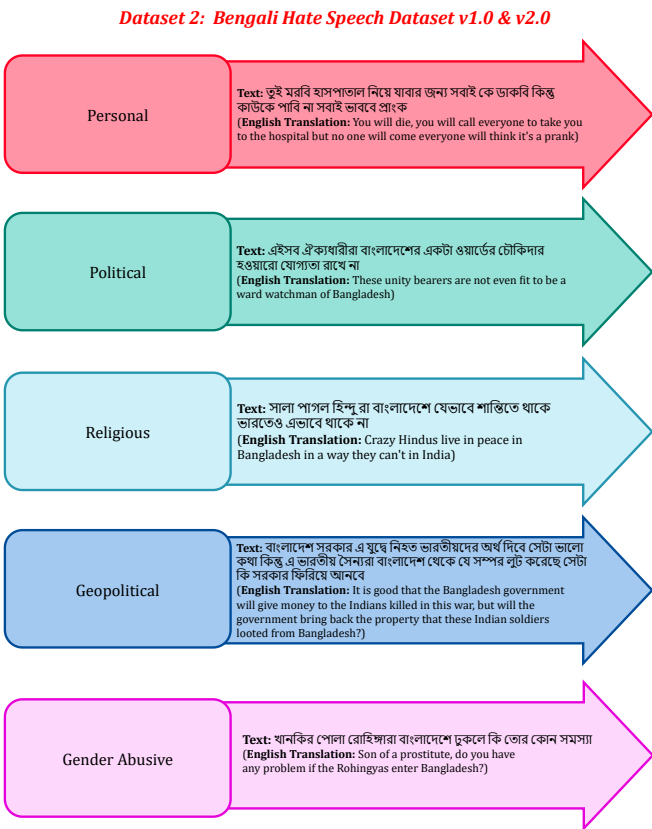


Figure 4. Detailed Visual Examples of Hate Speech Detection Categories in Dataset 2

4.3. Dataset 3: Bengali Hate Dataset

The dataset [2] comprises 15,000 Bengali posts collected from social media platforms such as Facebook and YouTube between January 2021 and April 2022. Initially, 110,000 posts were gathered, from which 8,500 posts containing Bengali profane words were filtered to focus on offensive content. To ensure class balance, an additional 8,500 non-offensive posts were selected from the original dataset.

After manual labeling and preprocessing to remove noise such as unidentified characters, symbols, and emojis, 15,000 posts were retained for analysis. The dataset annotation process involved 27 independent native Bengali labelers who meticulously categorized each sentence as either hate speech ('1') or non-hate ('0'), adhering to predefined guidelines. Figure 5 illustrating the overall distribution of comments and Figure 6 showcasing detailed examples of hate speech detection categories.

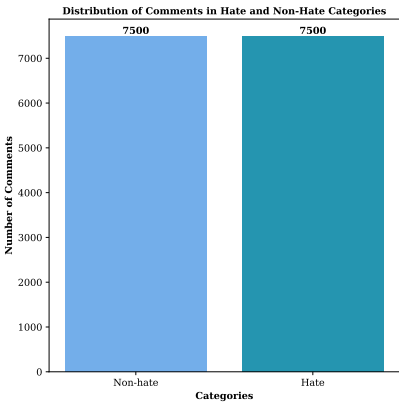


Figure 5. Visual Representation of Comment Distribution Across Datasets in Dataset 3

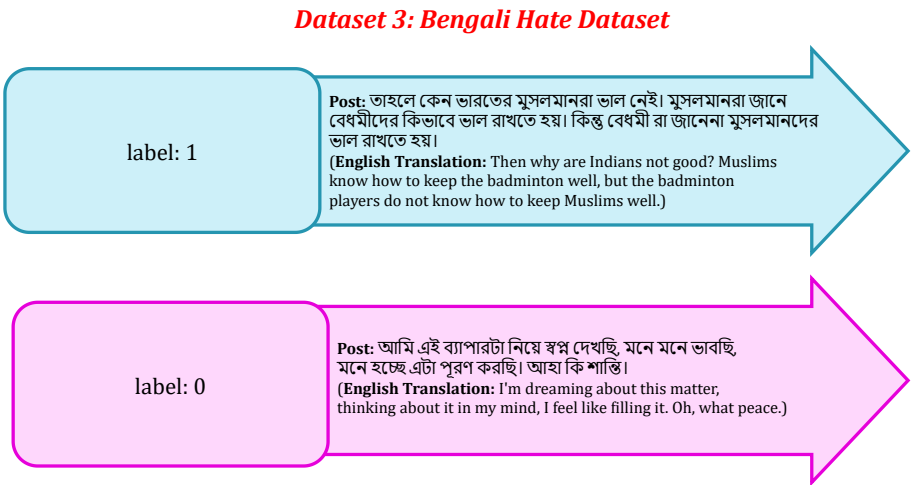


Figure 6. Detailed Visual Examples of Hate Speech Detection Categories in Dataset 3

5. Implementation Details

In this section, we detail our implementation of Bangla hate speech detection using advanced natural language processing techniques. Our approach harnesses the power of both pre-trained language models (PLMs) and large language models (LLMs) to achieve precise classification of hate speech and non-hate speech content in Bengali text. We begin by leveraging pre-trained language models (PLMs) to handle Bangla hate speech detection. This experiment involves tailored text preprocessing steps specific to Bengali, addressing class imbalance through effective strategies, fine-tuning procedures to adapt PLMs to the hate speech detection task, hyperparameter optimization for optimal performance, and evaluation using appropriate metrics. We also apply rigorous error analysis techniques and post-processing methods to refine prediction accuracy. In our second experiment, we explore the effectiveness of large language models (LLMs) in Bangla hate speech detection. This involves adapting LLMs to the Bengali language, incorporating sophisticated preprocessing steps, handling class imbalance challenges, fine-tuning LLMs on hate speech detection datasets, optimizing hyperparameters, and evaluating performance using standard metrics. Additionally, we investigate the potential of zero-shot and few-shot learning paradigms to further enhance model robustness across diverse datasets. The

code and supporting files for this study are publicly available for reference and use on GitHub at: <https://github.com/fatemafaria142/Bangla-Hate-Speech-Detection>

5.1. Experiment 1: Bangla Hate Speech Detection Using PLMs

5.1.1. Text Preprocessing

We used Dataset 1, Dataset 2, and Dataset 3, which undergo preprocessing to ensure compatibility with pre-trained language models (PLMs). A series of normalization steps are applied, specifically designed for Bengali text. These steps include handling whitespace by removing leading and trailing spaces and replacing multiple consecutive spaces with a single space, normalizing commas and other punctuation marks, correcting the placement of quotation marks, converting text to a consistent Unicode format, removing or replacing emojis depending on their relevance, normalizing numerical characters, identifying and removing English words unless contextually significant, removing common Bengali stop words, correcting common spelling mistakes, removing or replacing special characters and symbols that do not contribute to the text's meaning, and eliminating any other irrelevant information or noise present in the text.

5.1.2. Addressing Class Imbalance in Bangla Hate Speech Detection

Class imbalance, which can affect model performance, is a significant challenge in Bangla Hate Speech Detection. To address this issue, we employ class-balanced weights in the loss function. This approach adjusts the contribution of each class to the overall loss, ensuring that minority classes are given appropriate importance during training. The balanced weight is calculated based on the inverse frequency of each class in the dataset. We calculate class weights based on the inverse frequency of each class in the dataset. This function computes weights that are inversely proportional to class frequencies, helping to mitigate the effects of class imbalance.

5.1.3. Fine-tuning Procedure

Fine-tuning is a critical step in adapting PLMs for the specific task of Bangla hate speech detection. During this phase, the initialized PLMs are trained on a Bangla hate speech detection dataset using transfer learning techniques. This process involves updating the model parameters through gradient descent optimization algorithms to minimize the loss function. By fine-tuning, we leverage the pre-trained knowledge of the models while adapting them to the nuances of Bangla hate speech. To achieve this, we employ the AdamW optimizer, an extension of the Adam optimizer that includes weight decay. This helps prevent overfitting by regularizing the model parameters. The optimization process involves computing gradients of the loss function with respect to the model parameters and updating these parameters to reduce the loss. For the loss function, we use CrossEntropyLoss.

5.1.4. Optimization of Training Settings for Enhanced Model Performance

Hyperparameters are critical settings in machine learning that govern the training process and significantly impact model performance. These include parameters such as the learning rate, batch size, and number of epochs. Properly adjusting these hyperparameters during fine-tuning is essential to optimize performance and prevent overfitting. In our experiments, we tested batch sizes of 8, 16, and 32, and varied the number of epochs to 10, 15, 20, and 25. Additionally, we explored learning rates ranging from 0.01 to 0.001. By carefully tuning these parameters, we ensure that the fine-tuning process effectively adapts the pre-trained models to the specific task of Bangla hate speech detection, achieving optimal results.

5.1.5. Evaluation of Performance Metrics

The performance of fine-tuned PLMs is evaluated on a held-out test set using predefined evaluation criteria. Metrics such as Accuracy, Precision, Recall, and F1 score are utilized to objectively assess

model performance. Table 6, 7 and 8 presents comprehensive performance metrics, providing a clear understanding of the effectiveness of fine-tuned models for Bangla hate speech detection tasks.

5.1.6. Error Analysis for Insights into Model Performance

To gain deeper insights into our model's performance, we conduct a thorough error analysis. This step is crucial for identifying and understanding the specific instances where our models struggle, allowing us to pinpoint common patterns and challenges. By closely examining the misclassified instances, we can uncover underlying issues that may not be immediately apparent through overall performance metrics alone. Figure 13 provides a structured overview of the types of misclassifications, the number of instances for each type, the common patterns identified, the specific challenges faced by the models, and the potential improvements that can be implemented based on the analysis.

5.1.7. Post-processing for Enhancing Prediction Accuracy

In the final step of refining the model's predictions, post-processing techniques focus on enhancing accuracy and reliability by filtering out low-confidence predictions. After generating predictions, each is assessed based on its confidence score or probability estimate. A confidence threshold is set to determine the minimum acceptable level, adaptable to specific application needs and balancing precision and recall. Predictions meeting or surpassing this threshold are retained as high-confidence, reliable outputs. Those failing to meet it are filtered out or marked for further scrutiny, indicating potential uncertainty or ambiguity requiring additional review or correction within the model's output.

5.2. Experiment 2: Bangla Hate Speech Detection using LLMs

5.2.1. Data Selection

For our analysis, we randomly selected 400 data points from each of the "Dataset 1," "Dataset 2," and "Dataset 3" to assess the effectiveness of Zero Shot and Few Shot prompts. Specifically, we focused on data instances where the comments contained more than 5 words. Each label category, including Hate Speech (HS) and Non-Hate Speech (NH), comprises 100 instances from each dataset. In the case of "Dataset 2," which contains multi-label categories (Personal Hate, Political Hate, Religious Hate, Geopolitical Hate, and Gender Abusive Hate), we ensured that the selected samples represented a balanced mix of these categories. These samples were extracted from the training subsets of the respective datasets. Considering the cost implications associated with using the OpenAI GPT-3.5 Turbo API and Gemini 1.5 Pro API, we opted to limit our experiments to this subset of the dataset rather than the entire corpus.

5.2.2. Prompting Template

a) Zero-shot prompting for Bangla Hate Speech Detection

Zero-shot prompting for Bangla Hate Speech Detection involves guiding a model to classify text as either Hate Speech or Non-Hate Speech without prior training on such labels. In this context, hate speech includes offensive or derogatory language directed at individuals or groups, while non-hate speech conveys neutral or positive content. The model's objective is to accurately classify the text based solely on its content.

Figure 7 illustrates the application of zero-shot prompting for Bangla Hate Speech Detection using the "Dataset 1" dataset. The dataset consists of premise-text pairs labeled as either Hate Speech or Non-Hate Speech. The figure compares the performance of two models: GPT-3.5 Turbo and Gemini 1.5 Pro. Both models classify the text based solely on its content without prior training on the specific labels in the dataset.

Figure 9 focuses on the multi-label prompting approach for the "Dataset 2" dataset, guiding the model to predict one or more specific hate speech categories associated with a given text. The

categories include Personal Hate, Political Hate, Religious Hate, Geopolitical Hate, and Gender Abusive Hate.

Figure 11 Illustrates the application of Zero-shot prompting in Bangla Hate Speech Detection using "Dataset 3". Similar to Tables 6 and 7, this Table 8 demonstrates how the model can classify text accurately as hate speech ('1') or non-hate ('0') without any training data. It highlights the model's ability to generalize from pre-existing knowledge to classify hate speech in Bengali texts.

b) Few-shot prompting for Bangla Hate Speech Detection:

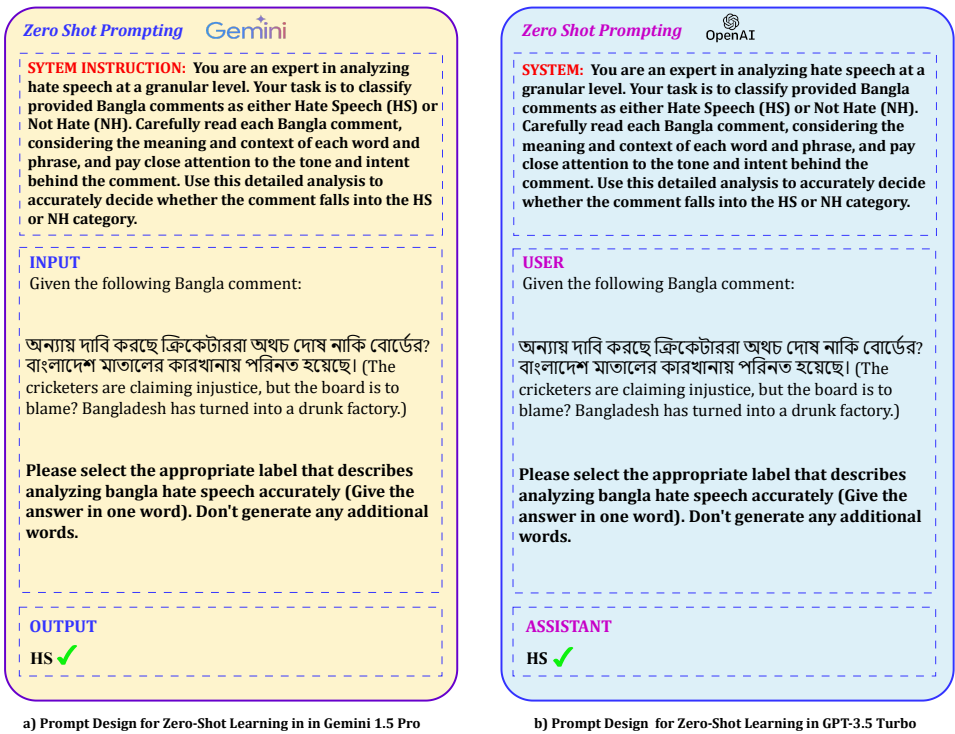
Few-shot prompting in Bangla Hate Speech Detection involves providing the model with a limited number of example texts labeled as Hate Speech or Non-Hate Speech to guide its predictions regarding the content expressed in subsequent texts. The model aims to determine the hate speech label associated with each text based on the provided examples, without extensive training on hate speech-labeled data. For instance, in a 5-shot scenario, the model receives five labeled examples for training before making predictions on new data, while in a 10-shot scenario, it receives ten examples, and so forth. In a 15-shot scenario, the model receives fifteen labeled examples, enhancing its ability to understand and classify the content expressed in the text accurately. The model utilizes the provided examples to inform its predictions, improving its classification performance as the number of examples increases.

Figure 8 showcases the application of Few-shot prompting in Bangla Hate Speech Detection using "Dataset 1". This Table 10–12 presents the labeled text examples used for training the model with a limited number of example texts labeled as Hate Speech (HS) or Non-Hate Speech (NH). The model, without extensive training on hate speech-labeled data, accurately predicts the hate speech label associated with each text based on the provided examples.

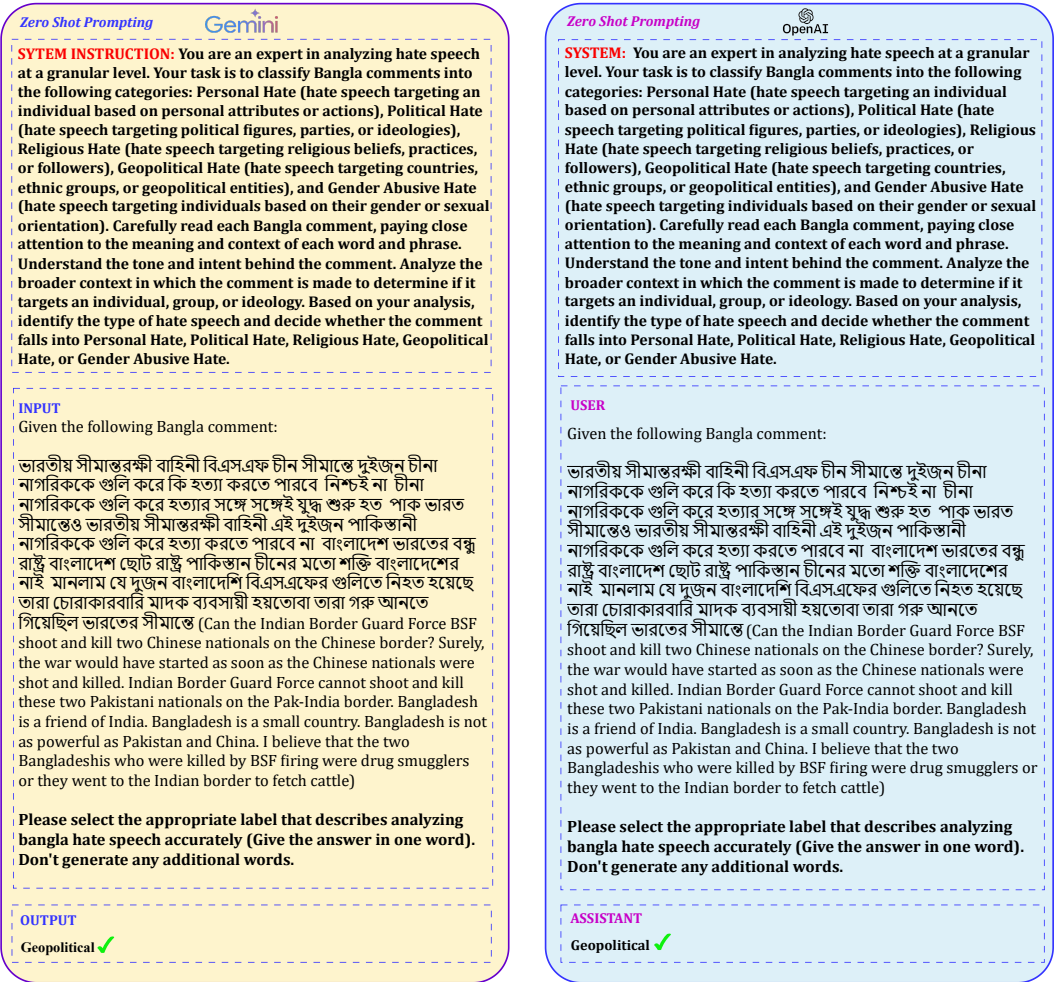
Figure 10 demonstrates the Few-shot capabilities of Gemini Pro in Bangla Hate Speech Detection using "Dataset 2". It illustrates how Gemini Pro effectively makes accurate predictions with a small number of training instances. The Table 10–12 showcases the model's ability in a multi-label few-shot prompting scenario, where it is trained on a limited number of example texts. Each example is labeled with one or more hate speech categories, including Personal Hate, Political Hate, Religious Hate, Geopolitical Hate, and Gender Abusive Hate. Gemini Pro uses this training to accurately predict and classify hate speech categories in subsequent texts, showcasing its robustness in handling diverse forms of hate speech in Bengali language.

Figure 12 illustrates the application of Few-shot prompting in Bangla Hate Speech Detection using "Dataset 3". This figure demonstrates how the model learns to classify text accurately as hate speech ('1') or non-hate ('0') with minimal training data. The Table 10–12 showcases the model's ability in a multi-label few-shot prompting scenario. It highlights the model's ability to generalize from a small number of labeled examples to classify hate speech in Bengali texts.

Dataset 1: BD-SHS



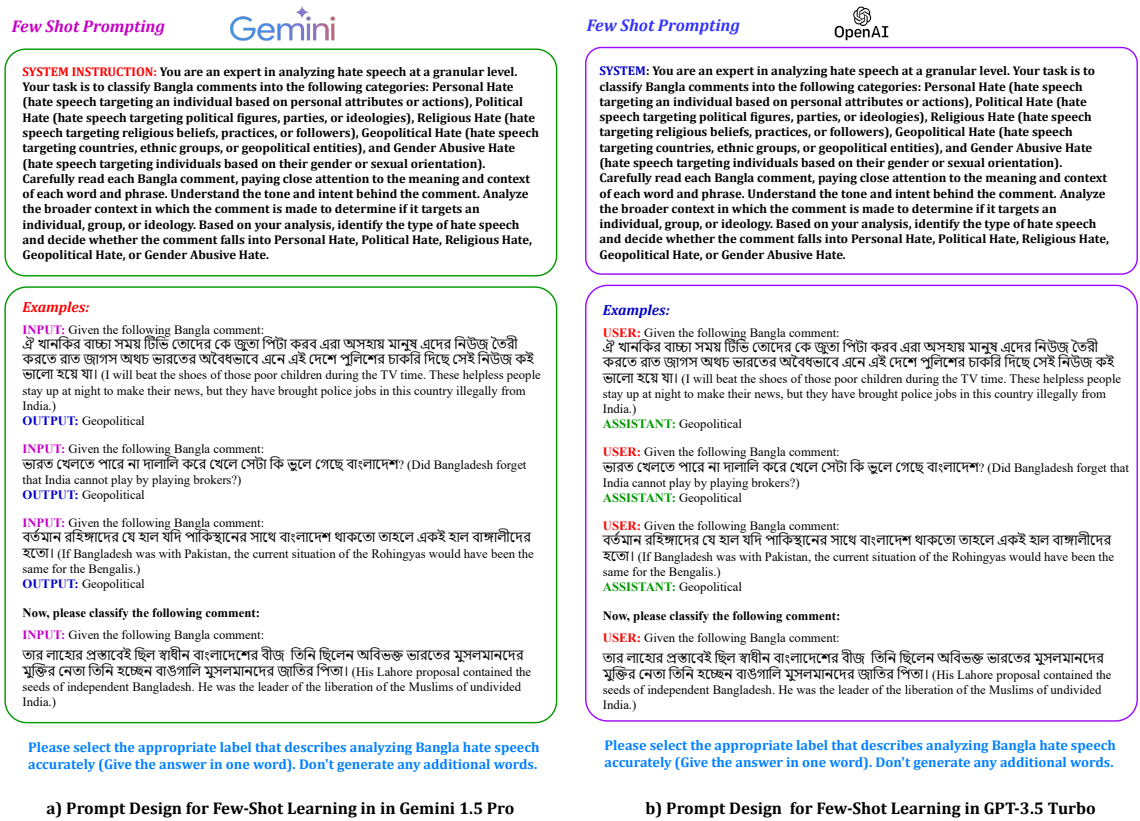
Dataset 2: Bengali Hate Speech Dataset v1.0 & v2.0



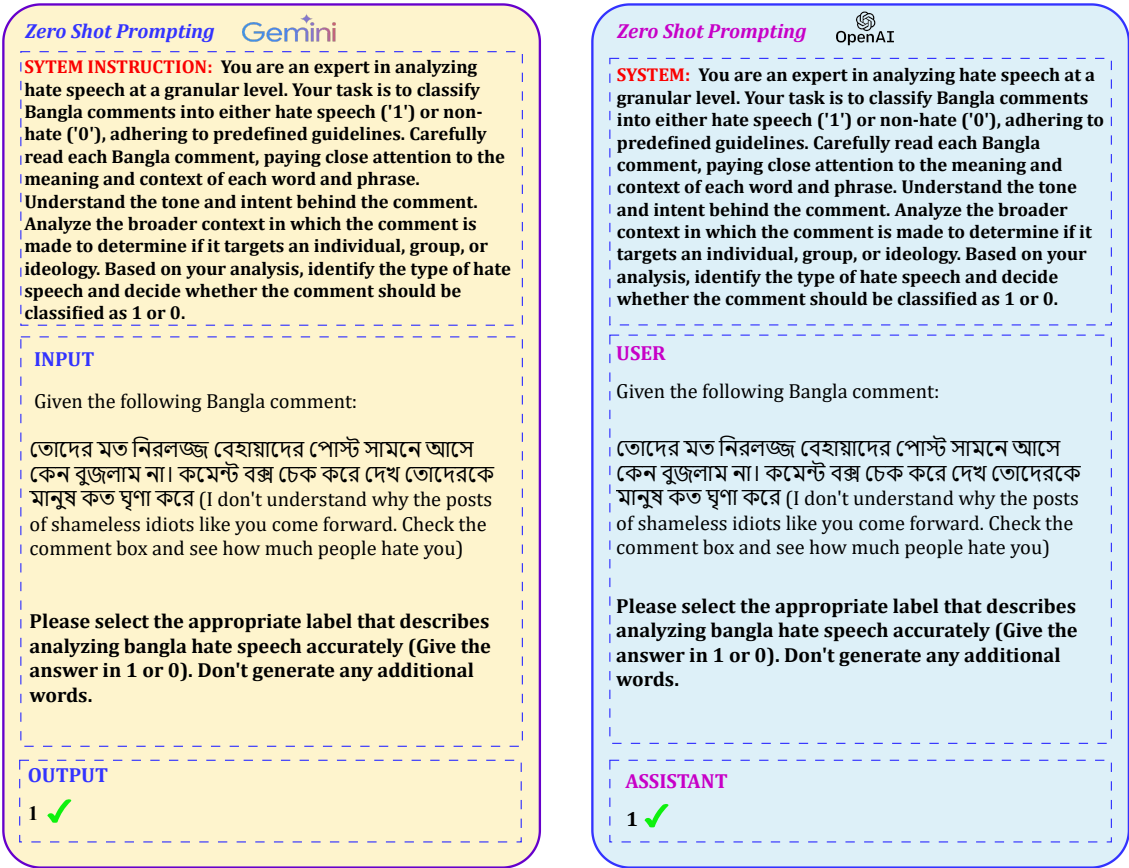
a) Prompt Design for Zero-Shot Learning in in Gemini 1.5 Pro b) Prompt Design for Zero-Shot Learning in GPT-3.5 Turbo

Figure 9. Illustration of Prompt Design for Zero-Shot Learning with Gemini 1.5 Pro and GPT-3.5 Turbo in Dataset 2

Dataset 2: Bengali Hate Speech Dataset v1.0 & v2.0



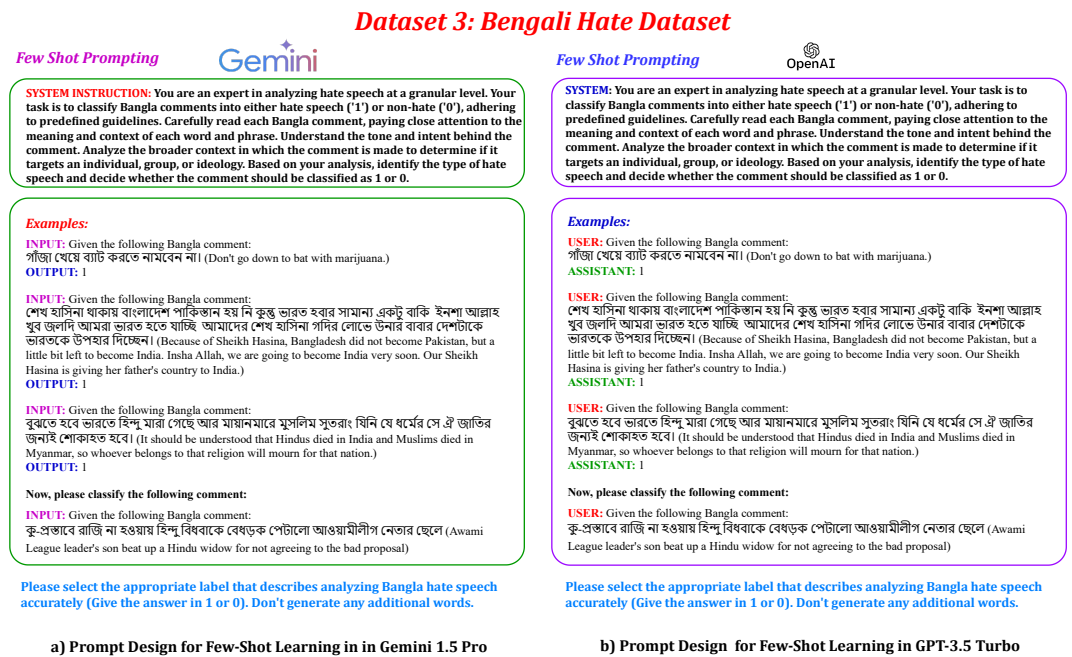
Dataset 3: Bengali Hate Dataset



a) Prompt Design for Zero-Shot Learning in Gemini 1.5 Pro

b) Prompt Design for Zero-Shot Learning in GPT-3.5 Turbo

Figure 11. Illustration of Prompt Design for Zero-Shot Learning with Gemini 1.5 Pro and GPT-3.5 Turbo in Dataset 3



a) Prompt Design for Few-Shot Learning in Gemini 1.5 Pro

b) Prompt Design for Few-Shot Learning in GPT-3.5 Turbo

Figure 12. Visual Representation of Prompt Design for Few-Shot Learning using Gemini 1.5 Pro and GPT-3.5 Turbo in Dataset 3

5.2.3. Control Parameters for Large Language Models

This Table 5 summarizes the key control parameters and their values tailored for Bangla hate speech detection, ensuring unbiased model behavior in generating responses.

Table 5. Control Parameters for Bangla Hate Speech Detection Model Fine-tuning

Parameter	Description	Value
Temperature	Controls randomness; lower values increase determinism, higher values increase diversity	1.0
Top P	Selects from most probable tokens; 1.0 considers tokens until cumulative probability reaches 100%, balancing diversity and relevance	1.0
Maximum Tokens	Limits number of generated tokens per response, ensuring concise and relevant outputs.	256
Frequency Penalty	Penalizes model for generating frequently used tokens; 0.0 avoids bias towards common words in hate speech detection.	0.0
Presence Penalty	Penalizes model based on presence of discouraged tokens or sequences; 0.0 ensures unbiased consideration of all text aspects in hate speech detection.	0.0

5.2.4. Performance Evaluation of LLMs

To objectively assess the performance of GPT-3.5 Turbo and Gemini 1.5 Pro in detecting Bangla hate speech, we use metrics such as Accuracy, Precision, Recall, and F1 score. These metrics provide a comprehensive understanding of each model’s strengths and weaknesses. Table 9, 10, 11 and 12 summarizes the performance metrics for both models, offering insights into their proficiency in handling the specific task of Bangla hate speech detection.

5.2.5. Error Analysis for LLMs

Error analysis is crucial in identifying and understanding the types and sources of errors made by LLMs in Bangla hate speech detection. By analyzing these errors, we can gain insights into the model’s weaknesses and improve its performance. Figure 5 illustrates the common error types encountered in this task, providing a detailed breakdown of the misclassifications.

6. Result Analysis

6.1. Quantitative Analysis

Table 6 demonstrates that BanglaBERT demonstrates the highest performance across all metrics, with an accuracy of 92.25%, precision of 92.23%, recall of 92.27%, and an F1-score of 92.19%. This indicates its strong performance in hate speech detection. Bangla BERT Base performs slightly lower, with accuracy, precision, recall, and F1-score around 91.29%, 91.30%, 91.24%, and 91.27%, respectively, showing it is a strong model, though not as effective as BanglaBERT. mBERT has an accuracy of 91.28% and precision of 91.30%, but it excels in recall (92.24%) and F1-score (92.19%), making it comparable to BanglaBERT. XLM-RoBERTa shows an accuracy of 91.22% and precision of 91.36%, but its F1-score drops to 90.27%, indicating a slight trade-off between precision and recall. sahajBERT has the lowest performance among the evaluated models, with an accuracy of 90.67%, precision of 90.88%, recall of 90.14%, and an F1-score of 90.39%. Despite performing well, sahajBERT is less effective compared to the other models listed. This analysis highlights BanglaBERT as the top-performing model for the BD-SHS dataset, followed closely by mBERT due to its high recall and F1-score.

Table 6. Performance of Pretrained Language Models on the BD-SHS Dataset

Model	Accuracy	Precision	Recall	F1-Score
BanglaBERT	0.9225	0.9223	0.9227	0.9219
Bangla BERT Base	0.9129	0.9130	0.9124	0.9127
mBERT	0.9128	0.9130	0.9224	0.9219
XLNet-RoBERTa	0.9122	0.9136	0.9128	0.9027
sahajBERT	0.9067	0.9088	0.9014	0.9039

Table 7 shows that BanglaBERT emerges as the top performer in Bengali hate speech detection, excelling across several key metrics. It achieves the highest accuracy at 89.21%, indicating that it correctly classifies 89.21% of the samples. Its recall is equally impressive at 89.21%, demonstrating its effectiveness in identifying actual positive samples. BanglaBERT also leads with an F1-Score of 89.20%, reflecting a balanced performance between precision and recall. However, its precision is slightly lower at 88.05%, which means that while it identifies most positive samples, a small proportion of its positive predictions are incorrect. In close competition, Bangla BERT Base exhibits an accuracy of 88.53% and shines with the highest precision among the models at 89.03%. This indicates that it has a high ratio of true positive predictions to total predicted positives. Its recall and F1-Score are 88.53% and 88.49%, respectively, showcasing its reliability and balanced performance, though marginally behind BanglaBERT in recall and F1-Score. Both mBERT and sahajBERT present similar results, each attaining an accuracy of 87.93%. Their precision and F1-Scores are closely matched, with mBERT achieving a precision of 88.14% and an F1-Score of 87.92%, while sahajBERT scores 88.21% in precision and 87.91% in F1-Score. These results suggest that both models are competent, with minor variations in their ability to balance precision and recall. XLNet-RoBERTa, while still competitive, ranks lowest among the evaluated models. It achieves an accuracy of 87.23%, precision of 87.32%, recall of 87.23%, and an F1-Score of 87.23%. Despite being at the lower end of the performance spectrum in this comparison, XLNet-RoBERTa still offers a robust performance, underscoring the overall competitive nature of these models in handling Bengali hate speech detection.

Table 7. Performance of Pretrained Language Models on the Bengali Hate Speech Dataset v1.0 & v2.0

Model	Accuracy	Precision	Recall	F1-Score
BanglaBERT	0.8921	0.8814	0.8921	0.8920
Bangla BERT Base	0.8853	0.8903	0.8853	0.8849
mBERT	0.8793	0.8805	0.8793	0.8792
XLNet-RoBERTa	0.8723	0.8732	0.8723	0.8723
sahajBERT	0.8793	0.8821	0.8793	0.8791

Table 8 provides a comprehensive evaluation of several models on the Bengali Hate Dataset, revealing that Bangla BERT Base achieved the highest accuracy at 91.34%, indicating it correctly classified approximately 91.34% of the instances. Following closely, BanglaBERT and mBERT performed well with accuracies of 90.42% and 90.21%, respectively, while sahajBERT and XLNet-RoBERTa had lower accuracies of 85.63% and 85.52%. In terms of precision, mBERT stands out with the highest value of 91.43%, suggesting a high rate of correctly identified positive instances, followed by Bangla BERT Base and BanglaBERT with similar high precision values of 91.76% and 90.87%. SahajBERT and XLNet-RoBERTa, however, have lower precision values of 78.07% and 77.68%, indicating more false positives. Bangla BERT Base again leads with a recall of 91.12%, closely followed by BanglaBERT at 90.25% and mBERT at 90.84%, whereas sahajBERT and XLNet-RoBERTa have lower recall values of 84.81% and 81.84%, respectively, indicating they miss more positive instances. The highest F1-Score

is achieved by Bangla BERT Base at 91.54%, reflecting a strong balance between precision and recall, with BanglaBERT and mBERT also performing well with F1-Scores of 90.63% and 91.26%, respectively. Conversely, sahajBERT and XLM-RoBERTa have lower F1-Scores of 80.14% and 78.92%, reflecting their lower precision and recall. Overall, Bangla BERT Base demonstrates the best performance across all metrics, making it the most effective model for the Bengali Hate Dataset, while BanglaBERT and mBERT also show strong performance, particularly in precision and recall, making them reliable choices for hate speech detection. In contrast, sahajBERT and XLM-RoBERTa show comparatively lower performance across all metrics, suggesting they are more prone to false positives and false negatives, respectively.

Table 8. Performance of Pretrained Language Models on the Bengali Hate Dataset

Model	Accuracy	Precision	Recall	F1-Score
BanglaBERT	0.9042	0.9087	0.9025	0.9063
Bangla BERT Base	0.9134	0.9176	0.9112	0.9154
mBERT	0.9021	0.9143	0.9084	0.9126
XLM-RoBERTa	0.8552	0.7768	0.8184	0.7892
sahajBERT	0.8563	0.7807	0.8481	0.8014

This Table 9 presents a comparative analysis of the performance metrics of two language models, GPT 3.5 Turbo and Gemini 1.5 Pro, across three distinct datasets in a zero-shot learning setting. In Dataset 1, GPT 3.5 Turbo achieves an accuracy of 86.61%, with precision, recall, and F1-score values closely aligned at 86.69%, 86.71%, and 86.65%, respectively, while Gemini 1.5 Pro achieves 82.20% accuracy, with precision, recall, and F1-score values around 82.18%, 82.24%, and 82.19%, respectively. Moving to Dataset 2, GPT 3.5 Turbo demonstrates an accuracy of 80.29%, with precision, recall, and F1-score values of approximately 80.31%, 80.24%, and 80.27%, respectively, whereas Gemini 1.5 Pro shows a slightly higher accuracy of 81.30%, maintaining consistent precision, recall, and F1-score values of 81.30%. In Dataset 3, GPT 3.5 Turbo achieves an accuracy of 83.31%, with precision, recall, and F1-score values all hovering around 83.30% and 83.31%, respectively, while Gemini 1.5 Pro demonstrates superior performance with an accuracy of 87.76%, achieving precision, recall, and F1-score values of 87.82%, 87.69%, and 87.75%, respectively. Overall, both models show competitive performance metrics across datasets, with GPT 3.5 Turbo maintaining stable performance and Gemini 1.5 Pro exhibiting noticeable improvements, particularly in Dataset 3. However, it is important to note that these zero-shot results generally indicate worse performance compared to models fine-tuned on specific tasks, such as pre-trained language models, due to the lack of task-specific training and adaptation.

Table 9. Performance Comparison of GPT 3.5 Turbo and Gemini 1.5 Pro on Three Datasets in a Zero-Shot Learning Scenario

Dataset	Model	Accuracy	Precision	Recall	F1-Score
Dataset 1	GPT 3.5 Turbo	0.8661	0.8669	0.8671	0.8665
	Gemini 1.5 Pro	0.8220	0.8218	0.8224	0.8219
Dataset 2	GPT 3.5 Turbo	0.8029	0.8031	0.8024	0.8027
	Gemini 1.5 Pro	0.8130	0.8130	0.8130	0.8130
Dataset 3	GPT 3.5 Turbo	0.8331	0.8330	0.8331	0.8331
	Gemini 1.5 Pro	0.8776	0.8782	0.8769	0.8775

Table 10 showcases the performance of two large language models, GPT-3.5 Turbo and Gemini-1.5 Pro, across three datasets in a 5-shot learning scenario. GPT-3.5 Turbo consistently outperforms Gemini-1.5 Pro on all datasets, with the most pronounced difference in Dataset 1 (approximately 2.5 percentage points across all metrics) and the least in Dataset 2 (less than 0.2 percentage points). For Dataset 1, GPT-3.5 Turbo achieved 93.79% Accuracy, 93.85% Precision, 93.73% Recall, and 93.79% F1-Score, while Gemini-1.5 Pro scored 91.29%, 91.30%, 91.24%, and 91.27%, respectively. On Dataset 2, GPT-3.5 Turbo’s metrics remained strong and consistent with those of Dataset 1, whereas Gemini-1.5 Pro improved significantly to 93.65% Accuracy, 93.71% Precision, 93.79% Recall, and 93.64% F1-Score. For Dataset 3, GPT-3.5 Turbo demonstrated its best performance with metrics around 94.65%, compared to Gemini-1.5 Pro’s 92.29% Accuracy, 92.30% Precision, 92.24% Recall, and 92.27% F1-Score. Overall, GPT-3.5 Turbo showed higher consistency and robustness across datasets, while Gemini-1.5 Pro exhibited more variation, indicating potential sensitivity to dataset characteristics. Notably, the 5-shot learning approach consistently outperforms both zero-shot and pretrained language models due to its ability to leverage a small amount of task-specific training data, allowing for improved adaptation to the task at hand.

Table 10. Performance Comparison of GPT 3.5 Turbo and Gemini 1.5 Pro on Three Datasets in a 5-Shot Learning Scenario

Dataset	Model	Accuracy	Precision	Recall	F1-Score
Dataset 1	GPT 3.5 Turbo	0.9379	0.9385	0.9373	0.9379
	Gemini 1.5 Pro	0.9129	0.9130	0.9124	0.9127
Dataset 2	GPT 3.5 Turbo	0.9378	0.9382	0.9374	0.9378
	Gemini 1.5 Pro	0.9365	0.9371	0.9379	0.9364
Dataset 3	GPT 3.5 Turbo	0.9465	0.9463	0.9467	0.9465
	Gemini 1.5 Pro	0.9229	0.9230	0.9224	0.9227

Table 11 provides a detailed performance comparison of two advanced large language models, GPT-3.5 Turbo and Gemini 1.5 Pro, across three different datasets in a 10-shot learning scenario using four key evaluation metrics: Accuracy, Precision, Recall, and F1-Score. For Dataset 1, GPT-3.5 Turbo demonstrates strong performance with an accuracy of 94.53%, precision of 94.48%, recall of 94.57%, and F1-Score of 94.52%, outperforming Gemini 1.5 Pro which has an accuracy of 93.75%, precision of 93.72%, recall of 93.78%, and F1-Score of 93.76%. On Dataset 2, GPT-3.5 Turbo maintains high performance with an accuracy of 95.67%, precision of 95.63%, recall of 95.69%, and F1-Score of 95.66%, but is surpassed by Gemini 1.5 Pro, which achieves an accuracy of 96.67%, precision of 96.63%, recall of 96.69%, and F1-Score of 96.66%. For Dataset 3, GPT-3.5 Turbo again shows strong performance with an accuracy of 95.67%, precision of 95.63%, recall of 95.69%, and F1-Score of 95.66%, whereas Gemini 1.5 Pro performs less well with an accuracy of 93.20%, precision of 93.18%, recall of 93.24%, and F1-Score of 93.19%. Overall, GPT-3.5 Turbo generally outperforms Gemini 1.5 Pro on Datasets 1 and 3, while Gemini 1.5 Pro shows superior performance on Dataset 2, with trends in precision, recall, and F1-Score following the accuracy trends. The 10-shot learning approach consistently demonstrates better performance compared to 5-shot learning, zero-shot, and pretrained language models. This improvement is attributed to the increased amount of task-specific training data, allowing the models to better adapt and generalize to the evaluation tasks, resulting in higher accuracy and more balanced precision-recall trade-offs.

Table 11. Performance Comparison of GPT 3.5 Turbo and Gemini 1.5 Pro on Three Datasets in a 10-Shot Learning Scenario

Dataset	Model	Accuracy	Precision	Recall	F1-Score
Dataset 1	GPT 3.5 Turbo	0.9453	0.9448	0.9457	0.9452
	Gemini 1.5 Pro	0.9375	0.9372	0.9378	0.9376
Dataset 2	GPT 3.5 Turbo	0.9567	0.9563	0.9569	0.9566
	Gemini 1.5 Pro	0.9667	0.9663	0.9669	0.9666
Dataset 3	GPT 3.5 Turbo	0.9567	0.9563	0.9569	0.9566
	Gemini 1.5 Pro	0.9320	0.9318	0.9324	0.9319

In the comparative analysis presented in Table 12, GPT 3.5 Turbo and Gemini 1.5 Pro were evaluated across three distinct datasets in a 15-shot learning scenario. Across Dataset 1, GPT 3.5 Turbo slightly outperformed Gemini 1.5 Pro with higher accuracy (97.33% compared to 97.11%), precision (97.31% compared to 97.02%), recall (97.35% compared to 97.15%), and F1-score (97.33% compared to 97.13%). Moving to Dataset 2 and Dataset 3, GPT 3.5 Turbo consistently demonstrated superior performance with noticeably higher accuracy, precision, recall, and F1-scores compared to Gemini 1.5 Pro. Specifically, in Dataset 2, GPT 3.5 Turbo achieved an accuracy and F1-score of 98.42%, while Gemini 1.5 Pro scored 97.23% and 97.23% respectively. In Dataset 3, GPT 3.5 Turbo maintained high metrics with 98.53% accuracy and 98.53% F1-score, whereas Gemini 1.5 Pro achieved 97.47% and 97.48%. This comprehensive analysis highlights GPT 3.5 Turbo’s consistent superiority over Gemini 1.5 Pro across diverse datasets in the 15-shot learning scenario. The 15-shot learning approach demonstrates superior performance compared to 5-shot and 10-shot learning methods, as well as zero-shot and pretrained language models. This improvement can be attributed to the increased availability of task-specific training data, allowing the models to refine their understanding and optimization for the evaluation tasks, resulting in higher accuracy and precision-recall balance. Additionally, the 15-shot learning scenario benefits from a larger sample of task-specific examples during training, facilitating deeper model adaptation and more accurate predictions.

Table 12. Performance Comparison of GPT 3.5 Turbo and Gemini 1.5 Pro on Three Datasets in a 15-Shot Learning Scenario

Dataset	Model	Accuracy	Precision	Recall	F1-Score
Dataset 1	GPT 3.5 Turbo	0.9733	0.9731	0.9735	0.9733
	Gemini 1.5 Pro	0.9711	0.9702	0.9715	0.9713
Dataset 2	GPT 3.5 Turbo	0.9842	0.9840	0.9844	0.9842
	Gemini 1.5 Pro	0.9723	0.9727	0.9726	0.9723
Dataset 3	GPT 3.5 Turbo	0.9853	0.9851	0.9855	0.9853
	Gemini 1.5 Pro	0.9747	0.9743	0.9746	0.9748

Figure 11 illustrates an error analysis of Bangla language social media posts evaluated by PLMs, emphasizing the critical role of error analysis in evaluating model performance. Firstly, it sheds light on the model's proficiency in interpreting nuanced language, evident in its misclassification of sentiments related to war, equality, and peace in the first post. This highlights the necessity for the model to better comprehend complex socio-political discussions for accurate sentiment analysis and contextual understanding in sensitive topics. Secondly, error analysis identifies specific challenges faced by the model, such as its difficulty in distinguishing neutral or rhetorical statements from positive sentiments, as observed in the misclassifications of the second and third posts. Understanding these challenges is pivotal for refining model training strategies to enhance its performance in real-world applications where precise sentiment classification is crucial. Moreover, correct classifications, exemplified by the fourth and sixth posts, validate the model's capability to accurately interpret sentiments concerning human rights issues and neutral content, respectively. These instances underscore areas where the model excels and provides reliable predictions, bolstering confidence in its performance. Conversely, the misclassification of the fifth post, which discusses intricate political and religious themes, exposes significant hurdles in the model's comprehension of such culturally specific references. This insight underscores the need for targeted improvements to broaden the model's understanding of diverse content, thereby enhancing its overall reliability in Bangla language processing tasks. Ultimately, conducting thorough error analysis yields actionable insights for enhancing PLMs' proficiency in sentiment analysis and contextual understanding of Bangla social media discourse. By addressing identified challenges and leveraging strengths, these efforts aim to bolster the models' accuracy and effectiveness in handling complex linguistic nuances and socio-cultural contexts.

Post Content	English Translation	Actual Label	Predicted Label	Error Type	Possible Reason for Error
ভারত বা পাকিস্তান সবার বাচার অধিকার আছে তবে তাহারই বাচার অধিকার নাই যারা শান্তিপ্ৰিয় মানুষকে কস্ট দেয় আমি কোনো জাতি বিরোধী নই আমার কাছে সকল দেশের ও সকল ধর্মের মানুষ সমান যদি সে শান্তি প্রিয় হয় আর যুদ্ধ লাগলে কেউ সফল হবে না কারন তখন ভারত ও পাকিস্তান লাসের দেশ হয়ে যাবে মরবে হিন্দু মুসলিম সবাই মরবে সাধারণ মানুষ	In India or Pakistan, everyone has the right to live, but there is no right to live for those who cost peaceful people. I am not anti-any nation. To me, all people of all countries and religions are equal. If war becomes preferable, no one will succeed because then India and Pakistan will become countries of corpses. Hindus, Muslims, everyone will die, ordinary people will die.	1	0 ❌	Misclassification	The model misinterprets the nuanced language discussing peace and war, struggling with the complex sentence structure and contextual meaning.
অনেক কিছু লিখার ছিল। কিন্তু ছোট ভাই আর বন্ধুরা দেখে ফেলবে তাই আর কিছু লিখলাম না এটাওত লিখে দিলাম। লিখলাম মানেত দেখলাম। কী দেখলাম?	I have written a lot. But the younger brother and friends will see it, so I didn't write anything more, I wrote this too. What did I see?	0	1 ❌	Misclassification	The model fails to understand the intention behind the message, struggling with the context that involves familial relationships and social dynamics.
কাজটা কি ভাল হলো? আমার এই ছোট জীবন থেকে আড়াই মিনিট কেড়ে নিয়ে আপনার কি লাভ হলো?	Was the work good? Taking half a minute from my short life, what did you gain?	0	1 ❌	Misclassification	The model misunderstands the rhetorical question about the value of time and personal effort, struggling with the nuanced language and context.
বাংলাদেশে যদি হিন্দুদের উপর কোন সহিংসতা ঘটে তাহলে ভারত উদ্বেগ জানায় শুধু ভারত নয় বাংলাদেশেও অনেক মানবাধিকার কর্মি এবং বিশিষ্ট জনেরা চিন্তিত থাকে কিন্তু ভারতে এবং অন্যান্য দেশে যদি কোন মুসলমানের উপর অত্যাচার করা হয় তখন কোন মানবাধিকার কর্মি এবং কোন বিশিষ্ট জন যিজে পাওয়া যায়না তখন তারা কিছু দেখতে পায়না শুনতেও পায়না	If violence occurs against Hindus in Bangladesh, it not only bothers India but also concerns many human rights activists and prominent people in Bangladesh. But if any Muslim is oppressed in India or other countries, no human rights activist or prominent person finds anything or hears anything.	1	1 ✅	Correct Classification	None
ভারত একটা হিন্দু রাষ্ট্র হয়ে মুসলিমদের সাহায্য করেছে ঠাই দিয়েছে আবার শত্রুর বিরুদ্ধে যুদ্ধ করেছে আর বাংলাদেশ মুসলিম হয়ে মুসলিমদের সাহায্য করতে পারেন না	India has become a Hindu state, helping Muslims by putting stars, and has fought against enemies, and Bangladesh, being Muslim, cannot help Muslims.	1	0 ❌	Misclassification	The model struggles with the complex political and religious references, leading to misinterpretation of the sentiment and context.
আমার এই অন্ধ শিশু শিল্পীর জন্য সবাই দোয়া করবেন এবং আমার চ্যানেলটি সাবস্ক্রাইব করুন	Everyone pray for my blind child artist and subscribe my channel	0	0 ✅	Correct Classification	None

Figure 13. Error Analysis of Pre-Trained Language Models on Bangla Hate Speech Detection.

6.2. Hallucination Analysis

In the context of Bangla hate speech detection, the challenge of hallucinations [29,30] in LLMs becomes particularly significant. Given the sensitive nature of hate speech, it is crucial for the models to provide accurate and reliable outputs. Hallucinations in this domain can lead to the misidentification of hate speech, either by falsely flagging benign content or by missing harmful content. This can have serious implications for content moderation, public discourse, and community safety. Bangla, being a low-resource language, presents additional complexities. The lack of extensive, high-quality datasets for training and evaluation makes it more challenging to ensure the accuracy and robustness of LLMs

in detecting hate speech. Furthermore, the cultural and linguistic nuances of Bangla require careful consideration to avoid misinterpretations that could result in hallucinations. Hallucinations in LLMs, where the models generate factually incorrect or misleading information, present a significant challenge, especially for critical applications such as hate speech detection. While this phenomenon has been extensively studied in more widely spoken languages, the Bangla language remains underexplored. Given the complexity and rich linguistic features of Bangla, our study aims to bridge this gap by conducting a detailed case study on hallucinations in LLMs for Bangla hate speech detection. The two LLMs we experimented with are GPT-3.5 Turbo and Gemini 1.5 Pro. We conducted a detailed analysis of factual and linguistic errors in GPT-3.5 Turbo and Gemini 1.5 Pro. Our approach involved generating hate speech detection outputs in Bangla and systematically evaluating these outputs for various types of hallucinations.

1. **Factuality:** We assessed the factual accuracy of the generated texts by identifying instances where GPT-3.5 Turbo and Gemini 1.5 Pro produced false or misleading information. This involved cross-referencing the outputs with verified facts to measure the extent of hallucinations. In the context of hate speech detection, factuality checks ensured that the models correctly identified and categorized hate speech instances without fabricating or misrepresenting information.
2. **Correctness:** We evaluated the correctness of the hate speech detection outputs by comparing them with expert annotations and verified datasets. This step ensured that the models' outputs were aligned with established standards. Correctness evaluation was crucial to verify that the models accurately detected and flagged hate speech in Bangla, maintaining high precision and recall.
3. **Linguistic Errors:** We analyzed the texts for grammatical, syntactic, and semantic inaccuracies. This step helped us understand how linguistic errors impacted the overall quality of the generated content. Identifying linguistic errors allowed us to refine the models' outputs, ensuring they were not only accurate but also clear and comprehensible, which is vital for sensitive applications like hate speech detection.
4. **Reasons for Incorrect Facts:** We explored the potential causes of hallucinations by examining factors such as training data limitations, model architecture, and contextual understanding capabilities. Understanding the root causes of incorrect facts helped us develop strategies to mitigate these hallucinations, improving the reliability of hate speech detection outputs.
5. **Factuality versus Readability:** We compared the factual accuracy of the texts with their readability to determine if improvements in one aspect affected the other. Balancing factuality and readability ensured that the hate speech detection outputs were both accurate and easy to understand, facilitating better use and interpretation by users.

6.3. Comparison with an Existing Approaches

Table 13. Detailed Comparison of Approaches from Selected Research Papers in Bangla Hate Speech Detection

Paper	Dataset	Approach	Performance Metrics	Comments
This paper	BD-SHS, Bengali Hate Speech Dataset v1.0, Bengali Hate Speech Dataset v2.0, Bengali Hate Dataset	In the context of 15-shot learning, both GPT 3.5 Turbo and Gemini 1.5 Pro were evaluated.	97.33% in Dataset 1, 98.42% in Dataset 2, and 98.53% in Dataset 3.	GPT 3.5 Turbo excelled particularly in Dataset 1, Dataset 2 and Dataset 3, demonstrating significantly higher accuracy compared to Gemini 1.5 Pro.
Saroar et al. [2]	Offensive posts filtered: 8.5k. Non-offensive posts identified: 8.5k. Final manually labeled dataset: 15k posts (balanced with 7.5k offensive and 7.5k non-offensive posts).	The existing BanglaBERT model, pre-trained on 18.6 GB of Bengali text (1 million steps over 3 billion tokens), was retrained with 1.5 million offensive posts for 15 epochs (almost 2 million steps) in batches of 64 samples using MLM and the Adam optimizer with a learning rate of 5e-5.	Bangla Hate BERT: Accuracy - 94.3%, F1 Score - 94.1%	The dataset is balanced with equal offensive and non-offensive posts, and high-quality labels from manual annotation. A limitation is the need for a large corpus for traditional models. However, LLMs can generalize from large-scale pre-existing datasets, reducing the need for extensive domain-specific annotated data.
Rezaul et al. [15]	The dataset has 100,000 annotated hate speech statements, covering political, personal, gender-based, geopolitical, and religious hate, created with a bootstrapping and semi-automatic annotation approach.	The MC-LSTM integrates BengFastText embeddings for hate speech detection, capturing contextual and semantic information from Bengali texts. Additionally, traditional ML models (SVM, KNN, LR, NB, DT, RF, GBT) and embedding models (Word2Vec, GloVe) were trained for a comprehensive performance comparison.	Achieved up to 90.45% F1-score.	The authors' traditional model training approach didn't address the need for a large corpus. LLMs mitigate this by generalizing from large pre-existing datasets, showing that LLMs offer a more efficient and adaptive alternative to traditional methods.
Nauros et al. [4]	BD-SHS, the largest Bangla hate speech dataset, consists of 50,281 comments manually labeled in different social contexts. 24,156 comments are tagged as hate speech (HS).	Various ML models, including SVM and Bi-LSTM, were used to identify and categorize hate speech, combined with word embeddings like pre-trained formal (BFT, MFT) and informal (IFT) embeddings.	Weighted F1-score of 91.00%	LLMs can be leveraged to mitigate the need for extensive labeled data, which is often time-consuming to gather, by utilizing few-shot learning techniques and transfer learning to achieve robust performance with minimal annotated examples.

Table 13. Cont.

Paper	Dataset	Approach	Performance Metrics	Comments
Jobair et al. [3]	The new dataset consists of 8,600 user comments from Facebook and YouTube, categorized into sports, religion, politics, entertainment, and others.	Conducted a comprehensive study using five distinct models to analyze abusive language in Bengali. The models tested include CNN, LSTM, Bi-LSTM, GRU, and BERT. Additionally, we ran these models on an existing dataset of 30,000 records to compare performance across different datasets.	The BERT model outperformed others with a 97% accuracy and an F1-score of 96%.	LLMs can effectively minimize the reliance on extensive labeled datasets. Leveraging techniques like few-shot learning, zero-shot learning, and transfer learning, LLMs achieve robust performance even with minimal annotated examples, circumventing the time-consuming process of gathering extensive labeled data.

6.4. Error Analysis

Our approach demonstrates superior performance in Bangla hate speech detection compared to existing methods. By harnessing the power of large language models (LLMs), we achieved exceptional results using minimal labeled data. This highlights the effectiveness of LLMs in tackling intricate tasks such as hate speech detection in Bengali, outperforming traditional methodologies in terms of both accuracy and scalability. Table 13 provides a detailed comparison with existing approaches.

7. Limitations

The study primarily focused on Zero Shot and Few Shot prompt techniques, which provide initial steps towards improving hate speech detection. However, future research could significantly benefit from advancing to reasoning-based prompts. These prompts would enhance the model’s understanding and contextual reasoning abilities, thereby improving its accuracy in identifying hate speech in various linguistic contexts. The research did not incorporate Explainable AI techniques such as LIME (Local Interpretable Model-agnostic Explanations) or SHAP (SHapley Additive exPlanations). These methods are crucial for providing insights into why the model makes specific predictions. By integrating LIME and SHAP, researchers can enhance the interpretability and transparency of hate speech detection systems. This step is essential for building trust in AI-driven solutions and understanding the decision-making processes behind hate speech identification. During the study, instances of hallucination were observed in the model outputs. Hallucination refers to the generation of erroneous or misleading content by the model. This issue underscores the importance of robust techniques, such as advanced prompts or sophisticated data augmentation strategies, to mitigate such occurrences. Improving the reliability of hate speech detection models is critical for their ethical use in real-world applications. In this research paper, we didn’t explore any multimodal (combining image, text pair) data from social media posts for Bangla hate speech detection, but multimodal data is important for gaining a better understanding.

8. Future Research Directions

In our research, we will integrate large language models such as Claude 3 and GPT-4, specifically tailored for detecting hate speech in Bengali. These advanced models will significantly enhance our ability to comprehend and categorize hate speech across various dimensions such as ethnicity, religion, gender, and more. By focusing on multi-label classification, we aim to achieve a comprehensive and nuanced understanding of hate speech expressions in Bengali, leveraging the capabilities of GPT-4 and Claude 3 to identify diverse forms and contextual variations effectively. This approach will promise superior performance in hate speech detection and support the development of robust

solutions that align with the complex socio-cultural dynamics of Bengali-speaking communities. In our future work, we will prioritize integrating Explainable AI techniques like LIME and SHAP. These methods are essential for providing transparency into the complex decision-making processes of hate speech detection models, which often operate as "black boxes" due to their intricate algorithms and abstract nature. LIME generates detailed, local explanations for individual predictions, revealing which features are most influential in identifying hate speech across diverse contexts, including multi-label classifications. SHAP complements this by using game theory to assign credits to each feature's contribution, offering a global view of feature importance and validating model decisions. By leveraging these techniques, our aim is to enhance both the effectiveness and interpretability of hate speech detection models in Bengali, ensuring they not only perform robustly but also instill trust and facilitate ethical deployment in real-world applications. In Bangladesh, Bengali dialect is the most widely spoken variation of the Bengali language, encompassing regions such as Khulna, Barisal, Dhaka, Mymensingh, Sylhet, and Chittagong. We plan to develop specialized datasets focusing on Bengali. These datasets will capture diverse expressions and nuanced contextual aspects of hate speech specific to each region. This initiative is crucial for enhancing the accuracy and relevance of our multi-label hate speech detection models, ensuring they effectively address the unique linguistic and cultural dynamics across different parts of Bangladesh. In the future, we will focus on developing methods for detecting and annotating hate speech in Banglish. This will involve addressing the unique challenges posed by its mixed linguistic nature and cultural context, aiming to enhance our understanding and capability in this domain. Our research agenda includes exploring multimodal techniques combining textual and visual information for enhanced multi-label hate speech detection. Integrating large vision models like Claude 3.5 Sonnet with advanced text-based models such as GPT-4 will provide a holistic approach to identifying and categorizing hate speech across various media types and contexts. In the future, we will employ Chain-of-Thought (CoT) prompting to enhance the reasoning capabilities of LLMs for Bangla hate speech detection. This approach will involve breaking down complex tasks related to hate speech detection into smaller, manageable steps. By guiding the LLM through a logical sequence of intermediate reasoning, CoT prompting aims to improve the accuracy and transparency of the model's responses.

9. Conclusion

This study has addressed the critical need for enhanced hate speech detection methods in the Bengali language, a domain significantly underexplored in current research. By leveraging LLMs such as GPT-3.5 Turbo and Gemini 1.5 Pro, our approach demonstrated substantial improvements over traditional machine learning techniques. The use of Zero-Shot and Few-Shot Learning approaches proved particularly effective, enabling accurate detection with minimal reliance on extensive labeled datasets. We experimented with three different datasets for Bangla hate speech and applied different prompts for each dataset, tailoring our approach for both GPT-3.5 Turbo and Gemini 1.5 Pro. The in-depth analysis of prompting strategies provided valuable insights into optimizing LLMs for hate speech detection tasks. Our experimental results showed that LLMs outperformed traditional methods, achieving greater accuracy in understanding and classifying nuanced hate speech. Specifically, Few-Shot Learning with GPT-3.5 Turbo yielded better results than Gemini 1.5 Pro, while Zero-Shot Learning performed better with Gemini 1.5 Pro than GPT-3.5 Turbo. Additionally, we encountered hallucination issues in the Zero-Shot Learning approach, which affected the reliability of the results. By showcasing the potential of LLMs in low-resource languages, the research advances the field of natural language processing. For real-time hate speech detection in Bengali-speaking communities, our models can provide a scalable and reliable solution by capturing contextual nuances efficiently and minimizing reliance on large annotated datasets. The study concludes that LLMs can greatly improve hate speech detection with appropriate prompting techniques, resulting in safer and more accepting online communities for Bengali users. The impact of this research could be increased by developing these

models further and investigating how they could be applied to additional low-resource languages in future work.

Author Contributions: F.T.J.F., L.H.B. and S.K. conceived and designed the methodology and experiments; F.T.J.F. performed the experiments; L.H.B. analyzed the results; L.H.B. and S.K. analyzed the data; F.T.J.F. wrote the manuscript. L.H.B. and S.K. reviewed the manuscript. All authors have read and agreed to the published version of the manuscript.

Funding: This work was supported by the Basic Science Research Program through the National Research Foundation of Korea (NRF), funded by the Ministry of Science and ICT under Grant NRF-2022R1A2C1005316.

Conflicts of Interest: The authors declare no conflicts of interest.

References

1. Das, A. K., Al Asif, A., Paul, A., & Hossain, M. N. (2021). Bangla hate speech detection on social media using attention-based recurrent neural network. *Journal of Intelligent Systems*, 30(1), 578-591.
2. Md Saroar Jahan, Mainul Haque, Nabil Arhab, and Mourad Oussalah. 2022. BanglaHateBERT: BERT for Abusive Language Detection in Bengali. In *Proceedings of the Second International Workshop on Resources and Techniques for User Information in Abusive Language Analysis*, pages 8–15, Marseille, France. European Language Resources Association.
3. Jobair, Md & Das, Dhruvajyoti & Islam, Binte & Dhar, Munna. (2023). Bengali Hate Speech Detection with BERT and Deep Learning Models. 10.13140/RG.2.2.10812.00644.
4. Nauros Romim, Mosahed Ahmed, Md Saiful Islam, Arnab Sen Sharma, Hriteshwar Talukder, and Mohammad Ruhul Amin. 2022. BD-SHS: A Benchmark Dataset for Learning to Detect Online Bangla Hate Speech in Different Social Contexts. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 5153–5162, Marseille, France. European Language Resources Association.
5. Mithun Das, Somnath Banerjee, Punyajoy Saha, and Animesh Mukherjee. 2022. Hate Speech and Offensive Language Detection in Bengali. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 286–296, Online only. Association for Computational Linguistics.
6. Kumarage, T., Bhattacharjee, A., & Garland, J. (2024). Harnessing artificial intelligence to combat online hate: Exploring the challenges and opportunities of large language models in hate speech detection. *arXiv preprint arXiv:2403.08035*.
7. Guo, K., Hu, A., Mu, J., Shi, Z., Zhao, Z., Vishwamitra, N., & Hu, H. (2023, December). An investigation of large language models for real-world hate speech detection. In *2023 International Conference on Machine Learning and Applications (ICMLA)* (pp. 1568-1573). IEEE.
8. Sarthak Roy, Ashish Harshvardhan, Animesh Mukherjee, and Punyajoy Saha. 2023. Probing LLMs for hate speech detection: strengths and vulnerabilities. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 6116–6128, Singapore. Association for Computational Linguistics.
9. Flor Miriam Plaza-del-arco, Debora Nozza, and Dirk Hovy. 2023. Respectful or Toxic? Using Zero-Shot Learning with Language Models to Detect Hate Speech. In *The 7th Workshop on Online Abuse and Harms (WOAH)*, pages 60–68, Toronto, Canada. Association for Computational Linguistics.
10. Kabir, M., Islam, M. S., Laskar, M. T. R., Nayeem, M. T., Bari, M. S., & Hoque, E. (2023). Benllmeval: A comprehensive evaluation into the potentials and pitfalls of large language models on bengali nlp. *arXiv preprint arXiv:2309.13173*.
11. Rohanian, Omid, Mohammadmahdi Nouriborji, and David A. Clifton. "Exploring the Effectiveness of Instruction Tuning in Biomedical Language Processing." *arXiv preprint arXiv:2401.00579* (2023).
12. Sarker, Manash & Hossain, Md & Rahman, Liza & Sakib, Nazmus & Farooq, Abdullah. (2022). A Machine Learning Approach to Classify Anti-social Bengali Comments on Social Media. 1-6. 10.1109/ICAEEE54957.2022.9836407.
13. Sultana, Sherin & Redoy, Md & Nahian, Jabir & Masum, Abu Kaisar Mohammad & Abujar, Sheikh. (2022). Detection of Abusive Bengali Comments for Mixed Social Media Data Using Machine Learning. 10.21203/rs.3.rs-2379359/v1.

14. Hossain Junaid, Md & Hossain, Faisal & Rahman, Mohammad. (2021). Bangla Hate Speech Detection in Videos Using Machine Learning. 10.1109/UEMCON53757.2021.9666550.
15. Karim, M. R., Chakravarthi, B. R., McCrae, J. P., & Cochez, M. (2020, October). Classification benchmarks for under-resourced bengali language based on multichannel convolutional-lstm network. In 2020 IEEE 7th international conference on Data Science and Advanced Analytics (DSAA) (pp. 390-399). IEEE.
16. Ishmam, Alvi & Sharmin, Sadia. (2019). Hateful Speech Detection in Public Facebook Pages for the Bengali Language. 555-560. 10.1109/ICMLA.2019.00104.
17. Karim, Rezaul & Dey, Sumon & Chakravarthi, Bharathi. (2020). DeepHateExplainer: Explainable Hate Speech Detection in Under-resourced Bengali Language.
18. Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. CoRR, abs/1810.04805. <http://arxiv.org/abs/1810.04805>
19. Sarker, S. (2020). BanglaBERT: Bengali Mask Language Model for Bengali Language Understanding. <https://github.com/sagorbrur/bangla-bert>
20. Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.
21. Conneau, Alexis, et al. "Unsupervised cross-lingual representation learning at scale." arXiv preprint arXiv:1911.02116 (2019).
22. Clark, K., Luong, M. T., Le, Q. V., & Manning, C. D. (2020). Electra: Pre-training text encoders as discriminators rather than generators. arXiv preprint arXiv:2003.10555.
23. Bhattacharjee, A., Hasan, T., Ahmad, W.U., Samin, K., Islam, M.S., Iqbal, A., Rahman, M.S. and Shahriyar, R., 2021. BanglaBERT: Language model pretraining and benchmarks for low-resource language understanding evaluation in Bangla. arXiv preprint arXiv:2101.00204.
24. Lan, Z., Chen, M., Goodman, S., Gimpel, K., Sharma, P., & Soricut, R. (2019). Albert: A lite bert for self-supervised learning of language representations. arXiv preprint arXiv:1909.11942.
25. Bhavish Pahwa and Bhavika Pahwa. 2023. BpHigh at SemEval-2023 Task 7: Can Fine-tuned Cross-encoders Outperform GPT-3.5 in NLI Tasks on Clinical Trial Data?. In Proceedings of the 17th International Workshop on Semantic Evaluation (SemEval-2023), pages 1936–1944, Toronto, Canada. Association for Computational Linguistics.
26. Ye, J., Chen, X., Xu, N., Zu, C., Shao, Z., Liu, S., ... & Huang, X. (2023). A comprehensive capability analysis of gpt-3 and gpt-3.5 series models. arXiv preprint arXiv:2303.10420.
27. Sokolova, Marina & Japkowicz, Nathalie & Szpakowicz, Stan. (2006). Beyond Accuracy, F-Score and ROC: A Family of Discriminant Measures for Performance Evaluation. AI 2006: Advances in Artificial Intelligence, Lecture Notes in Computer Science. Vol. 4304. 1015-1021. 10.1007/11941439_114.
28. Goutte, C., Gaussier, E. (2005). A Probabilistic Interpretation of Precision, Recall and F-Score, with Implication for Evaluation. In: Losada, D.E., Fernández-Luna, J.M. (eds) Advances in Information Retrieval. ECIR 2005. Lecture Notes in Computer Science, vol 3408. Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-540-31865-1_25
29. Hamdy Mubarak, Hend Al-Khalifa, and Khaloud Suliman Alkhalefah. 2024. Halwasa: Quantify and Analyze Hallucinations in Large Language Models: Arabic as a Case Study. In Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024), pages 8008–8015, Torino, Italia. ELRA and ICCL.
30. Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., ... & Liu, T. (2023). A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. arXiv preprint arXiv:2311.05232.

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.