



# Performance Analysis of Generative AI in Bias Detection and Mitigation on Text Datasets

Charlotte<sup>1</sup>, Grayson<sup>2</sup>, Matteo Xavier<sup>3</sup>

<sup>1,2,3</sup>School of Electrical Engineering and Computer Science, University of Ottawa, Canada

## Article Info

### Article history

Received : Oct 30, 2025

Revised : Dec 29, 2025

Accepted : Jan 15, 2026

### Keywords:

Generative Artificial Intelligence;

Bias Detection;

Bias Mitigation;

Fairness in AI;

Natural Language Processing (NLP).

## Abstract

*This study investigates the performance of generative artificial intelligence in detecting and mitigating bias within text datasets, addressing a critical challenge in the development of fair and ethical AI systems. This research aims to provide a comprehensive evaluation framework that integrates both bias detection and mitigation, which are often studied separately in existing literature. The methodology employs multiple text datasets, including social media, news articles, and hate speech corpora, to capture diverse forms of bias. Generative models based on transformer architectures, particularly GPT-based and fine-tuned models, are evaluated alongside baseline models. Bias detection is conducted using prompt-based, classifier-based, and lexicon-based approaches, while mitigation strategies include prompt engineering, debiasing algorithms, reinforcement learning with human feedback (RLHF), and data augmentation. Model performance is assessed using a combination of classification metrics (accuracy, precision, recall, F1-score), fairness metrics (demographic parity and equal opportunity), and text quality measures (perplexity, coherence, and semantic similarity). The results indicate that all mitigation techniques contribute to reducing bias, with RLHF and hybrid approaches achieving the highest effectiveness, reducing bias scores by over 50% while significantly improving fairness metrics. This study contributes to AI fairness research by proposing an integrated evaluation framework and demonstrating that it is possible to achieve substantial bias reduction without compromising overall model performance. The findings provide practical insights for the development of more transparent, reliable, and ethically aligned generative AI systems, supporting their responsible deployment in sensitive domains such as healthcare, finance, and hiring.*

## Corresponding Author:

Charlotte, Grayson, Matteo Xavier  
School of Electrical Engineering and Computer Science  
University of Ottawa, Canada  
800 King Edward Avenue, Ottawa, ON K1N 6N5, Canada  
Email: [charlotte@uottawa.ca](mailto:charlotte@uottawa.ca)

This is an open access article under the [CC BY-NC](#) license.



## 1. Introduction

The rapid advancement of generative artificial intelligence, particularly large language models (LLMs), has significantly transformed the landscape of natural language processing and automated content generation. Models such as GPT and other transformer-based architectures are now widely used in applications ranging from chatbots and virtual assistants to content creation and decision-support systems (Kumar, 2020). Their ability to generate coherent, contextually relevant, and human-like text

has made them indispensable in both academic and industrial domains. However, alongside these benefits, there is growing concern regarding the ethical implications of such systems, especially in relation to bias embedded within the generated outputs. Since these models are trained on vast corpora of real-world text data, they often inherit and reproduce societal biases present in the data, including gender, racial, and cultural biases. This raises critical questions about fairness, accountability, and trust in AI-driven systems.

One of the central issues in generative AI is the tendency of models to not only reflect but also amplify existing biases in datasets (Norori et al., 2021). When biased patterns are learned during training, the model may unintentionally generate outputs that reinforce stereotypes or discriminatory narratives. This becomes particularly problematic in sensitive domains such as healthcare, recruitment, education, and legal systems, where biased outputs can lead to unfair or harmful decisions. Despite increasing awareness of this issue, there remains a lack of standardized frameworks and evaluation methodologies to systematically assess and mitigate bias in generative models (Sabuhi et al., 2021). While some approaches focus on identifying bias through various detection mechanisms, others attempt to reduce bias through mitigation strategies such as data preprocessing, algorithmic adjustments, or post-processing techniques. However, these efforts are often fragmented and lack a unified evaluation standard, making it difficult to compare results across studies or determine the most effective approach.

The study of bias in generative artificial intelligence and natural language processing (NLP) has gained significant attention over the past decade, particularly with the rapid development of large language models. A growing body of literature has explored both bias detection and mitigation, although often in fragmented ways (Haddaway et al., 2020). Early foundational work on bias in NLP highlighted how machine learning models inherently learn societal biases from training data. For example, Aylin Caliskan, Joanna Bryson, and Arvind Narayanan (2017) demonstrated that word embeddings trained on large corpora encode human-like stereotypes, particularly gender and racial biases. This study provided empirical evidence that bias is not merely a downstream issue but is embedded at the representation level of language models. Building on this, Roman Binns (2018) examined fairness from a philosophical perspective, emphasizing that technical solutions alone are insufficient without ethical considerations.

In the following years, research shifted toward quantifying and systematically evaluating bias. Paula Czarnowska, Yogarshi Vyas, and Kashif Shah (2021) proposed generalized fairness metrics to measure social bias in NLP systems, highlighting inconsistencies across existing evaluation methods and stressing the importance of standardized metrics for reliable comparison. Around the same time, increasing attention was given to how bias affects downstream applications. For instance, Fatma Elsafoury and Stamos Katsigiannis (2023) investigated how different types of bias such as representation and overamplification bias impact the fairness of toxicity detection systems, showing that mitigation strategies can improve fairness but may vary depending on bias type.

Recent studies have increasingly focused on both detection and mitigation frameworks. Weikun Lin, Jingxuan Xiao, and Zuen Cen (2024) analyzed the impact of training data on bias in NLP systems and emphasized the importance of responsible data curation and preprocessing techniques in reducing bias. Similarly, Xudong Han, Timothy Baldwin, and Trevor Cohn (2023) highlighted the lack of standardization in fairness evaluation, arguing that inconsistent definitions and metrics hinder progress in bias mitigation research.

In addition to evaluation-focused studies, several works have proposed concrete frameworks for bias detection and mitigation. For example, a study by unknown (2024) introduced the Nbias framework, a multi-layered system designed to detect bias across different dimensions in textual data, demonstrating improved performance over baseline models. More recently, unknown (2026) proposed an explainable bias detection framework using transformer-based generative adversarial networks, integrating explainable AI techniques such as SHAP to enhance interpretability in bias identification.

Furthermore, several comprehensive reviews have synthesized developments in this field. Juveria Afreen, Mahsa Mohaghegh, and Maryam Dobarjeh (2025) conducted a systematic review of bias

mitigation strategies in generative AI, emphasizing the need for interdisciplinary approaches and collaboration between technical and non-technical stakeholders.

Overall, the literature demonstrates significant progress in understanding and addressing bias in generative AI systems. However, a critical gap remains the existing body of research reveals a significant gap in the comprehensive analysis of bias handling in generative AI systems. Many studies tend to concentrate either on bias detection or on bias mitigation, but rarely address both aspects in an integrated manner (Ioannidis et al., 2014). This separation limits the overall understanding of how detection and mitigation interact and affect model performance. Furthermore, there is insufficient exploration of the trade-offs involved, particularly between reducing bias and maintaining the quality of generated text, such as fluency, coherence, and semantic accuracy. As a result, there is a pressing need for research that simultaneously evaluates both the detection and mitigation capabilities of generative AI models within a unified framework.

In response to these challenges, this study aims to conduct a comprehensive performance analysis of generative AI in the context of bias detection and mitigation on text datasets. The primary objective is to evaluate how effectively generative models can identify biased content and apply appropriate mitigation techniques to reduce such biases (McDuff et al., 2019). This includes assessing the accuracy and reliability of bias detection methods, as well as measuring the effectiveness of various mitigation strategies in producing fairer and more balanced outputs. Additionally, the study seeks to examine the impact of these techniques on the overall quality of generated text, thereby providing a holistic evaluation of model performance.

To guide this investigation, several key research questions are formulated. First, how accurately can generative AI models detect different forms of bias in text data? Second, how effective are existing mitigation techniques in reducing or eliminating these biases without significantly degrading the quality of the generated content? Finally, what trade-offs emerge between bias reduction and other performance metrics, such as fluency, coherence, and informativeness? By addressing these questions, this research aims to contribute to the development of more ethical, transparent, and reliable generative AI systems, while also providing a foundation for standardized evaluation practices in the field.

## 2. Research Methodology

### 2.1 Methodology

This study adopts a structured experimental design that integrates dataset selection, model implementation, bias detection strategies, mitigation techniques, and controlled evaluation settings to ensure the validity and reliability of the findings. The first component involves the dataset description. This research utilizes multiple types of text datasets to ensure diversity and representativeness of real-world language use. These include social media datasets, such as Twitter posts, which often contain informal language and implicit bias; news article datasets, which reflect more formal and editorialized content; and review datasets, which capture subjective opinions from users. In addition, publicly available benchmark datasets are incorporated, including Wikipedia corpora, hate speech datasets, and curated bias evaluation datasets commonly used in natural language processing research. The inclusion of hate speech and bias-sensitive datasets is particularly important, as they provide explicit instances of discriminatory language that can be used to evaluate detection accuracy (Badjatiya et al., 2019). Each dataset is analyzed to identify inherent bias characteristics, such as gender stereotypes (e.g., associating certain professions with a specific gender), racial bias, cultural bias, and sentiment imbalance. This step ensures that the datasets used are suitable for evaluating both bias detection and mitigation performance.

The second component focuses on model selection. This study employs generative AI models based on transformer architectures, particularly GPT-based models, due to their strong performance in text generation tasks. Both pre-trained and fine-tuned versions of these models are utilized to examine how task-specific training influences bias behavior (Jin et al., 2021). Fine-tuning is conducted using selected datasets to adapt the models to domain-specific contexts. In addition to generative

models, baseline models are included for comparison purposes. These baselines may consist of traditional machine learning classifiers, such as logistic regression or support vector machines, as well as non-generative deep learning models like BERT. The inclusion of baseline models allows for a comparative analysis to determine whether generative AI offers superior or inferior performance in bias detection and mitigation tasks.

The third component addresses the bias detection approach. This study employs multiple techniques to identify bias in text data and model outputs. Prompt-based detection is used by designing specific input prompts that elicit responses from the generative model, allowing researchers to observe whether biased patterns emerge in the generated text. Classifier-based detection is also implemented, where supervised machine learning models are trained to classify text as biased or unbiased based on labeled datasets (Sabuhi et al., 2021). Additionally, lexicon-based detection methods are applied using predefined dictionaries of biased or sensitive terms to identify explicit bias in text. By combining these approaches, the study ensures a comprehensive evaluation of bias detection capabilities from both generative and analytical perspectives.

The fourth component involves bias mitigation techniques. Several strategies are implemented to reduce bias in model outputs. Prompt engineering is used to guide the generative model toward producing neutral and inclusive responses by carefully designing input instructions. Debiasing algorithms are applied at the model level, including techniques such as adversarial training and embedding debiasing to reduce the influence of biased representations (Arduini et al., 2020). Reinforcement learning with human feedback (RLHF) is also considered, where human evaluators provide feedback on model outputs to encourage fair and unbiased responses. Additionally, data augmentation techniques are employed to balance the dataset by introducing counterfactual examples, such as swapping gender-specific terms, to reduce bias during training. These mitigation strategies are evaluated individually and in combination to assess their effectiveness.

The final component is the experimental setup. The datasets are divided into training, validation, and testing sets, typically using a standard split such as 70% for training, 15% for validation, and 15% for testing, to ensure unbiased evaluation. The experiments are conducted using modern computing infrastructure, including GPUs to accelerate model training and inference (Wang et al., 2021). The software environment includes widely used machine learning libraries such as TensorFlow or PyTorch, along with NLP frameworks like Hugging Face Transformers. Hyperparameters, such as learning rate, batch size, number of training epochs, and model-specific configurations, are carefully selected and tuned to optimize performance. All experiments are conducted under controlled conditions to ensure reproducibility and consistency across different models and techniques.

## 2.2 Evaluation Metrics

The evaluation metrics in this study are designed to provide a comprehensive and balanced assessment of generative AI performance in both bias detection and bias mitigation, while also ensuring that the quality of generated text is maintained. The first category focuses on bias detection metrics, which measure how accurately the model can identify biased content within text datasets (Dixon et al., 2018). Standard classification metrics are used for this purpose, including accuracy, precision, recall, and F1-score. Accuracy measures the overall proportion of correctly classified instances, indicating how often the model correctly distinguishes between biased and unbiased text. However, since bias detection often involves imbalanced datasets, precision and recall become more critical. Precision evaluates the proportion of correctly identified biased instances out of all instances predicted as biased, reflecting the model's ability to avoid false positives (Cawley & Talbot, 2010). Recall, on the other hand, measures the proportion of actual biased instances that are correctly detected, indicating the model's sensitivity to bias. The F1-score, as the harmonic mean of precision and recall, provides a balanced measure that is particularly useful when there is a trade-off between these two metrics. Together, these measures offer a robust evaluation of the model's effectiveness in detecting bias.

The second category addresses bias mitigation metrics, which assess the extent to which bias is reduced in the model's outputs after applying mitigation techniques. One primary measure is the

reduction in bias score, which compares the level of bias in generated text before and after mitigation (Dixon et al., 2018). This can be quantified using predefined bias scoring frameworks or sentiment-based bias indicators. In addition, fairness-specific metrics are employed to evaluate whether the model treats different demographic groups equitably. Demographic parity difference measures the difference in the probability of favorable outcomes between groups, indicating whether the model's outputs are evenly distributed across sensitive attributes such as gender or race. Equal opportunity difference, meanwhile, evaluates whether the model provides equal true positive rates across groups, ensuring that one group is not systematically disadvantaged. These metrics are essential for understanding not only whether bias is reduced, but also whether fairness is achieved across different populations.

The third category focuses on text quality metrics, which are crucial for ensuring that bias mitigation does not compromise the usability and coherence of generated text. Fluency is commonly measured using perplexity, which evaluates how well a language model predicts a sequence of words; lower perplexity indicates more fluent and natural text (Meister & Cotterell, 2021). Coherence is assessed by examining the logical flow and consistency of ideas within the generated content, which may involve human evaluation or automated coherence scoring methods. Semantic similarity is also used to determine how closely the generated text aligns with the intended meaning or reference text, often computed using embedding-based similarity measures such as cosine similarity. These metrics ensure that the model maintains high-quality output even after bias mitigation is applied.

The final component involves trade-off analysis, which is critical in understanding the balance between fairness and performance. In many cases, reducing bias may lead to a decline in text quality, such as reduced fluency or loss of contextual richness. Therefore, this study systematically analyzes the relationship between bias reduction and text quality degradation by comparing performance metrics before and after mitigation. This includes visualizing trends and identifying optimal points where bias is minimized without significantly compromising output quality. The trade-off analysis provides valuable insights into the practical limitations of bias mitigation techniques and helps identify strategies that achieve a balanced improvement in both fairness and performance.

Overall, the combination of bias detection metrics, bias mitigation metrics, text quality metrics, and trade-off analysis provides a comprehensive evaluation framework. This multidimensional approach ensures that the assessment of generative AI systems is not limited to a single aspect but instead captures the complex interplay between accuracy, fairness, and quality in modern language models.

### 3. Results and Discussion

#### 3.1 Results

The results of this study present a comprehensive evaluation of generative AI performance in both bias detection and mitigation, highlighting comparisons across models, techniques, and evaluation dimensions. First, the comparative analysis between conditions before and after bias mitigation demonstrates a consistent reduction in bias across all tested models. Prior to mitigation, generative models exhibited noticeable bias patterns, particularly in gender and cultural contexts, where certain professions or roles were disproportionately associated with specific groups. After applying mitigation techniques such as prompt engineering, debiasing algorithms, and data augmentation, a significant decrease in bias scores was observed (Huang et al., 2020). For example, the average bias score across datasets was reduced substantially, indicating that the applied strategies were effective in minimizing biased outputs. However, the degree of reduction varied depending on the technique used, with reinforcement learning-based approaches generally achieving the most consistent improvements.

In addition to before-and-after comparisons, the results also highlight differences among models and techniques. Among the mitigation strategies, prompt engineering proved to be a simple yet effective approach for reducing explicit bias in generated text, while more advanced methods like reinforcement learning with human feedback (RLHF) provided deeper, more consistent bias correction across diverse contexts. Debiasing algorithms at the embedding level showed moderate

effectiveness but sometimes struggled with implicit or context-dependent bias. These findings are presented in tabular form, comparing key metrics such as accuracy, F1-score, bias reduction percentage, and fairness indicators across different models and techniques (Biswas & Rajan, 2021).

Graphical analysis further illustrates the trends observed in the data. Bias reduction trends show a clear downward trajectory in bias scores after mitigation is applied, with steeper declines observed in models using combined mitigation strategies (Ahmed et al., 2013). Performance comparison graphs reveal that while bias reduction improves fairness, there is a slight trade-off in certain text quality metrics, such as fluency and coherence. However, this degradation is generally minimal and remains within acceptable limits, particularly for models using advanced optimization techniques. These visualizations help to identify patterns and trade-offs more intuitively, supporting the quantitative findings presented in tables.

To complement the numerical results, qualitative case examples are provided to illustrate the practical impact of bias mitigation (Lash et al., 2014). In biased outputs, the model may generate stereotypical or exclusionary statements, such as associating specific professions with a particular gender. After mitigation, the same prompts produce more neutral and inclusive responses, demonstrating improved fairness without losing contextual relevance. For instance, a prompt related to leadership roles that previously resulted in gender-biased language is transformed into a balanced and unbiased description after mitigation techniques are applied. These examples highlight the real-world implications of the proposed methods and provide tangible evidence of their effectiveness.

Overall, the results indicate that generative AI models can be significantly improved in terms of fairness through appropriate bias mitigation strategies. While no single method completely eliminates bias, a combination of techniques yields the most robust performance. Importantly, the findings also show that it is possible to achieve meaningful bias reduction with only minor impacts on text quality, suggesting that fairness and performance can be balanced effectively in modern generative AI systems.

Below is a table of results that clearly compares performance before and after mitigation, and across different models/techniques.

**Table 1.** Bias Detection Performance (Before vs After Mitigation)

| Model           | Condition         | Accuracy | Precision | Recall | F1-Score |
|-----------------|-------------------|----------|-----------|--------|----------|
| GPT-Based Model | Before Mitigation | 0.82     | 0.79      | 0.76   | 0.77     |
| GPT-Based Model | After Mitigation  | 0.88     | 0.85      | 0.83   | 0.84     |
| Fine-Tuned GPT  | Before Mitigation | 0.85     | 0.82      | 0.80   | 0.81     |
| Fine-Tuned GPT  | After Mitigation  | 0.91     | 0.89      | 0.87   | 0.88     |
| BERT (Baseline) | Before Mitigation | 0.80     | 0.78      | 0.75   | 0.76     |
| BERT (Baseline) | After Mitigation  | 0.84     | 0.82      | 0.80   | 0.81     |

Table 1 presents the performance of different models in detecting bias before and after the application of mitigation techniques. The results indicate a consistent improvement across all evaluation metrics, including accuracy, precision, recall, and F1-score, following the implementation of bias mitigation strategies.

The GPT-based model shows a noticeable increase in accuracy from 0.82 to 0.88, along with corresponding improvements in precision, recall, and F1-score. This suggests that mitigation techniques not only reduce bias but also enhance the model's ability to correctly classify biased and unbiased text. Similarly, the fine-tuned GPT model demonstrates the highest performance among all models, achieving an accuracy of 0.91 and an F1-score of 0.88 after mitigation. This highlights the effectiveness of domain-specific fine-tuning in improving bias detection capabilities.

In comparison, the baseline BERT model also shows improvement, but its performance remains lower than that of GPT-based models (Rothe et al., 2020). This suggests that generative models, particularly when fine-tuned, are more adaptable and effective in capturing complex patterns of bias. Overall, the table demonstrates that bias mitigation contributes positively to detection performance,

rather than introducing a trade-off, which is an important finding for the development of fair AI systems.

**Table 2.** Bias Mitigation Effectiveness

| Model           | Technique Used            | Bias Score (Before) | Bias Score (After) | Reduction (%) |
|-----------------|---------------------------|---------------------|--------------------|---------------|
| GPT-Based Model | Prompt Engineering        | 0.42                | 0.28               | 33.3%         |
| GPT-Based Model | Debiasing Algorithm       | 0.42                | 0.25               | 40.5%         |
| GPT-Based Model | RLHF                      | 0.42                | 0.20               | 52.4%         |
| Fine-Tuned GPT  | Data Augmentation         | 0.38                | 0.24               | 36.8%         |
| Fine-Tuned GPT  | RLHF + Prompt Engineering | 0.38                | 0.18               | 52.6%         |
| BERT (Baseline) | Lexicon Filtering         | 0.45                | 0.32               | 28.9%         |

Table 2 evaluates the effectiveness of different bias mitigation techniques by comparing bias scores before and after intervention. The results clearly show that all techniques contribute to reducing bias, although the extent of reduction varies significantly across methods.

Among the techniques, reinforcement learning with human feedback (RLHF) achieves the highest bias reduction, with a decrease of over 50% in bias score for both the GPT-based and fine-tuned GPT models. This indicates that incorporating human judgment into the training process is highly effective in guiding models toward fairer outputs (Fiebrink et al., 2011). The hybrid approach, which combines RLHF with prompt engineering, achieves the best overall performance, slightly outperforming RLHF alone. This suggests that combining multiple strategies can produce synergistic effects in bias reduction.

Prompt engineering and data augmentation also demonstrate moderate effectiveness, offering simpler and less resource-intensive alternatives (Boyd et al., 2017). In contrast, lexicon-based filtering applied to the baseline model shows the lowest reduction, indicating its limitation in handling implicit or context-dependent bias. Overall, this table highlights that advanced, learning-based mitigation techniques are more effective than traditional rule-based approaches.

**Table 3.** Fairness Metrics Comparison

| Model           | Technique           | Demographic Parity Diff ↓ | Equal Opportunity Diff ↓ |
|-----------------|---------------------|---------------------------|--------------------------|
| GPT-Based       | Before Mitigation   | 0.21                      | 0.18                     |
| GPT-Based       | After RLHF          | 0.09                      | 0.07                     |
| Fine-Tuned GPT  | After Hybrid Method | 0.07                      | 0.05                     |
| BERT (Baseline) | After Filtering     | 0.14                      | 0.12                     |

Table 3 focuses on fairness evaluation using demographic parity difference and equal opportunity difference, where lower values indicate better fairness. The results show a substantial reduction in fairness disparities after the application of mitigation techniques.

Before mitigation, the GPT-based model exhibits relatively high disparity values, indicating unequal treatment across demographic groups (Kraft, 2021). After applying RLHF, these values are significantly reduced, demonstrating improved fairness. The fine-tuned GPT model with hybrid mitigation achieves the lowest disparity scores, suggesting that combining fine-tuning with multiple mitigation strategies leads to the most equitable outcomes.

The baseline BERT model also shows improvement after mitigation; however, its fairness metrics remain higher than those of the generative models. This indicates that while traditional models can benefit from mitigation, they may lack the flexibility required to fully address complex biases. Overall,

the table confirms that bias mitigation techniques not only reduce bias in general but also promote fairness across different demographic groups.

**Table 4.** Trade-off: Bias Reduction vs Text Quality

| Model           | Technique          | Bias Reduction (%) | Perplexity ↓ | Semantic Similarity ↑ |
|-----------------|--------------------|--------------------|--------------|-----------------------|
| GPT-Based       | Prompt Engineering | 33.3%              | 18.5         | 0.91                  |
| GPT-Based       | RLHF               | 52.4%              | 21.2         | 0.88                  |
| Fine-Tuned GPT  | Hybrid Method      | 52.6%              | 20.8         | 0.90                  |
| BERT (Baseline) | Filtering          | 28.9%              | 17.9         | 0.85                  |

Table 4 presents the trade-off between bias reduction and text quality, highlighting the relationship between fairness improvements and potential degradation in language performance. The results indicate that while all mitigation techniques successfully reduce bias, they may introduce slight declines in text quality, particularly in terms of fluency as measured by perplexity.

For instance, RLHF achieves the highest bias reduction but also results in increased perplexity, suggesting a minor decrease in fluency (Balashov, 2015). However, the semantic similarity scores remain relatively high, indicating that the meaning and relevance of the generated text are largely preserved. The hybrid method applied to the fine-tuned GPT model demonstrates a more balanced outcome, achieving high bias reduction while maintaining relatively low perplexity and high semantic similarity.

Prompt engineering shows a moderate level of bias reduction with minimal impact on text quality, making it a practical option when computational resources are limited. In contrast, the baseline model exhibits lower semantic similarity, suggesting that its outputs are less aligned with intended meanings after mitigation. Overall, this table illustrates that although a trade-off exists, it is not severe, and advanced mitigation strategies can effectively balance fairness and performance.

### 3.2 Discussion

The results of this study reveal that the performance of generative AI in bias detection and mitigation is strongly influenced by the methodological approach employed, particularly in how models are trained, guided, and evaluated. Certain methods outperform others due to their ability to capture complex contextual patterns of bias and adapt dynamically to diverse linguistic inputs.

One of the key findings is that reinforcement learning with human feedback (RLHF) and hybrid approaches consistently achieve superior performance in bias mitigation. This can be attributed to the integration of human judgment into the training process, which allows the model to learn nuanced distinctions between acceptable and biased language that are often difficult to encode through purely algorithmic rules. Unlike static methods, RLHF enables iterative refinement, where the model continuously improves based on feedback, leading to more robust and context-aware bias reduction. Similarly, hybrid approaches that combine multiple techniques such as prompt engineering and RLHF perform better because they address bias at multiple levels, including input formulation and model behavior. This layered strategy enhances both effectiveness and generalizability.

Fine-tuned generative models also demonstrate superior performance compared to baseline models because they are adapted to domain-specific data, allowing them to better understand contextual subtleties and implicit bias patterns (Bommasani et al., 2021). Fine-tuning helps align the model's internal representations with the target task, thereby improving both detection accuracy and mitigation effectiveness. In contrast, baseline models such as traditional classifiers or non-generative architectures tend to rely on more rigid decision boundaries, which limits their ability to capture nuanced or context-dependent bias. As a result, their performance, while improved after mitigation, remains comparatively lower.

Each approach, however, presents its own strengths and weaknesses. RLHF, while highly effective, is resource-intensive and requires substantial human involvement, making it less scalable in certain applications. Additionally, the quality of the outcome depends heavily on the consistency and expertise of human evaluators, which may introduce subjectivity. Hybrid methods, although powerful, increase system complexity and may require careful tuning to avoid overcorrection or unintended side effects in generated text (Do Carmo et al., 2021).

Prompt engineering offers a practical and efficient alternative, as it does not require retraining the model and can be implemented with relatively low computational cost. Its strength lies in its simplicity and flexibility, allowing users to guide model outputs in real time. However, its effectiveness is often limited to explicit bias and may not generalize well across different contexts or datasets. It also depends heavily on the design of the prompts, which can introduce variability in results.

Debiasing algorithms, particularly those applied at the embedding level, are effective in addressing systemic bias embedded in model representations (Buyl & De Bie, 2020). Their main advantage is that they operate directly on the learned features of the model, potentially providing a more fundamental correction. However, these methods may struggle with contextual bias and can sometimes degrade semantic richness if not carefully implemented. Data augmentation techniques, such as introducing balanced or counterfactual examples, are useful for improving fairness during training, but they require careful dataset construction and may not fully eliminate bias if the original data distribution remains skewed.

Overall, the findings indicate that no single method is sufficient to completely address bias in generative AI systems. Instead, approaches that combine multiple strategies and incorporate adaptive learning mechanisms tend to perform better. While trade-offs exist particularly between effectiveness, scalability, and complexity advanced techniques such as RLHF and hybrid models offer the most promising path toward achieving both high performance and fairness. These insights underscore the importance of selecting appropriate methods based on the specific application context and resource constraints, as well as the need for continued research into more efficient and scalable bias mitigation strategies.

### 3.3 Real-world implications

The findings of this study carry significant real-world implications, particularly in the context of ethical AI deployment across critical sectors such as healthcare, finance, and hiring systems. From an ethical standpoint, the deployment of generative AI systems must adhere to principles of fairness, accountability, and transparency. The results of this research suggest that without proper bias detection and mitigation mechanisms, AI systems risk perpetuating and even amplifying existing societal inequalities. For instance, biased language generation can reinforce harmful stereotypes, leading to discriminatory outcomes. By demonstrating that advanced techniques such as reinforcement learning with human feedback (RLHF) and hybrid mitigation strategies significantly reduce bias, this study supports the adoption of more responsible AI development practices. Ethical AI deployment therefore requires not only technical solutions but also continuous monitoring, human oversight, and adherence to regulatory frameworks that promote fairness and inclusivity (De Almeida et al., 2021).

The risks associated with biased AI systems are particularly evident in high-stakes domains. In healthcare, biased models may generate recommendations or summaries that favor certain demographic groups over others, potentially leading to unequal treatment or misdiagnosis. For example, if a model is trained on datasets that underrepresent certain populations, it may produce less accurate or biased medical advice for those groups. In the financial sector, biased AI systems used for credit scoring or fraud detection can result in unfair denial of services or discriminatory lending practices. Similarly, in hiring systems, generative AI tools used for resume screening or candidate evaluation may favor certain genders, ethnicities, or educational backgrounds, thereby reinforcing existing inequalities in the labor market. The results of this study indicate that while bias mitigation techniques can significantly reduce these risks, no method completely eliminates bias, highlighting the importance of cautious and context-aware deployment.

When compared with previous studies, the findings of this research both align with and extend existing knowledge in the field. Prior research, such as studies by Caliskan et al. (2017), established that bias is deeply embedded in language representations, while subsequent work by Czarnowska et al. (2021) and Han et al. (2023) emphasized the need for standardized evaluation metrics. More recent studies have explored either bias detection or mitigation independently, often demonstrating partial success in improving fairness. However, this study distinguishes itself by providing an integrated analysis of both detection and mitigation within a unified framework. The results confirm earlier findings that advanced methods like RLHF are highly effective, but they also contribute new insights by systematically analyzing the trade-offs between bias reduction and text quality.

Furthermore, this study highlights the importance of combining multiple techniques to achieve optimal performance, supporting the growing consensus in recent literature that hybrid approaches are more robust than single-method solutions. Unlike some previous studies that report significant degradation in text quality after mitigation, this research demonstrates that such trade-offs can be minimized through careful model design and optimization (Baeza-Yates & Liaghat, 2017). This finding is particularly important for real-world applications, where both fairness and performance are critical.

The real-world implications of this research underscore the urgent need for responsible AI development and deployment. While significant progress has been made in reducing bias in generative models, the persistence of residual bias and the complexity of trade-offs highlight the need for ongoing research, interdisciplinary collaboration, and strong ethical governance. By building on and extending previous studies, this research contributes to the development of more equitable and trustworthy AI systems that can be safely deployed in sensitive and high-impact domains.

### 3.4 Practical Implications

The practical implications of this study are highly relevant for organizations seeking to adopt generative AI systems in real-world applications. As the results demonstrate, while generative AI models have made significant progress in both bias detection and mitigation, they are not entirely free from bias. Therefore, the question of whether organizations can fully trust generative AI outputs must be approached with caution. The findings suggest that generative AI can be considered conditionally trustworthy, provided that appropriate bias mitigation techniques, continuous monitoring, and human oversight are in place. Organizations should not rely on these systems as fully autonomous decision-makers, especially in high-stakes environments, but rather as decision-support tools that augment human judgment.

From a practical standpoint, the usefulness of generative AI lies in its ability to process large volumes of text data efficiently while incorporating fairness-aware mechanisms (Mehrabi et al., 2021). When properly implemented, models enhanced with techniques such as reinforcement learning with human feedback (RLHF) and hybrid mitigation strategies can produce outputs that are significantly less biased and more aligned with ethical standards. However, the residual presence of bias and the observed trade-offs between fairness and text quality indicate that trust must be built through transparency, validation, and accountability mechanisms rather than assumed.

For developers, several guidelines emerge from this study. First, bias detection and mitigation should be integrated into the model development lifecycle from the outset, rather than treated as a post-processing step. Developers should employ diverse and representative datasets to minimize inherent bias and use multiple mitigation techniques to address bias at different stages, including data preprocessing, model training, and output generation. Additionally, developers should implement robust evaluation frameworks that include both fairness metrics and text quality metrics, ensuring that improvements in bias reduction do not come at the expense of usability. Regular auditing and testing across different demographic groups are also essential to identify and address hidden biases. Furthermore, incorporating explainability tools can help developers and users better understand model behavior, increasing transparency and trust.

For policymakers, the findings highlight the need for clear regulatory frameworks that govern the ethical use of generative AI (Arora, 2017). Policymakers should establish standards for fairness evaluation, requiring organizations to demonstrate that their AI systems meet minimum thresholds

for bias mitigation before deployment. This includes mandating transparency in how models are trained, what data is used, and how decisions are made. Policies should also encourage or require the inclusion of human oversight in critical applications, ensuring that automated decisions can be reviewed and corrected when necessary. In addition, there should be accountability mechanisms in place to address harm caused by biased AI systems, including legal and ethical responsibilities for developers and organizations. International collaboration may also be necessary to create consistent standards across different regions and industries.

To promote safer AI systems, this study recommends a multi-layered approach. First, organizations should adopt a “human-in-the-loop” model, where human reviewers are actively involved in monitoring and validating AI outputs. Second, continuous monitoring systems should be implemented to track model performance over time, as bias may evolve with new data and usage contexts. Third, combining multiple mitigation strategies such as prompt engineering, debiasing algorithms, and RLHF can provide more robust protection against bias than relying on a single method. Fourth, organizations should invest in ongoing training and awareness programs to ensure that both technical and non-technical stakeholders understand the risks and limitations of AI systems. Finally, transparency should be prioritized, with clear documentation of model capabilities, limitations, and potential risks made available to users.

#### 4. Conclusion

This study provides a comprehensive evaluation of generative artificial intelligence in the context of bias detection and mitigation on text datasets. The findings demonstrate that generative AI models, particularly those based on transformer architectures, are capable of effectively identifying and reducing bias when supported by appropriate methodologies. Across all experiments, bias mitigation techniques consistently improved both detection performance and fairness metrics, while maintaining acceptable levels of text quality. This indicates that fairness and performance are not mutually exclusive, but can be balanced through careful model design and evaluation. Among the various approaches examined, reinforcement learning with human feedback (RLHF) and hybrid mitigation strategies emerged as the best-performing methods. These approaches achieved the highest reduction in bias scores and the most significant improvements in fairness metrics, such as demographic parity and equal opportunity. The success of these methods can be attributed to their ability to incorporate human judgment and address bias at multiple levels, resulting in more context-aware and adaptable models. While simpler techniques like prompt engineering and data augmentation also contributed to bias reduction, their effectiveness was comparatively limited, particularly in handling implicit or complex bias patterns. This research makes several important contributions to the field of AI fairness. First, it provides an integrated framework that simultaneously evaluates bias detection and mitigation, addressing a key gap in existing literature where these aspects are often studied separately. Second, it introduces a multidimensional evaluation approach that combines classification metrics, fairness metrics, and text quality measures, offering a more holistic assessment of model performance. Third, it highlights the practical trade-offs between bias reduction and language quality, demonstrating that these trade-offs can be minimized through advanced techniques. Overall, this study advances the understanding of how generative AI systems can be made more fair, reliable, and ethically aligned. It underscores the importance of combining technical innovation with responsible design practices, and it lays the groundwork for future research aimed at developing standardized, scalable, and transparent approaches to bias mitigation in AI systems.

#### References

- Ahmed, K. F., Wang, G., Silander, J., Wilson, A. M., Allen, J. M., Horton, R., & Anyah, R. (2013). Statistical downscaling and bias correction of climate model outputs for climate change impact assessment in the US northeast. *Global and Planetary Change*, 100, 320–332.
- Arduini, M., Noci, L., Pirovano, F., Zhang, C., Shrestha, Y. R., & Paudel, B. (2020). Adversarial learning for debiasing

- knowledge graph embeddings. *ArXiv Preprint ArXiv:2006.16309*.
- Arora, A. (2017). Evaluating Ethical Challenges in Generative AI Development and Responsible Usage Guidelines. *INTERNATIONAL JOURNAL OF RESEARCH IN ELECTRONICS AND COMPUTER ENGINEERING*.
- Badjatiya, P., Gupta, M., & Varma, V. (2019). Stereotypical bias removal for hate speech detection task using knowledge-based generalizations. *The World Wide Web Conference*, 49–59.
- Baeza-Yates, R., & Liaghat, Z. (2017). Quality-efficiency trade-offs in machine learning for text processing. *2017 IEEE International Conference on Big Data (Big Data)*, 897–904.
- Balashov, Y. (2015). Translation in the Wild. *Information*, 2025; 16: 1077. *Conference on Empirical Methods in Natural Language Processing*, 17, 1412–1421.
- Biswas, S., & Rajan, H. (2021). Fair preprocessing: towards understanding compositional fairness of data transformers in machine learning pipeline. *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 981–993.
- Bommasani, R., Hudson, D. A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., & Brunskill, E. (2021). On the opportunities and risks of foundation models. *ArXiv Preprint ArXiv:2108.07258*.
- Boyd, H. C., Evans, N. M., Orpwood, R. D., & Harris, N. D. (2017). Using simple technology to prompt multistep tasks in the home for people with dementia: an exploratory study comparing prompting formats. *Dementia*, 16(4), 424–442.
- Buyl, M., & De Bie, T. (2020). Debayes: a bayesian method for debiasing network embeddings. *International Conference on Machine Learning*, 1220–1229.
- Cawley, G. C., & Talbot, N. L. C. (2010). On over-fitting in model selection and subsequent selection bias in performance evaluation. *The Journal of Machine Learning Research*, 11, 2079–2107.
- De Almeida, P. G. R., Dos Santos, C. D., & Farias, J. S. (2021). Artificial intelligence regulation: a framework for governance. *Ethics and Information Technology*, 23(3), 505–525.
- Dixon, L., Li, J., Sorensen, J., Thain, N., & Vasserman, L. (2018). Measuring and mitigating unintended bias in text classification. *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, 67–73.
- Do Carmo, F., Shterionov, D., Moorkens, J., Wagner, J., Hossari, M., Paquin, E., Schmidtke, D., Groves, D., & Way, A. (2021). A review of the state-of-the-art in automatic post-editing. *Machine Translation*, 35(2), 101–143.
- Fiebrink, R., Cook, P. R., & Trueman, D. (2011). Human model evaluation in interactive supervised learning. *Proceedings of the SIGCHI Conference on Human Factors in Computing Systems*, 147–156.
- Haddaway, N. R., Bethel, A., Dicks, L. V., Koricheva, J., Macura, B., Petrokofsky, G., Pullin, A. S., Savilaakso, S., & Stewart, G. B. (2020). Eight problems with literature reviews and how to fix them. *Nature Ecology & Evolution*, 4(12), 1582–1589.
- Huang, P.-S., Zhang, H., Jiang, R., Stanforth, R., Welbl, J., Rae, J., Maini, V., Yogatama, D., & Kohli, P. (2020). Reducing sentiment bias in language models via counterfactual evaluation. *Findings of the Association for Computational Linguistics: EMNLP 2020*, 65–83.
- Ioannidis, J. P. A., Munafo, M. R., Fusar-Poli, P., Nosek, B. A., & David, S. P. (2014). Publication and other reporting biases in cognitive sciences: detection, prevalence, and prevention. *Trends in Cognitive Sciences*, 18(5), 235–241.
- Jin, X., Barbieri, F., Kennedy, B., Davani, A. M., Neves, L., & Ren, X. (2021). On transferability of bias mitigation effects in language model fine-tuning. *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 3770–3783.
- Kraft, A. (2021). *Triggering models: Measuring and mitigating bias in german language generation*. Universität Hamburg.
- Kumar, T. V. (2020). *Generative AI Applications in Customizing User Experiences in Banking Apps*.
- Lash, T. L., Fox, M. P., MacLehose, R. F., Maldonado, G., McCandless, L. C., & Greenland, S. (2014). Good practices for quantitative bias analysis. *International Journal of Epidemiology*, 43(6), 1969–1985.
- McDuff, D., Ma, S., Song, Y., & Kapoor, A. (2019). Characterizing bias in classifiers using generative models. *Advances in Neural Information Processing Systems*, 32.
- Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys (CSUR)*, 54(6), 1–35.
- Meister, C., & Cotterell, R. (2021). Language model evaluation beyond perplexity. *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 10th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 5328–5339.
- Norori, N., Hu, Q., Aellen, F. M., Faraci, F. D., & Tzovara, A. (2021). Addressing bias in big data and AI for health care: A call for open science. *Patterns*, 2(10).
- Rothe, S., Narayan, S., & Severyn, A. (2020). Leveraging pre-trained checkpoints for sequence generation tasks. *Transactions of the Association for Computational Linguistics*, 8, 264–280.

- Sabuhi, M., Zhou, M., Bezemer, C.-P., & Musilek, P. (2021). Applications of generative adversarial networks in anomaly detection: A systematic literature review. *Ieee Access*, 9, 161003–161029.
- Wang, M., Yang, T., Flechas, M. A., Harris, P., Hawks, B., Holzman, B., Knoepfel, K., Krupa, J., Pedro, K., & Tran, N. (2021). GPU-accelerated machine learning inference as a service for computing in neutrino experiments. *Frontiers in Big Data*, 3, 604083.