

# Natural Language Processing for Hate Speech Detection: A Review

**Mohini Chakarverti**

Assistant Professor, Bennett University, Uttar Pradesh, India.

## How to cite this paper:

Mohini Chakarverti, "Natural Language Processing for Hate Speech Detection: A Review", IJIRE-V4I02-266-272.

Copyright © 2023 by author(s) and 5<sup>th</sup> Dimension Research Publication.  
This work is licensed under the Creative Commons Attribution International License (CC BY 4.0).  
<http://creativecommons.org/licenses/by/4.0/>

**Abstract:** The variations in hate speech laws between nations, it is generally accepted that hate speech includes expressions of hostility or disparagement of an individual or a group because of a characteristic shared by that group, such as race, colour, national origin, sex, disability, religion, or sexual orientation. In this context, researchers have populated big databases from a variety of sources, which supported field research. In many of these studies, the topic of hate speech in various non-English languages and online groups has also been covered. The natural language processing has various phases which include pre-processing, feature extraction and classification. In this paper various schemes which is related to natural language processing for the hate speech detection is reviewed and analysed in terms of certain parameters

**Key Word:** NLP, Hate Speech Recognition, Cyber bullying

## I.INTRODUCTION

In the age of social computing, interpersonal connection is more obvious, particularly on social media sites and in online forums. People all around the world now have the opportunity to express and share their opinions instantly and widely thanks to micro blogging tools. Motivated, on one hand, by the platform's simple access and anonymity. On the other side, this provided a favourable atmosphere for the dissemination of violent and damaging content due to the user's desire to dominate discussion, share defend beliefs or argumentation, and possibly some business motivations [1]. Despite the variations in hate speech laws between nations, it is generally accepted that hate speech includes expressions of hostility or disparagement of an individual or a group because of a characteristic shared by that group, such as race, colour, national origin, sex, disability, religion, or sexual orientation. The spread of hate speech online has picked up new momentum, posing ongoing challenges for both policymakers and the research community. This is due in part to the variation in national hate speech legislation, the difficulty of placing boundaries on the constantly changing cyberspace, the increased need for individuals and societal actors to express their opinions and counterattacks from opponents, and the delay in manual check by internet operators [2].

### 1.1 Automatic Hate Speech Detection

Several studies on automatic textual hate speech identification have been conducted recently thanks to advancements in natural language processing (NLP) technology. Several well-known contests, like SemEval-2019 and 2020 and GermEval-2018, have held numerous competitions in an effort to improve automated hate speech detection. In this context, researchers have populated big databases from a variety of sources, which supported field research. In many of these studies, the topic of hate speech in various non-English languages and online groups has also been covered [3]. This prompted researchers to compare and contrast various processing pipelines, including feature set selection, Machine Learning (ML) techniques, classification algorithms, and so on. Examples of these include Naive Bayes, Linear Regression, Convolution Neural Network (CNN), LSTM, and BERT deep learning architectures.

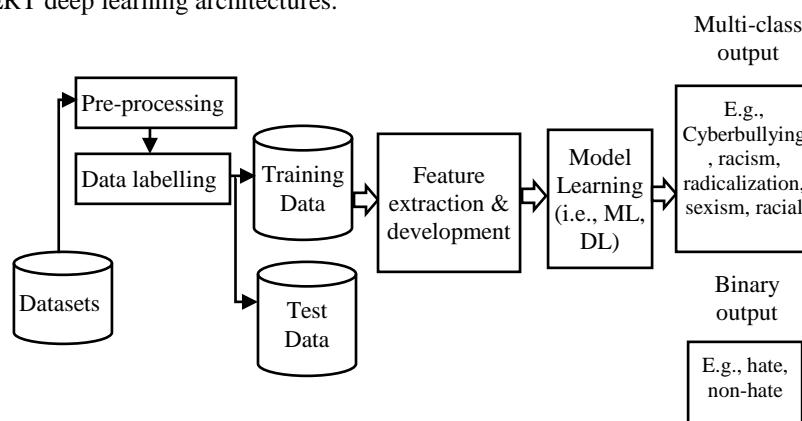


Figure 1: Typical automatic HS detection system pipeline

The general pipeline for the HS identification task, which is based on a text classification system, is shown in Figure 1. The following is a description of its primary elements:

- i. **Dataset collection and preparation:** The pipeline for HS detection starts with this stage. Social media sites (Facebook, Youtube, Twitter, etc.) are frequently used to collect datasets. Pre-processing is done in accordance with the structure and quality of the dataset [4]. Generally, this entails filtering and normalisation of textual inputs, which include, among other things, tokenization, stop word removal, misspelling correction, noise removal, stemming, and lemmatization. We'll also see that the dataset might be given to us right away, eliminating the need for collecting. The training and testing portions of the dataset should be separated during data preparation for the next machine learning stage.
- ii. **Feature Engineering:** The required characteristics are then retrieved from the textual inputs in the following phase of the analysis, transforming the unstructured text sequences into structured features. The TF-IDF, semantic, lexical, topic modelling, sentiment, BOW, and word embedding (FastText, GloVe, Word2Vec) are popular feature extraction methods. Dimensionality reduction is occasionally used to lessen the complexity in terms of time and memory. Principal component analysis (PCA) [5], linear discriminant analysis (LDA), non-negative matrix factorization (NMF), random projection, autoencoders, and t-distributed stochastic neighbour embedding (t-SNE) are a few examples of dimension reduction techniques.
- iii. **Model Training:** The training of a machine learning or deep learning model on the training dataset is one of the most important steps in the pipeline for text classification. Based on the needs of the task, a variety of classifiers, including RF, NB, LR, CNN, RNN, BERT, etc., can be modified. Typically, word embedding can be used in conjunction with another embedding layer in a neural network model to improve deep learning performance. The machine learning/deep learning model's output can be a binary decision (for example, hate speech versus non-hate speech) or a multi-class output where the model can distinguish between different types of hate speech and non-hate speech [6].

**Evaluation:** The performance of the machine learning/deep learning model is estimated in this last step of the text categorization pipeline. Accuracy, F1 score, precision, Matthews Correlation Coefficient (MCC), receiver operating characteristics (ROC), and area under the ROC curve are some of the evaluation measures utilised for this (AUC).

## 1.2 Natural Language Processing

The automatic processing of human languages is known as natural language processing, often known as computational linguistics. There are numerous definitions in use today, as NLP is a broad and heterogeneous science that is still relatively new. Natural language processing [7], in general, refers to a group of theoretically motivated computational techniques for analysing and portraying naturally occurring texts at one or more levels of linguistic analysis with the goal of achieving human-like language processing for a variety of tasks or applications. NLP is a branch of computer science and linguistics that studies how computers and human (natural) languages interact. Additionally, it is fueled by developments in machine learning (ML), which is a crucial component of artificial intelligence (AI).

The NLP approaches are designed in a way that allows the computer to comprehend orders given in natural language and act in accordance with them. It should be emphasised that spoken language and written language are the two categories into which natural language processing can be split. Given that speech accounts for the majority of human linguistic communication, written languages are generally [8] less important than speech in most contexts. Written language, however, is generally easier to understand than spoken language since spoken languages have to deal with a lot of background noise and ambiguity in the auditory stream. Language ambiguity makes natural language processing (NLP) a challenging topic in computer science. The "levels of language" approach [9] is the most illustrative way to describe what actually occurs within a Natural Language Processing system. People employ these levels to glean meaning from written or spoken languages. This is because language processing mostly relies on formal models or representation of knowledge connected to various levels. Additionally, by utilizing linguistic expertise, language processing applications set themselves apart from data processing systems. Four types of natural language processing techniques exist:

- i. **Symbolic techniques:** These methods carry out in-depth analyses of linguistic phenomena and are founded on the explicit representation of linguistic facts using well-known knowledge representation techniques and related algorithms.
- ii. **Statistical techniques:** These methods use a variety of mathematical techniques and frequently make use of large text corpora to create approximations of generalised models of linguistic phenomena based on examples of these phenomena provided [10] by the text corpora without the need for significant linguistic or outside-the-box knowledge.
- iii. **Connectionist approaches:** Similar to statistical techniques, these approaches create generalized models of linguistic phenomena. Connectionism, commonly referred to as "parallel distributed processing," "neural networks," or "neuro-computing," differs from other statistical techniques in that its models integrate statistical learning with a variety of representational theories. As a result, the connectionist representations enable the manipulation, inference, and change of logic rules.
- iv. **Hybrid approaches:** Hybrid approaches, also known as knowledge-driven and data-driven approaches, are being used by an increasing number of researchers [11]. It is clear from the foregoing that there are both similarities and distinctions between the methods. For instance, the assumptions, philosophical underpinnings, and proof sources vary between each strategy. Moreover, there are two groups into which the currently used text categorization techniques can be split: traditional techniques and deep learning techniques.

- Classical methods: These techniques combine statistical algorithms with manual feature engineering and guidelines. There are numerous methods for manually incorporating data instance features into feature vectors. Bag of words, word, and character n-grams are the most useful surface features for detecting hate speech, according to studies. The Support Vector Machine is the most often used algorithm in classifiers. For classification tasks, algorithms like Naive Bayes, Logistic Regression, and Random Forest are also employed [12].
- Deep learning methods: With the help of neural networks, these techniques automatically learn multiple layers of features from the input data. The term "deep learning" came into use in the early 2000s as a result of advancements in training methods and computer hardware that allowed for the training of progressively larger and deeper networks. Machine-learning methods that employ linear models and are trained over highly dimensional but sparse feature vectors have mostly dominated NLP techniques. Non-linear neural networks with dense inputs have, nevertheless, recently demonstrated success. The most widely employed networks are the Recurrent Neural Network (RNN) and the Convolutional Neural Network (CNN), which are primarily Long Short-Term Memory networks (LSTM). In the research, CNN is widely recognized as a network that works well as "feature extractors," whereas RNN is good at simulating problems involving learning orderly sequences [13].

Ensemble methods. Finally, in order to enhance the performance of the model, researchers frequently employ ensemble approaches. This approach aggregates the predictions by combining a number of separate independent models. In fact, numerous papers have demonstrated that applying an ensemble method yields excellent results on the training model and also greatly reduces the testing error.

## II. LITERATURE REVIEW

H. S. Alatawi, et.al (2021) focused on computing the feasibility to detect white supremacist hate speech in automatic way on Twitter when DL (deep learning) and NLP (natural language processing) methods were implemented [14]. Two DL algorithms were suggested for detecting hate speech. The initial algorithm made the deployment of BiLSTM (Bidirectional Long Short-Term Memory) with DSWEs (domain-specific word embeddings). These embeddings were extracted from white supremacist corpus for capturing the semantic of white supremacist slangs and coded words. BERT (Bidirectional Encoder Representations from Transformers) was employed in the second algorithm. According to the experiments, the initial algorithm offered F1-score of 0.75 and 0.80 using the latter one for detecting the hate speech on Twitter and a Storm front dataset.

A. Rodriguez, et.al (2022) introduced a new mechanism recognized as FADOHS in which DA (data analysis) was combined with NLP (natural language processing) methods for sensitizing all social media providers to the generality of hate on social media [15]. In particular, the algorithms of analyzing the sentiment and emotion exploited for analyzing the recent posts and comments on these pages. This mechanism aimed to process the post which were suspicious to involve dehumanizing words prior to be utilized in the clustering algorithm to accomplish further evaluation. The experimental outcomes demonstrated that the introduced mechanism enhanced the efficiency up to 10% as compared to the existing methods with regard to precision, recall, and F1 scores.

M. Bilal, et.al (2022) presented the annotation guidelines for RUHS (Roman Urdu Hate Speech) [16]. Thereafter, a novel RU-HSD-30K dataset was generated which a team of experts had annotated via the annotation rules. A context-aware technique was constructed on the basis of Bi-LSTM (Bidirectional Long Short-Term Memory) with an attention layer and a custom word2vec model was implemented to perform word embeddings. The fundamental aim was to evaluate the impact of lexical normalization of RU words on the efficiency of constructed technique. The experiments indicated that the constructed approach yielded an accuracy up to 87.5% and an F-Score around 88.5 in contrast to others. Furthermore, this approach was applicable on the unseen data.

B. Pariyani, et.al (2021) formulated an automated framework to detect the hate speech [17]. Diverse methods of NLP (Natural Language Processing) were employed to classify the hate speech on the basis of ML (Machine Learning) algorithms. The hate speech was classified from the tweets using ML algorithms namely SVM (Support Vector Machine), LR (Logistic Regression) and RF (Random Forest). The results generated from data without preprocessing exhibited that RF with BoW (bag of words) offered F1 Score of 65.80% and accuracy of 96.29%. When the data was preprocessed, gridsearch SVM with Tf-IDF performed well and attained F1 Score of 74.88% and accuracy of 96.68%.

H. Sohn, et.al (2019) established a MC-BERT system (multi-channel system in which three versions of Bidirectional Encoder Representations from Transformers): English, Chinese, and multilingual BERTs to detect the hate speech [18]. The training and test sentences were translated to the corresponding languages essential for dissimilar BERT models for investigating the usage of translations as additional input. Three datasets namely: 2019 SemEvalHatEval Spanish dataset, 2018 GermEval and 2018 EvalItaHaSpeeDe Italian dataset were applied to simulate the established system. The experimental results revealed that the established system performed more effectively in contrast to the traditional methods on these datasets.

S. W. A. M. D. Samarasinghe, et.al (2020) recommended a DL (deep learning) model in which two CNNs (convolution neural networks) algorithms employed for classifying a given text corpus as hateful or not [19]. Thereafter, in case of involvement of hate content text in the corpus, the text was classified again on the basis of its hate level occurred due to the authorities to make decisions. The text data was transformed into numerical vectors using FastText word embedding. The results proved that the recommended model offered an accuracy of 0.83 to classify the hate speech 0.60 to classify the hate level.

C. Baydogan, et.al (2021) projected a metaheuristic based automatic HSD (hate speech detection) system to generate the promising outcomes for detecting the hate speech [20]. This system employed ALO (Ant Lion Optimization) and MFO (Moth Flame Optimization) algorithms to detect the hate speech. First of all, the basic NLP (natural language processing)

stages were executed. Bag of Words (BoW), Term Frequency (TF), and document vector (Word2Vec) were employed to extract the attributes. After that, diverse data of real time was considered to quantify the projected algorithms. The projected algorithms were more effective as compared to the other methods with regard to different parameters such as accuracy, sensitivity, precision, and f-score. The results of experiments validated that the projected system offered accuracy of 92.1% with initial algorithm and 90.7% with second algorithm. Moreover, this system was effective for dealing with different social media and networking issues.

S. Khan, et.al (2022) designed BiCHAT which was a novel BiLSTM (Bidirectional Long Short-Term Memory) algorithm utilized with deep CNN (convolutional neural network) and HADL (Hierarchical ATtention-based deep learning) algorithm was presented to learn the tweet representation for detecting the hate speech [21]. The tweets were employed for input in this algorithm and underwent from the BERT layer. The BiLSTM made the deployment of convolutional encoded representation. In the end, a soft max layer exploited to assign label to the tweet as hateful or normal. Three datasets taken from Twitter were executed to train and compute the designed algorithm. The results confirmed the supremacy of the designed algorithm over the traditional methods concerning precision, recall, and f-score.

N. Badri, et.al (2022) investigated a technique called BiGRU Glove FT in which Glove and FastText word embedding were integrated and employed as input features and a BiGRU (Bidirectional-Gated Recurrent Unit) model was implemented for detecting the hate speech from social media websites [22]. The investigated technique worked effectively to detect the inappropriate content. This technique focused on detecting the hate speech on OLID dataset. For this, an effective learning process was put forward for classifying the text into offensive and normal language. The results depicted that the accuracy of investigated technique was measured 84%, precision was 87%, recall was 93%, and f1-score was 90% for detecting hate speech.

P. Sharmila, et.al (2022) suggested a novel PDHS (Pattern-based Deep Hate Speech) framework for detecting the presence of hate speech based on a cross-attention encoder with a DLA (dual-level attention) method [23]. Unlike the existing methods, this framework had not concatenated the attributes. This framework emphasized computing the dot product attention to attain superior representation when the irrelevant features were mitigated. The initial level of attention aimed to extract the aspect terms. To achieve this, the predefined POS (parts-of-speech) tagging was utilized. The second level was executed to extract the sentiment polarity so that a pattern was generated. The term frequency, parts-of-speech tag, and Sentiment Scores were considered to train the extracted patterns. The experimental results on Twitter Dataset depicted that the suggested framework was capable of learning effective attributes for enhancing the performance at lower training time and offered F1Score of 88%.

## 2.1 Comparison Table

Author	Year	Technique Used	Results	Limitations
H. S. Alatawi, et.al	2021	BiLSTM and BERT	According to the experiments, the initial algorithm offered F1-score of 0.75 and 0.80 using the latter one for detecting the hate speech on Twitter and a Storm front dataset.	The second algorithm of this approach was incapable of detecting intentionally misspellings and common slang from hate community and some of the datasets in the experiments were found imbalanced.
A. Rodriguez, et.al	2022	FADOHS	The experimental outcomes demonstrated that the introduced mechanism enhanced the efficiency up to 10% as compared to the existing methods with regard to precision, recall, and F1 scores.	The issue of misidentification was occurred. This mechanism was not considered data from comments and replies for accurately identifying the individuals who promoted hate speech.
M. Bilal, et.al	2022	Bi-LSTM with an attention layer	The experiments indicated that the constructed approach yielded an accuracy up to 87.5% and an F-	The efficacy of these DL (deep learning) algorithms was affected due to normalization.

			Score around 88.5 in contrast to others. Furthermore, this approach was applicable on the unseen data.	
B. Pariyani, et.al	2021	an automated framework	The results generated from data without preprocessing exhibited that RF with BoW (bag of words) offered F1 Score of 65.80% and accuracy of 96.29%. When the data was preprocessed, gridsearch SVM with Tf-IDF performed well and attained F1 Score of 74.88% and accuracy of 96.68%.	This framework was applicable only on twitter dataset, and detecting the hate speech from big data was a complex task.
H. Sohn, et.al	2019	A multi-channel system	The experimental results revealed that the established system performed more effectively in contrast to the traditional methods on these datasets.	This system was ineffective to mine the text in social media and knowledge transfer.
S. W. A. M. D. Samarasinghe, et.al	2020	DL (deep learning) model	The results proved that the recommended model offered an accuracy of 0.83 to classify the hate speech 0.60 to classify the hate level.	The amount of data sources used to extract, Sinhala text containing hate speech, was not enough. A small sized data set was employed in this work.
C. Baydogan, et.al	2021	a metaheuristic based automatic HSD (hate speech detection) system	The results of experiments validated that the projected system offered accuracy of 92.1% with initial algorithm and 90.7% with second algorithm. Moreover, this system was effective for dealing with different social	This system was not implemented on diverse datasets.

			media and networking issues.	
S. Khan, et.al	2022	BiCHAT	The results confirmed the supremacy of the designed algorithm over the traditional methods concerning precision, recall, and f-score.	This algorithm had not contained sentiment, content, and other profile-related features.
N. Badri, et.al	2022	BiGRU Glove FT	The results depicted that the accuracy of investigated technique was measured 84%, precision was 87%, recall was 93%, and f1-score was 90% for detecting hate speech.	The efficacy of the investigated technique was restricted to only the utilized dataset and it was not suitable for additional datasets.
P. Sharmila, et.al	2022	a new PDHS (Pattern-based Deep Hate Speech) framework	The experimental results on Twitter Dataset depicted that the suggested framework was capable of learning effective attributes for enhancing the performance at lower training time and offered F1Score of 88%.	The suggested framework was not utilized on unstructured social media datasets with regard to multi-modality attributes such as images and emojis.

### III.CONCLUSION

The NLP approaches are designed in a way that allows the computer to comprehend orders given in natural language and act in accordance with them. It should be emphasised that spoken language and written language are the two categories into which natural language processing can be split. The various schemes are proposed for the hate speech detection and those techniques can be natural language processing. In future natural language processing will be improved for the hate speech detection.

### References

- [1] N. Shawkat, J. Simpson and J. Saquer, "Evaluation of Different ML and Text Processing Techniques for Hate Speech Detection," 2022 4th International Conference on Data Intelligence and Security (ICDIS), Shenzhen, China, 2022, pp. 213-219
- [2] N. S. Mullah and W. M. N. W. Zainon, "Advances in Machine Learning Algorithms for Hate Speech Detection in Social Media: A Review," in *IEEE Access*, vol. 9, pp. 88364-88376, 2021
- [3] B. R. Amrutha and K. R. Bindu, "Detecting Hate Speech in Tweets Using Different Deep Neural Network Architectures," 2019 International Conference on Intelligent Computing and Control Systems (ICCS), Madurai, India, 2019, pp. 923-926
- [4] G. H. Panchala, V. V. S Sasank, D. R. HarshithaAdidela, P. Yellamma, K. Ashesh and C. Prasad, "Hate Speech & Offensive Language Detection Using ML & NLP," 2022 4th International Conference on Smart Systems and Inventive Technology (ICSSIT), Tirunelveli, India, 2022, pp. 1262-1268,
- [5] P. K. Roy, A. K. Tripathy, T. K. Das and X. -Z. Gao, "A Framework for Hate Speech Detection Using Deep Convolutional Neural Network," in *IEEE Access*, vol. 8, pp. 204951-204962, 2020
- [6] Rahul, V. Gupta, V. Sehra and Y. R. Vardhan, "Ensemble Based Hinglish Hate Speech Detection," 2021 5th International Conference on Intelligent Computing and Control Systems (ICICCS), Madurai, India, 2021, pp. 1800-1806
- [7] M. U. S. Khan, A. Abbas, A. Rehman and R. Nawaz, "HateClassify: A Service Framework for Hate Speech Identification on Social Media," in *IEEE Internet Computing*, vol. 25, no. 1, pp. 40-49, 1 Jan.-Feb. 2021
- [8] Rahul, V. Gupta, V. Sehra and Y. R. Vardhan, "Hindi-English Code Mixed Hate Speech Detection using Character Level Embeddings," 2021 5th International Conference on Computing Methodologies and Communication (ICCMC), Erode, India, 2021, pp. 1112-1118
- [9] U. A. N. Rohmawati, S. W. Sihwi and D. E. Cahyani, "SEMAR: An Interface for Indonesian Hate Speech Detection Using Machine Learning," 2018 International Seminar on Research of Information Technology and Intelligent Systems (ISRITI), Yogyakarta,

- Indonesia, 2018, pp. 646-651
- [10] O. Oriola and E. Kotzé, "Evaluating Machine Learning Techniques for Detecting Offensive and Hate Speech in South African Tweets," in *IEEE Access*, vol. 8, pp. 21496-21509, 2020
- [11] S. A. Kokatnoor and B. Krishnan, "Twitter Hate Speech Detection using Stacked Weighted Ensemble (SWE) Model," 2020 Fifth International Conference on Research in Computational Intelligence and Communication Networks (ICRCICN), Bangalore, India, 2020, pp. 87-92
- [12] H. Rathpisey and T. B. Adji, "Handling Imbalance Issue in Hate Speech Classification using Sampling-based Methods," 2019 5th International Conference on Science in Information Technology (ICSITech), Yogyakarta, Indonesia, 2019, pp. 193-198
- [13] K. Mnassri, P. Rajapaksha, R. Farahbakhsh and N. Crespi, "BERT-based Ensemble Approaches for Hate Speech Detection," GLOBECOM 2022 - 2022 IEEE Global Communications Conference, Rio de Janeiro, Brazil, 2022, pp. 4649-4654
- [14] H. S. Alatawi, A. M. Alhothali and K. M. Moria, "Detecting White Supremacist Hate Speech Using Domain Specific Word Embedding With Deep Learning and BERT," in *IEEE Access*, vol. 9, pp. 106363-106374, 2021
- [15] A. Rodriguez, Y. -L. Chen and C. Argueta, "FADOHS: Framework for Detection and Integration of Unstructured Data of Hate Speech on Facebook Using Sentiment and Emotion Analysis," in *IEEE Access*, vol. 10, pp. 22400-22419, 2022
- [16] M. Bilal, A. Khan, S. Jan and S. Musa, "Context-Aware Deep Learning Model for Detection of Roman Urdu Hate Speech on Social Media Platform," in *IEEE Access*, vol. 10, pp. 121133-121151, 2022
- [17] B. Pariyani, K. Shah, M. Shah, T. Vyas and S. Degadwala, "Hate Speech Detection in Twitter using Natural Language Processing," 2021 Third International Conference on Intelligent Communication Technologies and Virtual Mobile Networks (ICICV), Tirunelveli, India, 2021, pp. 1146-1152
- [18] H. Sohn and H. Lee, "MC-BERT4HATE: Hate Speech Detection using Multi-channel BERT for Different Languages and Translations," 2019 International Conference on Data Mining Workshops (ICDMW), Beijing, China, 2019, pp. 551-559
- [19] S. W. A. M. D. Samarasinghe, R. G. N. Meegama and M. Punchimudiyanse, "Machine Learning Approach for the Detection of Hate Speech in Sinhala Unicode Text," 2020 20th International Conference on Advances in ICT for Emerging Regions (ICTer), Colombo, Sri Lanka, 2020, pp. 65-70
- [20] C. Baydogan and B. Alatas, "Metaheuristic Ant Lion and Moth Flame Optimization-Based Novel Approach for Automatic Detection of Hate Speech in Online Social Networks," in *IEEE Access*, vol. 9, pp. 110047-110062, 2021
- [21] S. Khan, M. Fazil and A. R. Baig, "BiCHAT: BiLSTM with deep CNN and hierarchical attention for hate speech detection", *Journal of King Saud University - Computer and Information Sciences*, vol. 34, no. 7, pp. 4335-4344, 21 May 2022
- [22] N. Badri, F. Koubi and A. H. Chaibi, "Combining FastText and Glove Word Embedding for Offensive and Hate speech Text Detection", *Procedia Computer Science*, vol. 207, no. 12, pp. 769-778, 19 October 2022
- [23] P. Sharmila, K. S. M. Anbananthen, D. Chelliah, S. Parthasarathy and S. Kannan, "PDHS: Pattern-Based Deep Hate Speech Detection With Improved Tweet Representation," in *IEEE Access*, vol. 10, pp. 105366-105376, 2022