# Report of Reviewing Research Paper and Individual Presentation:

**Protiva Das**
**24166051**
**CSE713**
**Group:17**

**Paper Title:**

Nᴮɪᴀs**: A natural language processing framework for BIAS identification in text**

**Paper Link:**
https://www.sciencedirect.com/science/article/abs/pii/S0957417423020444

1. **Summary:**

The paper talks about "Nbias," a sophisticated tool designed to uncover and identify bias in text across various fields like social media, healthcare, and job recruitment. Biases in text can lead to unfair stereotypes and discrimination, emphasizing the need for tools that can ensure fair data analysis. Nbias uses advanced techniques to spot words and phrases that may indicate bias.

This research emphasizes the growing prevalence of bias in NLP applications and the consequential need for bias-aware methodologies. Given that biases often stem from the training datasets, which are riddled with inherent prejudices, Nbias seeks to facilitate the development of more equitable and unbiased AI systems. The proposed model utilizes annotated datasets and a semi-autonomous labeling process to enhance the efficiency and accuracy of bias detection.

The development of Nbias involves a comprehensive four-layer process: data collection from diverse sources, corpus construction with detailed annotations, model development using advanced neural network architectures, and a thorough evaluation against established benchmarks. The framework notably outperforms traditional models, achieving accuracy improvements of 1% to 8%. For instance, in evaluations, Nbias achieved impressive accuracy scores: 88.4% in social media, 90.6% in healthcare, and 91.8% in job hiring texts, outpacing competitors like TENER and BART-NER.

The methodology expands on previous techniques by introducing 'BIAS' as a new entity type within its named entity recognition system, allowing for more precise detection of biased phrases and terms. This approach not only enhances the reliability of bias detection in NLP tools but also sets a new standard in the field by filling the critical need for robust, scalable bias detection methods.

By providing an effective solution to a critical problem in NLP, this paper marks a significant step towards the creation of fair, unbiased AI systems, ensuring that automated decision-making processes are equitable across all sectors.

## 1.1 Motivation

The study reveals that text-to-image technology often reinforces stereotypes about different groups of people, highlighting a crucial issue in how these technologies are designed.

## 1.2 Contribution

The research identifies and highlights biases in AI systems, pointing out that current measures to prevent these biases are not sufficient.

## 1.3 Methodology

The approach involved using different types of written prompts to generate images and then analyzing these images to study how biases manifest depending on the nature of the prompt.

## 1.4 Conclusion

The research found that the type of prompt used—whether specific or general—did not significantly change the biased nature of the images produced.

## 2 Limitations

### 2.1 First Limitation

The study's reliance on specific prompts means that any inherent issues with these prompts could have influenced the results. The conclusions of this research are largely based on the specific set of prompts used. If these prompts were biased or flawed in any way, such as being too narrow or not fully representative, the results could be skewed. This dependency raises concerns about whether the findings would hold under different circumstances or with a broader array of prompts.

### 2.2 Second Limitation

The study used a limited number of analytical tools, suggesting that using different tools might reveal different results. Only a select few analytical tools were employed in the study, which may

limit the scope of the findings. Different tools could potentially detect different types of biases or provide deeper insights into the biases already identified. Incorporating a wider variety of tools could help in obtaining a more comprehensive understanding of biases in AI systems, enhancing the reliability and depth of the conclusions

## 3     Synthesis

The study on the Nbias model highlights its strong ability to pinpoint bias in texts related to social media, health, and job roles. The Nbias model effectively identifies biases in various types of text, such as social media posts, health advice, and job advertisements, with accuracy rates like 88.4% for social media bias detection. By selecting unfair stereotypes and biases, the model helps make AI tools more just and equitable. Enhancing the Nbias model could extend its application to broader sectors, ensuring that digital content across platforms is unbiased. This could contribute to promoting fairness in society. Future improvements in the model, guided by the study's findings, could help it achieve even greater accuracy and versatility, expanding its use in combating bias more effectively in diverse contexts. The study invites future researchers to explore new questions such as: How do biases differ across various AI applications and models? Can effective bias-mitigation techniques developed in one context be successfully applied in others? These inquiries are vital for creating more nuanced AI systems that are fair and impartial, particularly in high-stakes areas like recruitment, healthcare, and law enforcement.