



# IJRASET

International Journal For Research in  
Applied Science and Engineering Technology



---

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

---

**Volume:** 14    **Issue:** IV    **Month of publication:** April 2026

**DOI:** <https://doi.org/10.22214/ijraset.2026.81277>

[www.ijraset.com](http://www.ijraset.com)

Call:  08813907089

E-mail ID: [ijraset@gmail.com](mailto:ijraset@gmail.com)

# Ethical AI Analyzer

Nityam Kumar Tiwari, Neeraj Patel, Himanshu, Nistha Garg, Mr. Jagbeer Singh

Department of Computer Science and Engineering Meerut Institute Of Engineering and Technology, Meerut, Uttar Pradesh, India

**Abstract:** *Human biases significantly influence the efficiency and fairness of Natural Language Processing (NLP) systems. Social biases, encompassing gender, race, and caste frequently manifest within linguistic datasets and are subsequently potentially yielding inequitable results. Although the investigation of bias within textual data has expanded in recent years, existing resources remain constrained, fragmented and predominantly concentrated on specific categories such as gender or race with insufficient consideration given to caste or the interplay of multiple biases.*

*This study has two main goals: first, to investigate how large language models (LLMs) and transformer-based architectures can be used to find and categories different types of social bias in text.*

*Second, to analyze the challenges of limited datasets, uneven distribution of categories, and the difficulties of fair evaluation. The results highlights the need for more diverse and representative datasets, and they provide insights into creating more fair and inclusive Natural Language Processing.*

## I. INTRODUCTION

There is growing recognition that human biases are deeply embedded in language and can strongly affect the behavior and fairness of Natural Language Processing (NLP) systems. These biases often reflect broader social inequalities and when learned by machine learning models, can result in discriminatory outcomes in real-world application. As a result, research on bias detection, classification, and mitigation in NLP has gained increasing attention, with broader goal of developing more fair and responsible AI systems.

Despite this progress, bias detection in NLP remains insufficiently explored. Many existing datasets are small in scale, focus on a limited set of categories such as gender or race, or were originally designed for related tasks like hate speech detection rather than explicit bias identification.

Context-dependent forms of bias, particularly caste-based bias, are largely absent from mainstream resources.

These gaps raise an important question: How can NLP models be effectively trained to detect and classify multiple forms of social bias using limited and imperfect data?

To address this question, this study pursues three main objectives. First, it analyzes publicly available datasets related to bias and hate speech and uses them as training resources. Second, it introduces a custom-annotated dataset centered on caste bias, extending the scope beyond commonly studied categories. Third, it trains and evaluates transformer-based model and large language model (LLMs) on both binary (biased vs. non-biased) and multi-class (e.g. Gender, race, caste) classification tasks.

This work contributes to existing research by providing a comparative analysis of model performance across different bias categories, highlighting the impact of dataset composition and class imbalance. It also discusses the challenges associated with heterogeneous and under-representative data, emphasizing the need for more balanced and intersectional datasets.

In this study, bias is defined as unequal or unfair treatment based on social characteristics. We consider bias expressed through derogatory language, harmful stereotypes, and abusive or exclusionary statements targeting specific groups. The target categories examined include gender, race, caste, religion, disability, sexual orientation, nationality, and age. While bias and hate speech are closely related, they are not equivalent; therefore, resources from both domains are utilized.

We acknowledge the sensitive nature of studying social bias and recognize existing limitations, such as simplified representations of gender and caste and limited modeling of intersectionality due to dataset constraints. The remainder of the paper is organized as follows: Section 2 reviews related work, section 3 describes dataset construction and preprocessing, section 4 details the modeling and evaluation framework, section 5 discusses the result, and section 6 concludes with directions for future research.

## II. RELATED WORK

We begin our work by reviewing existing research on bias and hate speech detection in Natural Language Processing (NLP), with particular attention to commonly used models, datasets, and their limitations.

We also examine real-world examples of bias language drawn from online platforms such as Twitter (X) and LinkedIn which illustrate how social bias appears in everyday digital communication.

We observe that research on bias in NLP has expanded rapidly in recent years, especially with the rise of transformer-based architecture such as BERT and its variants. Prior studies have demonstrated that these models can encode and propagate harmful stereotypes present in their training data. Bias has been reported across a wide range of NLP tasks, including coreference resolution, sentiment analysis, dialogue generation, and syntactic parsing. These findings suggest that bias is not confined to specific applications but is instead a broader challenge affecting modern NLP systems as a whole.

We identify training data as a major contributing factor to this problem. Large-scale datasets collected from web crawls and social media platforms contain rich and diverse language, but they also include substantial amounts of offensive and discriminatory content. For instance, analyses of widely used corpora such as Common Crawl have shown that a noticeable portion of the data contains hate speech or sexually explicit material. Models trained on such data may therefore unintentionally learn biased associations. High-profile incidents involving systems such as Microsoft's Tay chatbot, Meta's Galactica, and other large language models further highlight the risk of deploying biased AI systems, even when mitigation strategies are applied.

At the same time, we note a growing body of work that focuses on using NLP models themselves to detect and classify bias in legal texts, study the representation of marginalized communities in Wikipedia, and examine disparities in literary reviews across different social groups. Much of this research relies on data collected from social media platforms, particularly Twitter, using keyword-based collection techniques. More recent studies have explored synthetic data generation to address data scarcity, although this data generation to address data scarcity, although this approach raises concerns about realism and representativeness.

Overall, we find that existing datasets are often limited in size, narrowly focused on specific bias categories, or originally designed for related tasks such as hate speech detection rather than explicit bias classification. These limitations motivate our work and underscore the need for broader, more representative datasets and systematic evaluations across multiple forms of social bias.

### III. METHODOLOGY

In this section, we describe the design and implementation of our system for detecting social bias in online textual content. Our methodology is organized into four main stages: (1) data collection and preprocessing, (2) definition of target bias categories, (3) model selection and training, and (4) inference and output generation.

#### A. Problem Definition

The objective of this work is to automatically identify biased language in user-generated online content, with particular emphasis on caste-based, gender-based, religious, and racial bias. Such forms of discrimination frequently appear in social media posts and professional discussions, either explicitly or in more subtle forms. Unlike conventional hate speech detection tasks, our approach aims to capture both overt hostility and implicit biased framing, even when harmful intent is expressed indirectly.

#### B. Data Collection and Processing

We collected textual data from publicly available sources, including platforms such as Twitter and LinkedIn, where biased language is commonly observed. In addition, we performed validation and analysis. To ensure ethical compliance and protect user privacy, all identifiable information was removed or anonymized.

The collected text was cleaned and normalized using standard preprocessing steps, including lowercasing, removal of punctuation, URLs, hashtags, mentions, and emojis, and tokenization using a transformer-based tokenizer with sub-word segmentation. The resulting corpus combines curated benchmark data with real-world examples of both biased and non-biased language, enabling robust and generalizable model training.

#### C. Target Bias Categories

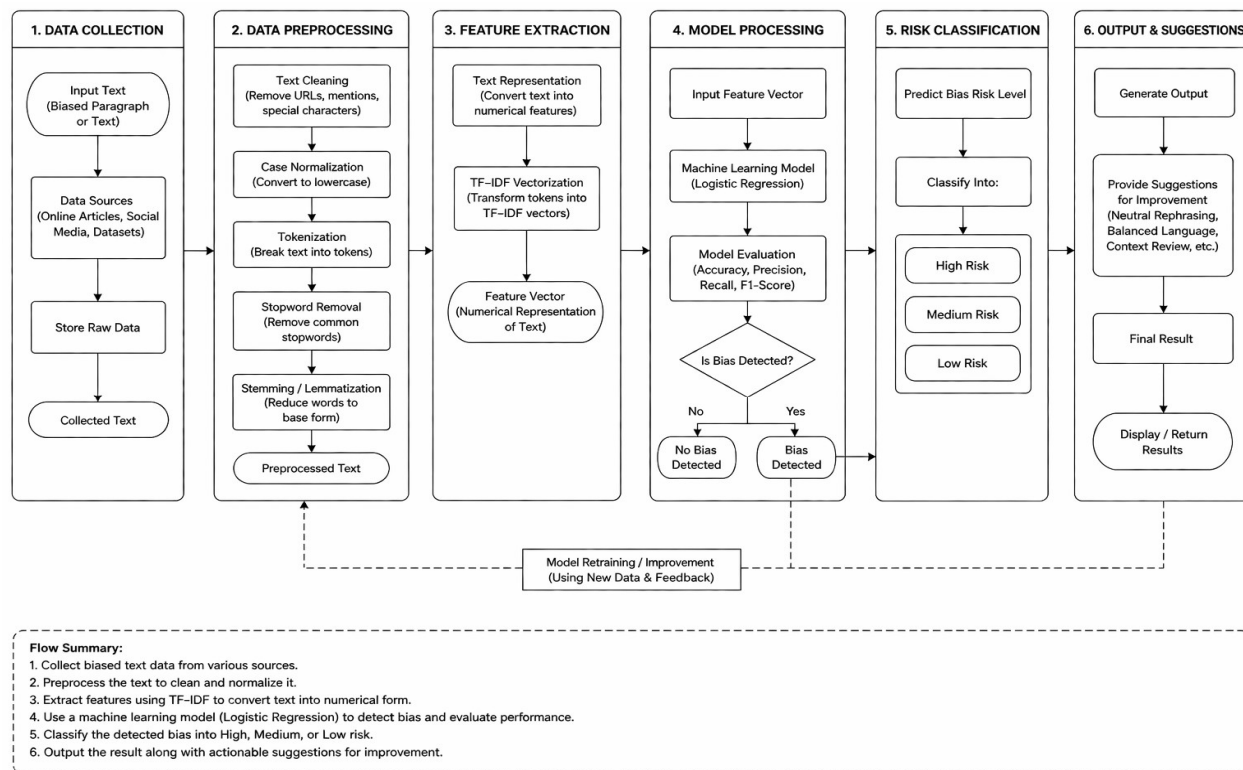
We categorized bias based on the social group being targeted. The primary target categories considered in this study include gender, caste, race and ethnicity, religion, and profession-related bias. Each sample was annotated with a binary label indicating whether it is biased or non-biased. When sufficient contextual or linguistic cues were available, the system further classified the specific types of bias.

#### D. Model Architecture

For model implementation, we employed an Emotion-Transformer architecture built on top of Distill BERT, selected for its balance between computational efficiency and strong contextual representation. The model was fine-tuned on our dataset for text classification. The final hidden layer output is passed through a pooling operation followed by a classification head, producing a binary prediction (biased vs. non-biased) and, when applicable, an associated bias category.

#### E. Flowchart

We illustrate the overall workflow of our system through this flowchart. We begin by taking textual input from the user and initially checking whether the content is biased or unbiased. If the text is found to be unbiased, we classify it as low or no risk and accept it as valid. When biased content is detected, we further process the text using NLP techniques to analyze its severity. We then assess the level of bias and categorize it into different risk levels. For high-risk cases, we suggest an alternative unbiased version of the text, while low and medium risk cases are handled accordingly. This flow allows us to both detect biased language and support bias mitigation in a structured manner.



### IV. MODEL TRAINING

#### A. Experimental Setup

To develop our bias detection system, we adopted the Emotion-Transfer model, which was originally proposed for emotion classification but can be effectively adapted to identify biased language in real-world online content, including tweets, LinkedIn posts, and other publicly available social media text. The architecture is built on Distill BERT, a distilled variant of BERT that offers a favourable balance between computational efficiency and contextual representation capability.

We implemented our approach using the open-source Emotion-Transformer framework and the corresponding Distill BERT model provided by the HuggingFace Transformers library. To assess the model's effectiveness on established benchmarks, we trained and evaluated it on subsets of the Davidson and MLMA datasets. These datasets were selected because their original studies reported clear performance baselines, such as Waseem-Hovy and DynGen, were not included in this comparison due to differences in task formulation and annotation schemes.

On the Davidson dataset, the original work reported an F1-score of 0.90 using a support vector machine with L2 regularization. Using our transformer-based configuration, trained for five epochs, we achieved an F1-score of 0.80. Similarly, for the MLMA dataset, which originally reported an F1-score of 0.43, our model achieved a comparable F1-score of 0.42 when trained for four epochs under the same configuration. These results indicate that the Emotion-Transformer performs competitively with existing approaches and is suitable for the deployment in real-world bias detection scenarios.

## V. RESULT AND DISCUSSION

In this section we, present the results obtained from our trained model and discuss the key observations that emerged during our experiments. As a team, we aimed not only to evaluate how accurately the model classifies biased content but also to understand how well it generalizes to different forms of social bias encountered on real-world platforms such as Twitter and LinkedIn.

### A. Performance on Benchmark Datasets

To establish a reliable performance baseline, we evaluated our model on two publicly available benchmark datasets: Davidson and MLMA. We selected these datasets because their original studies reported F1-scores, allowing us to make direct comparisons. On the Davidson dataset, we achieved an F1-score of 0.80, compared to the originally reported score of 0.90 using an SVM with L2 regularization. On the MLMA dataset, our best result was an F1-score of 0.42, which closely matches the original score of 0.43. As a team, we interpret these results as evidence that the Emotion-Transformer performs at a level comparable to existing methods on standard hate speech and toxicity datasets, validating its use for further bias detection tasks.

### B. Bias Detection in Real-World Text

We consider the application of our model to real-world data from Twitter and LinkedIn to be the most significant contribution of our work. While analyzing the outputs together, we observed that these platforms often contain implicit or context-dependent bias that does not rely on overtly offensive language. After fine-tuning the model using both benchmark data and curated real-world examples, we found that it consistently detected biased statements across several categories. We observed that gender bias was identified with relatively high precision, largely due to the availability of labeled examples. In contrast, we found caste-based bias to be more challenging, especially when expressed indirectly through subtle language or microaggressions. However, we noted that including manually curated caste-related examples noticeably improved the model's sensitivity. We also observed that religion- and race-based biases were easier to detect when explicit terms were present, whereas ambiguous or context-heavy cases remained difficult. Overall, we agreed that the model was more confident when handling highly polarized or toxic statements, while nuanced or sarcastic expressions common in professional settings such as LinkedIn were harder to classify reliably.

### C. Challenges Observed

Throughout our experiments, we collectively identified several challenges that influenced model performance. Class imbalance across bias categories, particularly for caste and disability, limited the model's ability to generalize. We also encountered annotation noise in real-world data, where bias was highly dependent on context and interpretation. Additionally, we observed that code-mixed language, especially Hindi-English content commonly found on Indian social media platforms, posed difficulties for models primarily trained on English text. Despite these challenges, we found that the system was still able to flag a substantial portion of biased content accurately and consistently, reinforcing its potential value for bias analysis, awareness, and moderation-related applications.

## VI. CONCLUSION AND FUTURE WORK

### A. Conclusion

In this study, we investigated the problem of social bias detection in online textual content by designing and evaluating a transformer-based system tailored to real-world data. We focused on identifying multiple forms of bias, including gender, caste, religion, and race, using a combination of benchmark datasets and curated examples collected from platforms such as Twitter and LinkedIn. Throughout our experiments, we observed that the proposed approach is effective in detecting explicit and strongly expressed bias, while also maintaining competitive performance on standard evaluation datasets.

At the same time, our analysis revealed important limitations, particularly when dealing with subtle, context-dependent, or underrepresented forms of bias such as caste-based discrimination and microaggressions in professional settings.

These findings reinforce the need for more balanced, diverse, and context-aware datasets, as well as evaluation strategies that go beyond traditional hate speech benchmarks. Overall, our work highlights both the potential and the current challenges of applying transformer-based models to bias detection and underscores the importance of continued research towards more inclusive and socially responsible NLP systems.

### B. Future Work

We have identified five specific areas where this architecture can be expanded. Building on the findings of this study, we plan to extend our work in several important directions. First, we aim to apply the proposed bias detection system in practical settings such as social media platforms, where it can assist in real-time content analysis and moderation. We also see potential in integrating the model into intelligent text input systems, including virtual keyboards and writing assistants, to provide bias-aware suggestions during text composition. In addition, we plan to explore its use within search engines, where identifying and reformulating biased queries could contribute to more neutral and inclusive information access.

## VII. ACKNOWLEDGEMENT

We would like to express our sincere gratitude to our supervisor, Mr. Jagbeer Singh, for his guidance and constructive feedback throughout this project. His support helped shape the direction and execution of our research. This work represents the collective effort of our team - Nityam Kumar Tiwari, Neeraj Patel, Himanshu and Nishtha Garg - who collaboratively designed, implemented and evaluated the Ethical AI Analyzer system.

We also thank the Department of Computer Science and Engineering, MIET, for providing the academic environment, infrastructure, and research support necessary to carry out this work. Finally, we acknowledge the publicly available datasets, open research resources, and prior studies in Ethical AI that informed and strengthened our approach.

Computer Science and Engineering at MIET for giving us the resources and lab support we needed. Finally, we acknowledge the developers behind the Gemini, Hugging Face and OpenAI APIs which made our experiments possible.

## REFERENCES

- [1] A. Founta, C. Djouvas, D. Chatzakoa, L. Tiadis, J. Blackburn, G. Stringhini, A. Vakali, M. Sirivianos, and N. Kuortellis, "Large scale crowdsourcing and characterization of Twitter abusive behavior," *Proc. Int. AAAI Conf. Web and Social Media*, vol. 12, pp. 491-500, 2018.
- [2] A. Field, S. L. Blodgett, Z. Waseem, and Y. Tsvetkov, "A survey of face, racism, and anti-racism in NLP," in *proc. 59th Annu. Meeting Assoc. Comput. Linguistics (ACL)*, pp. 1905-1925, Aug. 2021.
- [3] T. Bolukbasi, K. -W. Chang, J. Y. Zou, V. Saligrama, and A. T. Kalai, "Man is to computer programmer as woman is to homemaker? Debiasing word embeddings," *Advances in Neural Information Processing Systems*, vol. 29, pp. 4349-4357, 2016.
- [4] C. Basta, M. R. Costajussa, and N. Casas, "Evaluating the underlying gender bias in contextualized word embeddings," in *Proc. 1st Workshop on Gender Bias in NLP*, pp. 33-39, 2019.
- [5] W. Guo and A. Caliskan, "Detecting emergent intersectional biases: Contextualized word embeddings contain a distribution of human-like biases," *Proc. ACM Conf. Fairness, Accountability, and Transparency*, pp. 122-133, 2021.
- [6] M. Jiang and C. Fellbaum, "Interdependencies of gender and race in contextualized word embeddings," in *Proc. 2nd Workshop on Gender Bias in NLP*, pp. 17-25, Dec. 2020.
- [7] Y. C. Tan and L. E. Celis, "Assessing social and intersectional biases in contextualized word representation," *arXiv preprint arXiv: 1911.01485*, 2019.
- [8] S. Sharifirad, A. Jacovi, and S. Matwin, "Learning and understanding different categories of sexism using convolutional neural network filters," in *Proc. Widening NLP Workshop*, pp. 21-23, 2019.
- [9] P. Parikh, H. Abburi, P. Badjatiya, R. Krishnan, N. Chhaya, M. Gupta, and V. Verma, "Multi-label categorization of accounts of sexism using a neural framework," in *Proc. EMNLP-IJCNLP*, pp. 1642-1652, 2019.
- [10] H. Liu, W. Wang, Y. Wang, H. Liu, and J. Tang, "Mitigating gender bias for neural dialogue generation with adversarial learning," in *Proc. EMNLP*.
- [11] H. Liu, W. Wang, Y. Wang, H. Liu, and J. Tang, "Mitigating gender bias for neural dialogue generation with adversarial learning," in *Proc. EMNLP*, pp. 893-903, Nov. 2020.



10.22214/IJRASET



45.98



IMPACT FACTOR:  
7.129



IMPACT FACTOR:  
7.429



# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24\*7 Support on Whatsapp)