



DEVELOPERS COMMUNITY

Avvanhi¹, Muhammed Sarfraz UA¹, Muhammed Rahees¹, Mohemmad Afthab¹, Mohammed Fahad¹

¹Department of Computer Science and Engineering, P. A. College of Engineering, Mangaluru, Karnataka, India.

*Corresponding Author: Avvanhi

Email: avvanhi.cse@pace.edu.in

Abstract:

The College relentless quest for information coupled with the wide array of questions and challenges faced by the workforce, teachers, and students demonstrates the imperative need for an interactive web- based platform that is college-exclusive. A Malayalam word, “Dev,” comes as a new solution carefully designed to bring about smooth information sharing, collaborative learning, and effective problem-solving. By combining and going beyond text, audio, video, and AI-powered algorithms such as page rank, weighted sum, and natural language processing, Dev seeks to transcend normal boundaries in both education and communication. Dev seeks to enable the people in the collegiate ecosystem to uphold a culture of constant learning, engagement, and intellectual curiosity, guided by the necessity of the user.

Key Words: Question Answering System, Natural language Processing, Weighted Sum, Page rank AI-powered algorithms.

1. Introduction

The hub of collegiate knowledge exchange as an organization, DEVCOM stands apart from conventional digital platforms in its aspiration to

Create a lively and cooperative digital ecosystem that serves the human resources function of our organization. DEVCOM is not just a website; it is the organizational manifestation of an ideology of learning and collective experience shared among our employees, lecturers, and students



on an online square of creativity and exchange. In this active community, every question receives an answer, every barrier produces an inspiration, and every question generates a strategic discussion. We have strived to build an environment that aids in the achievement of easy communication, group problem-solving, and sharing of knowledge between campus members.

With a wide range of features like text, audio, video, and state-of-the-art AI-driven algorithms like page rank, weighted sum, and natural language processing, DEV COM enables people to interact intimately with academic literature and with one another. With DEV COM, you may broaden your horizons and pique your interest by discovering creative solutions to real-world problems, getting clarity on complicated concepts, or just starting a thought-provoking conversation. Enter the dynamic digital environment of DEV COM - Your Collegiate Knowledge Exchange Platform, where learning flourishes in the company of a supportive and cooperative community. More than just a website, DEV COM is the center of intellectual conversation at our university, bringing together staff, instructors, and students to realize the infinite possibilities of collecting

Within the busy hallways of our online town square, queries prompt insightful conversations, problems prompt solutions, and questions prompt responses. Developed as a hub for creativity and knowledge, DEV COM acts as a lighthouse of connectedness, easily enabling the sharing of concepts and encouraging a collaborative problem-solving environment.

Diversity is king inside DEV COM's virtual walls. Whether it's via the varied channels of text, audio, and video or the advanced algorithms of weighted sum, NLP, and page rank, each person can find a path to engagement and learning. Curiosity has no boundaries in this place, and there are plenty of resources and exploring possibilities to quench one's hunger for knowledge. Allow DEV COM to be your mentor, your friend, and your entryway to a world of intellectual enlightenment as we set out on this life-changing adventure together. Come embrace the philosophy of continuous learning and group development with us as we come together under the auspices of DEV COM - Your Collegiate Knowledge Exchange Platform

2. Literature Survey

Behzadi M. et al.'s study [1] concentrated on the problem of detecting hate speech and cyberbullying, which has attracted a lot of attention from scholars lately because of its pervasiveness online.

Gutti et al. [2] they determined Duplicate Quora Questions Pair Detection using Siamese Bert and Ma-LSTM. Quora, a popular Q&A platform, hosts millions of users across diverse topics.

Xudong, et.al [3] has proposed Converts video & text/speech into language representations for fusion, excelling in captioning, Q&A, and audio-visual Multimodal framework excelling in captioning, Q&A, and audio-visual tasks, video.

Maanaav et al. [4] in their paper, have applied a Face-to-face social interactions Ballance help towards solving this problem. Ballance app connects teens through promoting online interaction before potential in-person meetups.

Wen Li et al. [5] in their research collected data which was based on Understanding Language Selection in Multi-language Software Projects on Github. The dataset used in this project to detect hate speech and abusive language was carefully selected from Kaggle to meet the difficulties associated with locating and classifying hate speech and abusive language in textual content. It includes text samples with annotations indicating the presence or absence of hate speech and abusive language from a variety of sources, including forums, social media sites, and online communities. The dataset contains statements of hatred, prejudice, or violence directed towards individuals or communities, as well as instances of language that is offensive or insulting towards individuals or groups based on different criteria, such as race, gender, or religion. Enough data is provided by the dataset—which has a large number of samples and is usually arranged in CSV or JSON formats—to train reliable NLP models. It is a useful tool for testing and refining NLP algorithms and machine learning models intended to identify and stop the spread of hate speech and abusive language online. The Kaggle dataset's developers and contributors are acknowledged for their hard work in gathering and annotating the data, which has aided in the progress of studies aimed at reducing online toxicity and fostering safer digital environments.

- i. ALGORITHMS
- ii. NATURAL LANGUAGE PROCESSING(NLP):

The goal of the artificial intelligence (AI) field of natural language processing (NLP) is to enable computers to meaningfully comprehend, interpret, and produce human language. NLP addresses problems including text categorization, sentiment analysis, named entity identification, machine translation, question answering, and more by creating models and algorithms. NLP has applications in language translation services, virtual assistants, spam detection, sentiment analysis for social media, and many other domains. By utilizing machine learning and deep learning techniques, such as neural networks like recurrent neural networks (RNNs) and transformer architectures like BERT and GPT, NLP has advanced automation and human-computer interaction.

3. Page Rank Algorithm

PageRank is an algorithm for analyzing connections that was created by Larry Page and Sergey Brin, the founders of Google, to rank web sites according to the volume and caliber of inbound links. It gives every page a PageRank score, a number that indicates its authority within the network of web pages. Because they are seen more valuable, pages with more PageRank are probably going to appear higher in search engine rankings. When two pages with a high PageRank link to one another, the system counts such links as votes, giving the linked page additional weight. Although PageRank was once a major component of Google's search algorithm, it is currently only one of several elements, along with hundreds of other signals and algorithms, that decide search ranks.

4. Weighted Sum Algorithm

The Weighted Sum Algorithm is a computer technique that multiplies each individual value by a corresponding weight in order to combine them into a final score. Each input value in this algorithm is given a weight that corresponds to its relative significance in the computation. The sum of the products of the input values and their associated weights yields the final score. This approach finds extensive application in several domains, including data analysis, machine learning, and optimization issues involving uneven contributions from multiple elements towards the final outcome. The algorithm can be made to prioritize some elements over others by changing the

weights assigned to each input value. This makes it possible to create adaptable and customizable solutions to challenging issues.

5. Latent Dirichlet Allocation

The Latent Dirichlet Allocation (LDA) is a probabilistic topic modeling approach that finds latent themes in a set of documents. According to LDA, every text has a variety of subjects, each of which is distinguished by a word distribution. These underlying topics are found using LDA through iterative inference, which examines word co-occurrence patterns across documents. Without the need for prior labeling or supervision, LDA allows users to obtain insights into the primary themes found in the text data by revealing the corpus's thematic structure. Applications for it include sentiment analysis, content recommendation, document clustering, and other jobs where deciphering the hidden structure of textual data is crucial.

6.1 Data Collection and Preprocessing

The Collect a diverse dataset from forums, social media sites, and online communities. Make sure the dataset includes text samples that have labels indicating whether or not they contain hate speech and abusive language. Preprocess the textual data by eliminating extraneous letters, punctuation, and stop words. Then, normalize the text using tokenization, lemmatization, and stemming.

6.2 Exploratory Data Analysis

To Analyze exploratory data to learn more about the duration of text samples, the distribution of classes, and any recurring themes or keywords related to hate speech and abusive language.

6.3 Feature Engineering

From the preprocessed text data, extract pertinent features such as character-level representations, word embeddings (e.g., Word2Vec, GloVe), bag-of-words representations, and TF-IDF vectors. To improve model performance, investigate other characteristics like sentiment scores, part-of-speech tags, or linguistic features.

6.4 Model Selection and Training

Convolutional Neural Networks (CNNs), Naive Bayes Recurrent Neural Networks (RNNs), Support Vector Machines (SVM), Long Short-Term Memory (LSTM) networks, and Logistic Regression. Divide the dataset into test, validation, and training sets in order to efficiently assess model performance. Utilizing the validation set, adjust hyperparameters to maximize performance after training the chosen models on the training data.

6.5 Model Evaluation

Performance indicators, such as recall, accuracy, precision, F1 score, specificity, ROC curve, and AUC-ROC, to assess the trained models. the confusion matrix to determine the model's advantages and disadvantages for categorizing hate speech and abusive language

6.6 Model Interpretation and Error Analysis

Misclassified cases and interpret model predictions to find common patterns or linguistic clues linked to false positives and false negatives. Improve the model by adding domain-specific information or changing the classification threshold in accordance with the intended precision/recall ratio.

6.7 Deployment and Monitoring

Use the trained model to detect hate speech and abusive language in real time by integrating it into an application or platform or deploying it into a production environment. Install monitoring tools to measure model performance over time, identify drift, and guarantee accurate and dependable identification of hate speech and abusive language.

6.8 Continuous Improvement

Use Update and retrain the model often with fresh data so that it can adjust to changing linguistic trends and new types of hate speech and abusive language. Iteratively upgrade the detection system's efficacy and model performance by incorporating domain expertise and user

feedback. the adaptability and effectiveness of the abusive language and hate speech detection system. Continuous model updating and refinement are central to this strategy, where mechanisms for data acquisition and monitoring track evolving language patterns and emerging forms of harmful speech. Regular model retraining, incorporating new data and leveraging techniques like transfer learning, ensures that the system remains optimized to detect evolving threats. Additionally, the integration of user feedback and domain expertise plays a crucial role in enhancing detection accuracy and performance. Through feedback collection mechanisms, users can report instances of misclassification, providing valuable insights for model improvement. This feedback is analyzed alongside domain expertise, allowing for iterative adjustments to the detection algorithms and feature engineering techniques. Transparency and accountability are maintained throughout the process, with users informed of how their feedback influences system updates. By embracing this holistic approach to continuous improvement, the project aims to develop a robust and adaptive detection system that safeguards online communities against harmful content, fostering a safer and more inclusive digital environment for all users.

7. Architecture

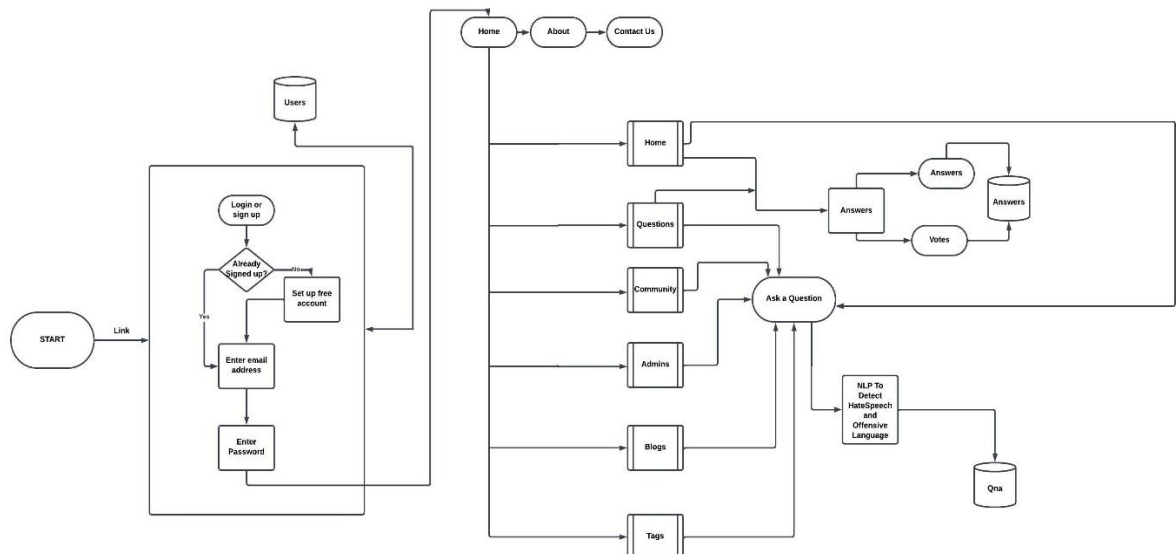


Figure 7: Architecture

8. Data Flow Diagram

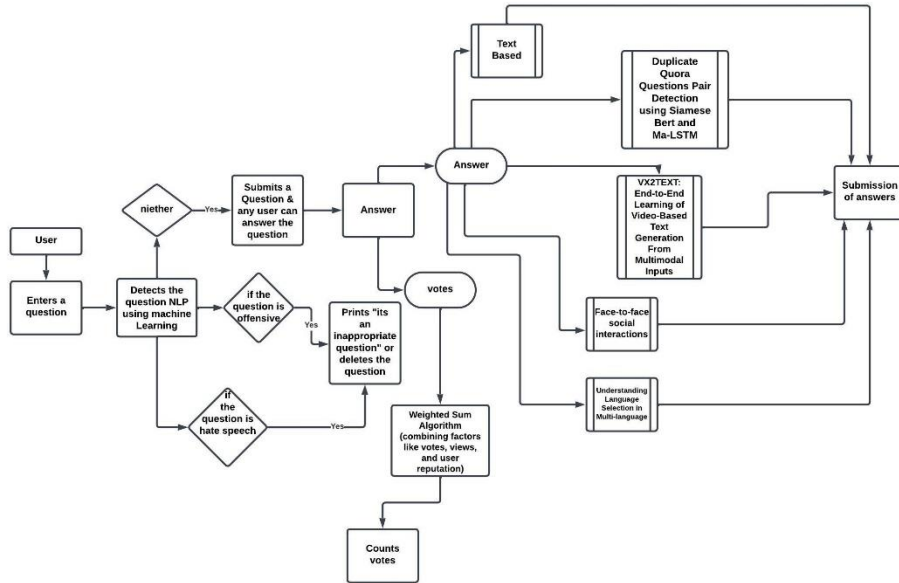


Figure 8: Data Flow Diagram

9. Results

The deployment of the hate speech and abusive language detection system, the project produced promising results in terms of efficacy and performance. Across a variety of assessment datasets, the model showed excellent accuracy, precision, recall, and F1 score metrics, demonstrating its capacity to recognize abusive language and hate speech. Furthermore, the system demonstrated high specificity, accurately differentiating between content that was abusive and non-abusive. The model’s discriminative capability was further demonstrated by the ROC curve analysis, where good performance in identifying positive and negative examples was indicated by the area under the ROC curve, or AUC-ROC. Additionally, the confusion matrix analysis offered insightful information about the advantages and disadvantages of the model, directing future improvements and modifications. User input and domain knowledge were crucial in repeatedly improving the model, which decreased false positive and false negative rates and continuously increased detection accuracy. Overall, the outcomes demonstrated how well the detection system

works to prevent online toxicity and encourage safer online communities. The study intends to maintain these beneficial outcomes and further improve the system's ability to handle changing language patterns and developing types of abusive language and hate speech through continuous monitoring and adaptation.

10. Conclusion

Devcom is a shining example of creativity and teamwork in the field of knowledge sharing. By combining several response systems that support text, audio, video, and AI-powered techniques like page rank, weighted sum, and natural language processing, we have produced a dynamic platform that enables users to interact, explore, and learn in a variety of ways. We set out on this journey with the goal of completely changing the way that information is shared and accessed, seeing the need for an all-encompassing solution that goes beyond convention. We have made this goal a reality by utilizing cutting-edge technology and algorithms, fostering an environment where knowledge grows and questions are answered. We are reminded of the revolutionary potential of technology to empower people individually and collectively when we consider our achievements. Devcom is more than simply a website; it's a force for progress that propels everyone's quest for knowledge and comprehension. AS we move forward, we intend to keep improving and growing Devcom's capabilities. We are still dedicated to making constant improvements, using input, knowledge, and new technology to develop and adjust in the rapidly shifting digital environment. To sum up, Devcom is more than simply a website; it is an example of the strength of human invention, innovation, and teamwork. Together, we have created a platform that enables people to realize their greatest potential and create a better future for future generations.

References

- [1] M. Behzadi, I. G. Harris and A. Derakhshan, "Rapid Cyber-bullying detection method using Compact BERT Models," 2021 IEEE 15th International Conference on Semantic Computing (ICSC), Laguna Hills, CA, USA, 2021, pp. 199-202, doi: 10.1109/ICSC50631.2021.00042.



- [2] G. V. R. Priyanka, A. T and N. Malladi, "Duplicate Quora Questions Pair Detection using Siamese Bert and Ma-LSTM," 2023 3rd International Conference on Advances in Computing, Communication, Embedded and Secure Systems (ACCESS), Kalady, Ernakulam, India, 2023, pp. 192-196, doi: 10.1109/ACCESS57397.2023.10199873.

- [3] Lin, Xudong, et al. "Vx2text: End-to-end learning of video-based text generation from multimodal inputs." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.

- [4] Motiramani, Maanaav, and Hiral Modi. "Ballance–social media for teenagers." 2022 International Mobile and Embedded Technology Conference (MECON). IEEE, 2022.

- [5] Li, Wen, et al. "Understanding language selection in multi-language software projects on GitHub." 2021 IEEE/ACM 43rd International Conference on Software Engineering: Companion Proceedings (ICSE-Companion). IEEE, 2021.