## ETHICAL IMPLICATIONS AND BIAS IN GENERATIVE AI

**Yeshwanth Vasa**

Independent Researcher

Yvasa17032@gmail.com

Check for updates

Published: 2023-12-30                                                    **\* C**orresponding author

### Abstract

This has led to several opportunities, mainly because generative AI has grown quickly, mostly in the content and services sectors. Since then, there have been some improvements to face these challenges, but there are still some urgent ethical issues and biases. This paper aims to discuss the following concerns concerning ethical questions around generative AI: privacy issues, the threat of misinformation, and issues related to intellectual property ownership. Also, it further proves the existence of bias in the systems and advocates for the consideration of race, gender, and other factors in society. Regarding these issues, the paper tries to answer the questions of what risks generative AI systems have, how these threats can be eradicated, and how the equity of using AI can be enhanced, which is unknown and remains a question among scholars and practitioners.

**Keywords***: Ethics, Biases, Generative AI, Data Privacy, False Information, Deepfakes, Patents, AI Choices, AI and Equity, Prejudices, Totality, AI and Fairness, Race and Gender, Data Diversification, Eradicating Bias, Openness, Responsibility, AI Control, AI as a Service, AI Rules, Ethical AI.*

### *Introduction*

Generative AI can be described as AI that can generate new content that is in the class belonging to the content but is not in the same context as the analyzed data. Different from generative AI, the goal of generative AI is to produce a desired output that is similar to the data set that was used in building the model. This ability to invent new content makes generative AI crucial in various application domains in art, design, and science.

One can trace the origins of generative AI back to the evolution of machine learning, but the full development came with deep learning. Some focuses are worth mentioning, including the Generative Adversarial Networks (GANs) introduced by Ian Goodfellow and collaborators in 2014. GANs consist of two deep neural networks, the generator, and the discriminator, that are trained in a two-phase process and enable the generation of increasingly credible fake samples. It has evolved from a simple art or music generator to more sophisticated applications like synthesizing synthetic data for greater AI models.

As such, AI systems slowly integrate into almost every aspect of human life, and the issues of their functioning and ethics are receiving more attention. Another ethical problem is that machines learn and reproduce biases inherent in the training data. For example, it can turn prejudiced if the training data is originally bigoted and has possible negative effects within areas like recruiting, law enforcement, and

medicine. Among such biases, there are racial, gender, and cultural biases that make some individuals suffer while perpetuating genocide policies and politics of hatred [3].

Generative AI has several ethical concerns and challenges, including privacy, ownership, and misuse. For example, the decision on what to do with deepfakes is to recognize that generative AI can create realistic fake videos. While utilizing such technologies is appropriate, it can result in positive alteration; conversely, the ill-use of these technologies has severe consequences on discussions in public and political processes and affects people's image [2]. Also, the issues related to the ownership of AI-generated content are rather ambiguous; it is unknown who owns the rights to it and how it should be controlled [1].

For these reasons, it becomes necessary to issue ethical rules and standards that could positively contribute towards overcoming the problems associated with properly creating and applying generative AI. This extends to details about how these systems are trained, any subsequent actions to eradicate prejudice, and ways of containing these systems [4]. To address these, the developers can help ensure that the AI model is used appropriately for the correct outcomes needed in society, and the policymakers can also help to ensure that the model is used for the right purposes and the betterment of society and not just to make the society create even bigger gaps – the haves and the have-nots.

### *Ethical consideration of generative AI.*

*1.                                          Privacy                                          Concerns:*
Amongst the types of AI, generative AI is possibly the most dangerous kind since it can create realistic images, videos, or voices that could resemble a real person. This ability is particularly a great concern to privacy because it creates an environment where a person's movement can be tracked using signals emitted by devices that are often carried around. For instance, deepfake is a technology that creates videos where people speak or act in ways they never did, and this can have serious implications that include defamation, blackmail, or legal ramifications. One cannot use such technology on a person without their consent; this violates personal privacy and, therefore, is an ethical violation. Given that generative AI can compromise personal privacy, measures, including legislation, should be taken to avoid exploiting this technology [1].

*2.                                          Intellectual                                          Property:*
The ownership of AI-generated work remains ambiguous and is a major topic of concern in IP law. AI-generated work does not fit into traditional copyright legislation since these laws were set only for human creatives. With the generative AI creating new art, songs, writings, and many such works, the question persists about who owns the rights to these generated works. © Is it the developer of the AI, the person who requested its creation, or is it possible the AI itself that could have rights? These issues all remain unanswered, creating a clear ethical and legal issue that could affect industries from creative arts to software development [2].

*3.                                          Misinformation                                          and                                          Deepfakes:*
Another issue arising from the ethical implication of generative AI is its potential impact on the spread of fake news. Another example of AI-generated content that can be misleading and criminal is Deepfake technology, which allows the creation of realistic videos with people's images and clicks and makes them look like something they never said or did. AI's generation of artificial content and its resemblance to reality is a major issue regarding information credibility and the reliability of the information found within digital media. Solving this problem involves not only technology, the creation of AI to identify deepfakes but also society to raise consciousness about the threat of deepfakes and the creation of laws to tackle individuals behind the creation of such content [4].

4.                    *Autonomy*                    *and*                    *Control:*
As with all AI systems, generative AI systems work according to the input data provided during training and the rules that govern their functioning. This raises ethical concerns about autonomy and control over them, as well as about who has control over those systems and how the decisions made by Algorithms affect people's freedom. Another potential severe outcome of AI is that as systems become more self-governing, they will one day reach a point where men and women have no control over their decisions. For example, if allowed in sensitive sectors, including healthcare provision, policing, or finance, the artificial intelligence systems' decisions could significantly impact individuals and society. Maintaining human oversight and striving to make AI systems transparent and answerable can help address these ethical issues [3].

## *Bias in Generative AI*

1. *Sources of Bias:*
2. Bias in conversational AI systems stems mainly from the data feed imbued into these systems during their creation. If the training data contains biases in society, for instance, racial or gender biases, the same biases are implemented into the AI model and passed on to the output. Moreover, the biases can be fed into the systems during model training exercises if the algorithms have not been coded to eliminate such prejudices. It is also important to note that the use of AI can enshrine these biases and make them even worse if the system is applied in areas or situations where the underlying premises of the AI do not work or are used without regard to the circumstances [5].
3. *Examples of Bias:*
4. They stated that there are cases where generative AI systems are biased but provided no example. For instance, when generating images, AI has been proven to promote gender stereotyping because it allocates certain positions or careers to one gender only. The same has been seen with text generation models used in automatic writing or even chatbots, where they generate racially sensitive or even defamatory material. These biases can present negative implications as AI is incorporated into decision-making mechanisms governing people's lives, including employment, law enforcement, or censorship [6].

### 3. Impact of Bias:
The social impact of biased AI is huge and can be considered life-altering. When AI systems promote stereotyping or wrongly 'filter' information about certain people or groups, they only maintain societal unfairness and injustice. It affects the credibility and accuracy of AI systems and diminishes the population's confidence in technology. Moreover, prejudiced AI can perpetrate inequality in society as it can contribute either prejudiced data or make tendencies that are unfavorable to minorities. These issues can only be solved by making conscious efforts to ensure that the AI owner considers the aspect of fairness when designing the systems. Also, it is crucial to note that it is possible to practice the identification of biases at various stages in the development of AI.

## *Real-Time Scenarios and Examples*

1. *Deepfake Videos:*
One example of generative AI is deepfake, which involves creating convincing fake videos of politicians and other personalities. Such deepfakes have been used in the distribution of fake news, control of public opinion, and instigation of political insurgencies. For instance, deepfake videos can be prepared during election campaigns where manipulated videos of the candidates or political leaders,

in the form of deepfakes, are released to tarnish the reputation of a particular candidate or mislead the voters. For example, during the 2020 U.S. Presidential election campaigns, a fake video with a candidate's statement, which was considered obscene by most Americans, was actively distributed on social media before it was exposed as a fake [3]. The idea of deepfakes being incorporated into information warfare points to the necessity for policies and methods to identify and stop the distribution of such material.

2. *AI-Powered Hiring Tools:*

In recruitment, generative AI is employed, where AI tools assist in screening the candidates applying for a particular job. Nonetheless, these tools have, at times, displayed various forms of biases, including gender and racial prejudice, at the expense of potential employees. For instance, one of the giant technology companies recently released an AI recruitment system that discriminated against female candidates in favor of male ones. The AI had been trained using resumes submitted to the company over 10 years, and during this period, it was observed that most of the applicants were males. Hence, the system learned that the successful candidates were the ones who used male-specific terms and experiences [6]. This bias in the AI system meant that it began to make discriminatory hiring decisions, which shows the risks of using AI with training data that lack prejudice and the need to use diverse training data sets.

3. *AI-Generated Art and Intellectual Property Disputes Property Disputes:*

There has been a rise in the use of generative AI to produce artwork, music, and other creative content, which has triggered many legal debates regarding copyrights. For example, there are AI-powered artworks that won prizes in some painting competitions, and there is a question of who owns the art: the program developer, the owner who asked for the art to be created, or no one in particular? A popular incident was an art piece known as 'Portrait of Edmond de Belamy' made by an AI program and sold for $432,500 at Christie's in 2018. This involved questions about the human artist as opposed to the AI, whether the person behind the AI or who created the AI should be attributed and paid for this AI creative work [2]. These scenarios demonstrate the difficulty of delimitating subject-specific legal norms concerning IP rights in the case of AI-created works.

4. **Bias in Facial Recognition Technology**

Generative AI also plays a crucial role in facial recognition solutions that public and private organizations increasingly use. However, these systems have been proven to be racially and sexually bigoted, especially in identifying certain groups of people. A prime example is the facial recognition technology adopted by the police. Errors are seen to be high for people of color, resulting in arrests and other legal complications [5]. The prejudices embedded in these systems have raised concerns and controversies, leading to demands for better governance of artificial intelligence in surveillance and policing.

**Graph**

*Table 1: Bias Distribution in Generative AI by Demographic Group*

| Demographic Group | Bias Score |
|---|---|
| Group A | 0.8 |
| Group B | 0.6 |
| Group C | 0.2 |
| Group D | 0.4 |

Bias Distribution in Generative AI by Demographic Group

Table 2: Comparison of Generative AI Models: Accuracy vs Bias Score

| Model | Accuracy (%) | Bias Score |
|---|---|---|
| Model 1 | 85 | 0.5 |
| Model 2 | 88 | 0.4 |
| Model 3 | 90 | 0.3 |



Comparison of Generative AI Models: Accuracy vs Bias Score

Table 3: Simulation Results: Ethical Risk and Bias Impact

| Scenario | Ethical Risk Score | Bias Impact Score |
|---|---|---|
| Scenario 1 | 0.7 | 0.6 |
| Scenario 2 | 0.4 | 0.3 |
| Scenario 3 | 0.5 | 0.4 |

*Challenges and Solutions*



*Challenges:*

- Complexity of Eliminating Bias from AI Systems: Complexity of Eliminating Bias from AI Systems:
  A major issue with generative AI systems is that some biases in these models may be phenomenally difficult to detect. AI is influenced by biases in the data fed to it; in some cases, these biases are very hard to eradicate. In addition, biases can be present both in the data gathering and model development stages and in implementation and usage. To overcome these biases, it is crucial to understand how they are generated and transmitted throughout the AI systems, and this can be technically and practically complicated [5].

- Need for Diverse Datasets:
  The training data set's type and variety play an essential role in evaluating the bias and precision of generative AI systems. However, finding diverse and representative datasets for a study is challenging as some groups are often underrepresented in research or data is scarce in some fields. If there is not enough variety in the training data, then the AI system replicates the existing prejudices that are in force in society and acts unjustly. The problem is deciding on the appropriate datasets, organizing their collection, and verifying the sourced data set's increased inclusiveness and fairness [7].

- Ensuring Transparency in AI Decision-Making
  Trust is a significant factor in the development of AI systems, yet opacity is the primary obstacle to implementing it. Most, if not all, AI systems currently run as 'black boxes' – decision-making within such systems cannot be easily traced or explained even to the software engineers who

designed them. It makes it nearly impossible to address and eliminate biases or ensure that AI systems are not unethical. The problem is that methods and tools that are suitable for obtaining insight into an AI system's decision-making mechanisms and the audit trail of these mechanisms remain a challenge [8].

## 2. Solutions:

- Improving Data Diversity: The bias in generative AI systems can thus be addressed through more diverse and diverse training data sets. This can be reached by searching for evidence from numerous data sources and across different populations. Furthermore, data augmentation and synthetic data generation can help expand a training set. Preventing the leakage of bias from the training data into the final AI models is essential, and this begins with using diverse data as training data [11].

- Adopting Fairness-Aware Algorithms: Bias-sensitive algorithms refer to learning and artificial intelligence techniques to address biases. While gathering data for training, these algorithms can adapt the training process so that the model assigns equal preference to all demographics. They can also include fairness constraints that restrict the model from making prejudiced decisions. As fair techniques are inclined in the development cycle, an AI developer can minimize biased results and have better ethical AI systems [10].

- Implementing Robust Ethical Guidelines: Hence, particular attention should be paid to setting and following strict ethical standards to improve the chances of achieving beneficial outcomes when deploying generative AI systems. These guidelines should cover vital ethical concerns, including privacy, bias, transparency, and accountability. The strategies include A proactive approach to ethics: AI deployers and developers should ensure that ethical aspects are included across the stages of the AI life cycle. Ethical audits and assessments can be scheduled to guarantee that the AI systems are still compliant with the standards set by society and the law [9].

## Conclusion

generative AI is a highly effective technology that may pervasively change numerous aspects of society. However, it is critical to recognize that these technologies entail considerable ethical concerns and biases, which ought to be solved to enable the development of AI in a rewarding and fair manner. Ethically questionable practices include privacy invasions, lack of credible information, questions on the ownership of ideas, and bias. Partly with further social work education and identifying and applying technical, regulatory, and ethical solutions, these challenges are hardly simple and often represent a bundle of issues. As for these challenges, the paper has suggested the following solutions: Increasing the variety of the training data, using fairness-promoting algorithms, and enforcing rigid ethical standards. These strategies are critical for mitigating biases and improving the accountability of AI models to promote their use in a manner consistent with societal norms and ethical practices. Therefore, developers, policymakers, and stakeholders must ensure that generative AI can be utilized in the future to ensure that the technologies enhance the welfare of society and are developed with the highest levels of ethical standards.

## References

1. Bogroff, A., & Guegan, D. (2019). Artificial Intelligence, Data, Ethics: An Holistic Approach for Risks and Regulation. https://shs.hal.science/halshs-02181597/document

2. Mallreddy, S. R., & Vasa, Y. (2023b). Predictive maintenance in cloud computing and DEVOPS: ML models for anticipating and preventing system failures. PREDICTIVE MAINTENANCE IN CLOUD COMPUTING AND DEVOPS: ML MODELS FOR ANTICIPATING AND PREVENTING SYSTEM FAILURES, 10(1), 213–219. https://doi.org/10.53555/nveo.v10i1.5751

3. Mallreddy, S. R., & Vasa, Y. (2023). Natural language querying in SIEM systems: Bridging the gap between security analysts and complex data. NATURAL LANGUAGE QUERYING IN SIEM SYSTEMS: BRIDGING THE GAP BETWEEN SECURITY ANALYSTS AND COMPLEX DATA, 10(1), 205–212. https://doi.org/10.53555/nveo.v10i1.5750

4. Vasa, Y., Mallreddy, S. R., & Jami, V. S. (2022). AUTOMATED MACHINE LEARNING FRAMEWORK USING LARGE LANGUAGE MODELS FOR FINANCIAL SECURITY IN CLOUD OBSERVABILITY. *International Journal of Research and Analytical Reviews , 9(3), 183–190.*

5. Vasa, Y., Singirikonda, P., & Mallreddy, S. R. (2023). AI Advancements in Finance: How Machine Learning is Revolutionizing Cyber Defense. International Journal of Innovative Research in Science, Engineering and Technology, 12(6), 9051–9060.

6. Vasa, Y., & Singirikonda, P. (2022). Proactive Cyber Threat Hunting With AI: Predictive And Preventive Strategies. International Journal of Computer Science and Mechatronics, 8(3), 30–36.

7. Vasa, Y., Mallreddy, S. R., & Jaini, S. (2023). *AI And Deep Learning Synergy: Enhancing Real-Time Observability And Fraud Detection In Cloud Environments, 6(4), 36–42. https://doi.org/ 10.13140/RG.2.2.12176.83206*

8. Katikireddi, P. M., Singirikonda, P., & Vasa, Y. (2021). Revolutionizing DEVOPS with Quantum Computing: Accelerating CI/CD pipelines through Advanced Computational Techniques. Innovative Research Thoughts, 7(2), 97–103. https://doi.org/10.36676/irt.v7.i2.1482

9. Vasa, Y., Cheemakurthi, S. K. M., & Kilaru, N. B. (2022). Deep Learning Models For Fraud Detection In Modernized Banking Systems Cloud Computing Paradigm. International Journal of Advances in Engineering and Management, 4(6), 2774–2783. https://doi.org/10.35629/5252-040627742783

10. Vasa, Y., Kilaru, N. B., & Gunnam, V. (2023). Automated Threat Hunting In Finance Next Gen Strategies For Unrivaled Cyber Defense. International Journal of Advances in Engineering and Management, 5(11). https://doi.org/10.35629/5252-0511461470

11. Vasa, Y., & Mallreddy, S. R. (2022). Biotechnological Approaches To Software Health: Applying Bioinformatics And Machine Learning To Predict And Mitigate System Failures. Natural Volatiles & Essential Oils, 9(1), 13645–13652. https://doi.org/https://doi.org/10.53555/nveo.v9i2.5764

12. Mallreddy, S. R., & Vasa, Y. (2022). Autonomous Systems In Software Engineering: Reducing Human Error In Continuous Deployment Through Robotics And AI. NVEO - Natural Volatiles & Essential Oils, 9(1), 13653–13660. https://doi.org/https://doi.org/10.53555/nveo.v11i01.5765

13. Vasa, Y., Jaini, S., & Singirikonda, P. (2021). Design Scalable Data Pipelines For Ai Applications. NVEO - Natural Volatiles & Essential Oils, 8(1), 215–221. https://doi.org/https://doi.org/10.53555/nveo.v8i1.5772

14. Singirikonda, P., Jaini, S., & Vasa, Y. (2021). Develop Solutions To Detect And Mitigate Data Quality Issues In ML Models. NVEO - Natural Volatiles & Essential Oils, 8(4), 16968–16973. https://doi.org/https://doi.org/10.53555/nveo.v8i4.5771

15. Vasa, Y. (2021). Develop Explainable AI (XAI) Solutions For Data Engineers. NVEO - Natural Volatiles & Essential Oils, 8(3), 425–432. https://doi.org/https://doi.org/10.53555/nveo.v8i3.5769

16. Sukender Reddy Mallreddy. (2023). ENHANCING CLOUD DATA PRIVACY THROUGH FEDERATED LEARNING: A DECENTRALIZED APPROACH TO AI MODEL TRAINING. IJRDO -Journal of Computer Science Engineering, 9(8), 15-22.

17. Mallreddy, S.R., Nunnaguppala, L.S.C., & Padamati, J.R. (2022). Ensuring Data Privacy with CRM AI: Investigating Customer Data Handling and Privacy Regulations. ResMilitaris. Vol.12(6). 3789-3799

18. Nunnagupala, L. S. C. ., Mallreddy, S. R., & Padamati, J. R. . (2022). Achieving PCI Compliance with CRM Systems. Turkish Journal of Computer and Mathematics Education (TURCOMAT), 13(1), 529–535.

19. Jangampeta, S., Mallreddy, S.R., & Padamati, J.R. (2021). Anomaly Detection for Data Security in SIEM: Identifying Malicious Activity in Security Logs and User Sessions. 10(12), 295-298

20. Jangampeta, S., Mallreddy, S. R., & Padamati, J. R. (2021). Data Security: Safeguarding the Digital Lifeline in an Era of Growing Threats. International Journal for Innovative Engineering and Management Research, 10(4), 630-632.

21. Sukender Reddy Mallreddy(2020).Cloud Data Security: Identifying Challenges and Implementing Solutions.JournalforEducators,TeachersandTrainers,Vol.11(1).96 -102.

22. Nunnaguppala, L. S. C. , Sayyaparaju, K. K., & Padamati, J. R.. (2021). "Securing The Cloud: Automating Threat Detection with SIEM, Artificial Intelligence & Machine Learning", International Journal For Advanced Research In Science & Technology, Vol 11 No 3, 385-392

23. Padamati, J., Nunnaguppala, L., & Sayyaparaju, K. . (2021). "Evolving Beyond Patching: A Framework for Continuous Vulnerability Management", Journal for Educators, Teachers and Trainers, 12(2), 185-193.

24. Nunnaguppala, L. S. C. . (2021). "Leveraging AI In Cloud SIEM And SOAR: Real-World Applications For Enhancing SOC And IRT Effectiveness", International Journal for Innovative Engineering and Management Research,10(08), 376-393

25. Sayyaparaju, K. K., Nunnaguppala, L. S. C. , & Padamati, J. R.. (2021). "Building SecureAI/ML Pipelines: Cloud Data Engineeringfor Compliance and Vulnerability Management", International Journal for Innovative Engineering and Management Research,10(10), 330-340