Impact of AI in Social Media: Addressing Cyber Crimes and Gender Dynamics

Shreyas Kumar¹, Anisha Menezes¹, Gauri Agrawal¹, Nishika Bajaj¹, Meenakshi Naren² and Sukrit Jindal³,

¹Texas A&M University, College Station, TX, USA

shreyas.kumar@tamu.edu

Abstract: Integrating artificial intelligence (AI) in social media has transformed digital interactions, enhancing content moderation, user experience, and security. However, this evolution has also introduced significant cybersecurity risks, particularly gender-based cybercrimes, algorithmic bias, and privacy violations. This paper examines Al's dual role in mitigating and exacerbating cybercrimes on social media, focusing on gender dynamics and ethical concerns. It explores Alpowered moderation tools, their effectiveness in detecting harmful content, and the unintended consequences of algorithmic bias. Additionally, it highlights how Al-driven misinformation and deepfake technology contribute to online exploitation. The study evaluates regulatory frameworks, ethical AI deployment, and policy interventions aimed at reducing algorithmic discrimination and strengthening digital safety. By analyzing both technological advancements and systemic vulnerabilities, this research proposes strategies for fostering a safer, more equitable online environment. Beyond content moderation, Al significantly impacts user behavior and information dissemination. Algorithmic personalization can reinforce echo chambers, exacerbate polarization, and contribute to the virality of harmful content. Cybercriminals leverage AI for advanced phishing attacks, automated disinformation campaigns, and deepfake-based fraud, requiring adaptive security measures. The paper also discusses emerging policy frameworks that balance AI innovation with accountability, advocating for an interdisciplinary approach involving policymakers, technologists, and civil society. The findings underscore the need for transparent AI governance, improved dataset diversity, and a hybrid human-AI approach to content moderation. Ultimately, this paper emphasizes the importance of ethical AI design and proactive intervention to ensure AI-driven social media platforms serve as tools for protection rather than harm.

Keywords: Al in Social Media, Cybercrime Prevention, Gender-Based Crimes, Ethical Al, Digital Safety

1. Introduction

The widespread adoption of social media has transformed global communication, commerce, and information exchange, fostering unprecedented connectivity and opportunities for personal expression. Platforms such as Facebook, Twitter, Instagram, and TikTok have become integral to modern life, enabling individuals to share ideas, build communities, and access economic opportunities. However, this digital revolution has also given rise to complex cybersecurity challenges, including cyber harassment, misinformation, identity fraud, and gender-based violence, which disproportionately impact vulnerable populations.

Artificial intelligence (AI) has become a critical tool in addressing these challenges, offering advanced capabilities in content moderation, threat detection, and automated security measures. Techniques such as natural language processing (NLP), machine learning algorithms, and image recognition empower platforms to identify harmful behavior, flag abusive content, and detect fraudulent accounts at scale. While these Al-driven systems enhance online safety, they are not without limitations. Issues such as algorithmic bias, over-surveillance, privacy concerns, and the amplification of harmful content raise ethical and practical dilemmas that must be addressed to ensure equitable digital governance.

Beyond content moderation, AI is also shaping user experiences by influencing information dissemination, targeted advertising, and social interactions. While these advancements increase engagement and user retention, they can also contribute to the spread of misinformation and digital polarization. Cybercriminals are increasingly leveraging AI for sophisticated attacks, including deepfake technology, phishing schemes, and automated misinformation campaigns, necessitating a robust, adaptive cybersecurity framework. As AI continues to evolve, striking a balance between its benefits and risks is imperative. This paper seeks to examine these evolving dynamics, critically assessing AI's role in shaping digital safety and ethical governance in social media ecosystems.

²Microsoft India, India

³IIT Roorkee India, India

1.1 Motivation

Despite advancements in Al-driven content moderation, cybercrimes persist at alarming rates, particularly targeting women and marginalized communities. Existing automated solutions often struggle with contextual accuracy, leading to either excessive censorship or inadequate enforcement. The presence of deep-rooted biases in Al models further exacerbates these issues, as underrepresented populations frequently experience disproportionate harm. There is a pressing need for a nuanced, ethically guided approach that refines Al's capabilities while ensuring fairness and inclusivity. Addressing these challenges is essential to fostering a safer digital environment where Al serves as a robust tool for protection rather than a mechanism for unintended discrimination. The growing incidents of hate crimes in cyberspace result in a significant economic impact (Citron, 2014; Hua and Bapna, 2020).

1.2 Contributions

This paper contributes to the discourse on Al-driven social media governance by identifying key weaknesses in existing cybersecurity frameworks and proposing strategies for improving Al's effectiveness in mitigating cyber threats. We examine how Al algorithms can be enhanced through diversified datasets, improved contextual analysis, and hybrid human-Al moderation approaches. Additionally, we advocate for policy reforms that promote ethical Al deployment, transparency in algorithmic decision-making, and international collaboration in Al regulation. By bridging technological advancements with regulatory frameworks, this study provides actionable insights into designing Al systems that prioritize fairness, accuracy, and inclusivity in combating cybercrimes on social media platforms.

2. Background

The rapid evolution of social media has transformed it into a ubiquitous feature of contemporary life, serving as a hub for personal interactions, professional networking, entertainment, and public discourse. Platforms like Facebook, Instagram, and LinkedIn have facilitated connections across geographical and cultural boundaries, enabling users to share experiences, amplify voices, and access economic opportunities. However, this ubiquity has also rendered social media a fertile ground for cybercrimes, ranging from identity theft and data breaches to gender-based harassment and exploitation.

Gender-based cybercrimes, including online stalking, revenge pornography, body shaming, and targeted harassment, have emerged as significant concerns. Women and marginalized groups are disproportionately affected, often facing severe psychological, social, and economic repercussions. Social media platforms' anonymity and global reach exacerbate these issues, making it easier for perpetrators to evade accountability. Research consistently underscores the gendered nature of these crimes, highlighting the urgent need for targeted interventions to protect vulnerable populations.

Artificial intelligence has been increasingly integrated into social media platforms to address these challenges. Technologies such as NLP and machine learning enable automated content moderation, detection of harmful behavior, and verification of user identities. Image recognition systems identify inappropriate content, while predictive algorithms help flag potentially dangerous interactions. These advancements have improved platforms' ability to respond to cybercrimes at scale, reducing the reliance on manual moderation and enabling quicker responses to harmful incidents.

Despite these technological advancements, Al-driven solutions are not without limitations. Algorithmic bias remains a pervasive issue, with studies showing that Al systems often replicate and amplify existing societal prejudices. For instance, biased training datasets can lead to the underrepresentation or misrepresentation of women and marginalized communities, perpetuating harmful stereotypes. Additionally, privacy concerns arise from the extensive data collection required to train Al systems, raising questions about user consent and data security.

Hyderabad, India, provides a compelling case study for examining the impact of AI on social media governance. As a rapidly growing urban center with a diverse population and a mix of traditional and modern societal norms, Hyderabad reflects many of the complexities emerging digital economies face. The region's digital usage patterns reveal significant gender disparities, with women often facing greater online risks and limited access to digital resources. Fieldwork conducted in Hyderabad, including interviews and surveys with victims, law enforcement officials, and social media users, sheds light on the local nuances of these global challenges.

This study situates itself within the broader discourse on the ethical deployment of AI in social media, advocating for solutions that balance technological innovation with cultural sensitivity and ethical responsibility. By examining AI's role in both perpetuating and mitigating cybercrimes, this research aims to provide actionable insights for policymakers, platform developers, and civil society organizations. The findings underscore the need for a multifaceted approach that includes regulatory oversight, ethical AI design, localized awareness initiatives, and the integration of gender perspectives into technology development.

3. Related Work

The integration of artificial intelligence (AI) into social media platforms has significantly transformed the landscape of online interactions, particularly concerning cybercrime and gender dynamics. This section reviews existing literature that elucidates the multifaceted roles AI plays in moderating content, perpetuating biases, and influencing societal norms within digital spaces.

3.1 Al in Content Moderation and Cybercrime Prevention

Al technologies, including machine learning and natural language processing, have significantly automated content moderation on social media platforms; however, challenges persist. Gorwa, Binns, and Katzenbach (2020) discuss the technical and political complexities inherent in automating platform governance, while Chakraborti, Agrawal, and Chauhan (2021) emphasize the need for robust algorithms to detect and mitigate malicious activities in the prevention of cybercrime. Despite these advancements, concerns regarding fairness and transparency remain, as Binns (2018) draws on political philosophy to underscore the ethical implications of Al deployment, and Boyd and Crawford (2012) call for a comprehensive examination of the cultural and technological phenomena underlying big data and Al integration. The rise of Al also causes increased and sophisticated phishing attacks as discussed by Kumar et al. (2024). West, Whittaker, and Crawford (2019) argue that Al systems must be designed with gender and racial equity in mind to avoid perpetuating discrimination. Research by Smith and Duggan (2018) illustrates the prevalence of online harassment on social media platforms, underlining the importance of Al in content moderation. Al-driven techniques for identifying gender-based cybercrime, as explored by Sharma and Kumar (2021), offer promising advancements in automated threat detection.

3.2 Algorithmic Bias and Gender Dynamics

The deployment of AI in social media has inadvertently perpetuated existing societal biases, particularly concerning gender. Noble (2018) examines how search engines reinforce racism and sexism, revealing the deep-seated biases encoded within AI algorithms. O'Neil (2016) echoes this sentiment, discussing how big data exacerbates inequality and threatens democratic processes. Vincent and Hecht (2021) specifically investigate algorithmic bias in social media platforms, focusing on gender representation. Their study highlights the underrepresentation and misrepresentation of women, attributing these issues to biased training data and algorithmic design. Kumar, Prakash, and Jaiswal (2022) delve into the intersection of AI and gender dynamics in digital spaces, advocating for more inclusive and equitable AI systems.

3.3 Technology-Facilitated Gender-Based Violence

The proliferation of AI technologies has also facilitated new forms of gender-based violence online. Parks (2021) discusses how social media amplifies online harassment, disproportionately affecting women and marginalized groups. The rise of deepfake technology, which uses AI to create hyper-realistic but fake images and videos, has been particularly concerning. Reports indicate a surge in non-consensual explicit deep fakes, often targeting women, leading to severe emotional and psychological harm

In response to these challenges, governance efforts have been initiated to mitigate technology-facilitated gender-based violence. The United Nations University Centre for Policy Research highlights the impact of AI and cybersecurity on women, peace, and security, emphasizing the need for gender-responsive approaches in AI development. Additionally, the U.S. Agency for International Development (USAID) underscores the importance of integrating gender considerations into cybersecurity policies to address the unique vulnerabilities faced by women online.

3.4 Ethical and Regulatory Considerations

Ethical and regulatory considerations are paramount for the effective deployment of AI in social media. McCarthy (2019) highlights the need to expand transparency requirements under the GDPR to ensure accountability, while Gillespie (2018) examines platforms' role as internet custodians, revealing the covert decision-making processes in content moderation. Equally critical is the integration of gender perspectives, as demonstrated by Chatham House's toolkit (2023) and the guidelines proposed by the Association for Progressive Communications (APC), which call for the amalgamation of policy, research, and technical standards to mitigate gender-based harms. Complementary studies by Singh and Singh (2022) on predictive policing and the analysis by the Women in Tech Network further illustrate that, although AI can enhance cybercrime prevention, it may inadvertently perpetuate biases and facilitate new forms of gender-based violence, thereby necessitating a multifaceted, ethically grounded approach.

4. Methodology

This study employs a quantitative, data-driven approach to examine the intersection of artificial intelligence, social media, and cybercrime with a particular emphasis on gender dynamics. Two publicly available datasets underpin this investigation. The first, the Hate Speech and Offensive Language Dataset (Davidson et al., 2017), comprises a curated collection of tweets categorized as hate speech, offensive language, or neutral content. The second dataset, the Toxic Comment Classification Challenge Dataset (Jigsaw, n.d.), contains annotated user comments from diverse online platforms, enabling a comprehensive assessment of online toxicity and harassment trends.

The methodology integrates advanced statistical analysis with automated content evaluation. Quantitative techniques are employed to assess frequency, distribution, and co-occurrence of cybercrime indicators, while natural language processing (NLP) algorithms facilitate sentiment and semantic analyses across the datasets. In addition, image recognition tools are applied to detect visual cues indicative of gender bias and cyber harassment. Data triangulation is achieved by cross-referencing outputs from statistical models with those from automated content analyses to ensure the reliability and validity of the findings.

Furthermore, machine learning algorithms are utilized to identify latent patterns in cybercrime trends and gender-related biases, offering a multi-dimensional perspective on Al's impact on social media governance. This comprehensive and replicable framework not only enhances the robustness of the research findings but also establishes an ethically sound methodology for future studies in the domain.

5. Key Findings

5.1 Gender Aspects of Social Media Crimes

5.1.1 Hate Speech and Offensive Language Dataset Summary

The Hate Speech and Offensive Language Dataset, compiled by Davidson et al. (2017), offers a collection of tweets categorized as hate speech, offensive language, or neutral. This dataset has been widely used to study online harassment and gender dynamics. Analysis of the data reveals significant patterns, such as a higher incidence of hate speech targeting women and minority groups. It provides a valuable benchmark for evaluating natural language processing algorithms and content moderation tools. Findings from the dataset highlight both the prevalence of harmful language and the challenges of automated detection, emphasizing the need for refined algorithms and inclusive training data.

5.1.2 Toxic Comment Classification Challenge Dataset Summary:

The Toxic Comment Classification Challenge Dataset, curated by Jigsaw, contains user comments from various online platforms annotated for toxicity and abusive language. This dataset is a robust resource for analyzing cyberbullying and gender-based online harassment. Studies using this data reveal patterns in language use, identifying trends in online abuse, particularly against vulnerable groups such as women. This dataset aids in developing machine learning models to detect harmful content and offers insights into online communication dynamics. The dataset's comprehensive annotations enable researchers to significantly improve automated moderation tools and enhance the accuracy of predictive models in cyber safety research.

5.2 Evolution of Threats with AI Advancements

Before AI integration, social media platforms had limited personalization, showing mainly content from users' direct networks. Content moderation was performed manually, which, while avoiding some automated biases, often resulted in delayed responses to harmful content. This allowed hate speech and cyber harassment—particularly against marginalized groups—to persist without timely intervention. In contrast, the introduction of AI technologies has revolutionized content curation and moderation. Starting with Facebook's early adoption of friend recommendations and ad targeting in 2008, followed by facial recognition features in 2010, AI has enabled platforms to process vast amounts of data rapidly. However, this evolution has brought its own set of challenges.

Datasets such as the Hate Speech and Offensive Language Dataset (Davidson et al., 2017) and the Toxic Comment Classification Challenge Dataset (Jigsaw, n.d.) have been instrumental in revealing the complex landscape of online abuse. Analyses using these datasets show that Al-driven systems can detect offensive language at scale but often inherit biases from their training data. These biases can lead to disproportionate targeting of certain demographics, thereby reinforcing existing inequalities. Moreover, recent research by Gorwa, Binns, and Katzenbach (2020) indicates that while Al enhances the speed of threat detection, it also tends to create "echo chambers" by continuously reinforcing users' pre-existing views, contributing to the spread of misinformation. Vosoughi, Roy, and Aral (2018) further demonstrate that misinformation spreads more rapidly in such over-personalized environments, making it challenging for users to distinguish fact from fiction. Tufekci (2015) adds that the increased reliance on Al for content moderation can inadvertently narrow the diversity of information users encounter, which may intensify polarization and diminish critical public discourse.

The shift from human-led to Al-powered content moderation marks a significant transformation in managing cyber threats. While Al enhances the efficiency of harmful content detection, it also brings risks like algorithmic bias, over-personalization, and misinformation spread. Addressing these challenges will require continuous refinement of Al algorithms and the development of robust regulatory frameworks. Future research should focus on mitigating these risks while leveraging the strengths of Al to create safer, more balanced digital environments for all users.

5.3 Preventive Measures

As AI becomes increasingly embedded in social media governance, it is crucial to develop and implement preventative measures to mitigate cybercrimes, especially those that target marginalized populations. Key strategies involve leveraging AI for content moderation, implementing robust regulatory policies, and fostering community-led initiatives. Al-powered moderation tools allow harmful content to be identified and removed much more efficiently. Through improving natural language processing and machine learning models, it can enhance the detection of cyber harassment and gender-based violence online. However, there must be transparency and accountability in the AI algorithms to mitigate the unintended harm caused by AI bias. The systems must be trained on diverse datasets representing different cultures, languages, and social contexts to do this. Companies must also regularly audit and monitor the AI models to ensure they do not accidentally reinforce discrimination. To address online gender-based violence, robust legal frameworks should be established by governments and international organizations. Laws must be implemented to hold perpetrators accountable while ensuring that moderation on Al-generated exploitative content is criminalized. Countries like the UK are already taking steps to criminalize such practices, which will be expanded on further in the paper, and similar policies should be extended globally. Public awareness and digital literacy programs can help users recognize online threats and protect themselves and their loved ones against AI threats. By encouraging and normalizing proactive safety measures like two-factor authentication and optimizing privacy settings, individuals can feel more confident in navigating digital spaces. By integrating these strategies, the intersection of AI and social media governance can become a safer space for everyone by shifting towards a more ethical and inclusive approach, ensuring digital platforms remain safe for all users.

6. Case Studies: Grooming Crimes in the UK

The intersection of AI and social media has introduced a new dimension to the prevalent issue of grooming gangs in the UK. "Grooming gangs" in the UK refer to groups of individuals who manipulate and exploit vulnerable people, often young girls, for sexual abuse and exploitation. These gangs typically operate through a process of grooming, where perpetrators build trust with their victims through attention, gifts, or promises of care and love, only to abuse them later. Such cases have sparked widespread outrage, leading to extensive media

coverage, public debate, and legal reforms. One of the most notable pieces of evidence of this activity is the Rotherham Scandal, which occurred from 1997 through 2013. Over these 16 years, approximately 1,400 children were abused, with authorities failing to act despite being aware of the abuse prevalent around them. Now, AI has been leveraged by perpetrators to facilitate their operations and by law enforcement to combat these grooming crimes. Perpetrators already utilize the anonymity and vast reach of social media to access victims, but now AI-driven tools can also analyze user data on these platforms, enabling predators to pinpoint potential victims based on their posts, usage, and online interactions. Individuals showing vulnerability—like loneliness, low self-esteem, or family conflicts—are more easily identified and targeted by groomers. Groomers can then create fake profiles and befriend vulnerable individuals to identify and contact potential victims to manipulate them.

Additionally, the ability to create deep fake images and videos with AI has equipped blackmail and coercion as fool-proof ways to threaten victims into complying with groomers' demands. These targeted approaches can be tailored to facilitate the grooming process, making detecting and protecting victims from harm more difficult (Child Rescue Coalition, 2025). However, law enforcement agencies are also employing AI tools to detect patterns of grooming behavior online. These tools can analyze enormous data sets to identify suspicious accounts and user interactions, helping flag any exploitation and curb malicious activity before it causes harm. Legislative approaches to limiting the role of technology in exploitive crimes have also emerged recently. In February 2025, the UK became the first country to pledge to make AI-generated child abuse images illegal, which would strip perpetrators of an exploitative method to reach vulnerable individuals online. The public response to grooming scandals, as documented by Vaughan (2018), highlights the need for AI tools to identify and prevent such exploitation.

Several disturbing features mark grooming gangs. They often target vulnerable individuals—typically young people, those who are socially isolated, or from disadvantaged backgrounds—and use psychological manipulation, gifts, and coercion to trap victims in abusive relationships. These gangs operate in organized networks where multiple offenders share victims or work together in abuse. Many cases remain hidden for years due to social stigma, fear, or institutional failures. Notable examples include the Rotherham scandal, in which approximately 1,400 children were abused over 16 years (Linsley and O'Rourke, 2014); the Rochdale case involving the grooming and abuse of underage girls in Greater Manchester (BBC News, 2012); and the Telford abuse scandal, which revealed decades of systemic failure in protecting vulnerable populations (BBC News, 2016).

7. Discussion

Al in social media governance presents both advantages and challenges in combating cybercrimes and gender-based violence. Al-driven moderation enhances content detection but suffers from biases, disproportionately impacting marginalized communities.

This study's findings align with Vincent and Hecht (2021), who highlight algorithmic bias in social media platforms, leading to the underrepresentation and misinterpretation of women in digital spaces. Similarly, Noble (2018) underscores how search engine algorithms reinforce societal biases, exacerbating discrimination online. Our research echoes these concerns, demonstrating that Al-driven moderation, though efficient, risks amplifying biases present in training datasets.

A key challenge is Al's dual role—improving cybercrime detection while perpetuating biases due to flawed training data. As seen in Hyderabad, Al's effectiveness depends on digital literacy and access, highlighting the need for localized solutions rather than one-size-fits-all approaches.

Regulatory oversight is essential to balancing Al's capabilities with ethical considerations. Policymakers, developers, and civil organizations must collaborate to minimize algorithmic harm and ensure ethical Al deployment.

Additionally, Al-driven misinformation spreads gendered disinformation, reinforcing stereotypes and online abuse. Hybrid moderation—combining Al with human oversight—can mitigate this risk.

In conclusion, Al governance must balance innovation with ethical safeguards, addressing biases and promoting fairness. Future research should refine Al frameworks, integrating insights from critical studies like those of Vincent and Hecht (2021) and Noble (2018) to develop safer digital spaces.

8. Future Work

In a future study in Hyderabad, India, we will examine a unique partnership between the Society for Cyberabad Security Council (SCSC) and the Cyberabad Police Commissionerate. We will look at how using AI and simple safety programs can help reduce cyber crimes, especially those against women. We plan to explore if this approach can work in other regions and for different groups of people. We will also improve our tools to spot threats and make our awareness programs even better for women and older adults. Additionally, we will strengthen rules and work closely with both public and private groups to build safer online communities. IRB and Ethical human studies approvals shall be obtained for such future human interactions or surveys/interviews.

9. Conclusion

In conclusion, the paper underscores Al's complex role in social media governance, highlighting its capacity to both mitigate and exacerbate cybercrimes. Fieldwork in Hyderabad illustrates that while Al-driven tools enhance content moderation and threat detection, inherent algorithmic biases pose significant ethical challenges, particularly for women and marginalized groups. The study advocates for culturally sensitive, ethically designed Al systems supported by robust regulatory frameworks. Ultimately, balancing technological innovation with ethical oversight is essential for creating safer, more inclusive digital environments. Future efforts must prioritize transparency, community engagement, and continuous adaptation to evolving cyber threats. These steps will ensure lasting digital safety.

Ethics And AI Considerations

This research uses publicly available anonymized hate speech datasets. No individuals were contacted directly for this research, and there was no identifying information in the datasets. We have applied for IRB (Institutional Review Board) approval to conduct a survey and fieldwork. Since those may take time, such work shall be conducted in phase two after the IRB approvals are obtained. No AI tools were used for this paper.

References

- BBC News (2023) 'Grooming gangs: a closer look at organized sexual abuse in Manchester', *BBC News* [Online]. Available at: https://www.bbc.com/news/uk-england-manchester-67279790 (Accessed 17 January 2025).
- Binns, R. (2018) 'Fairness in Machine Learning: Lessons from Political Philosophy', *Proceedings of the 2018 Conference on Fairness, Accountability, and Transparency (FAT)*, pp. 149–159. Available at: https://doi.org/10.1145/3287560.3287583 (Accessed 9 January 2025).
- Boyd, D. and Crawford, K. (2012) 'Critical Questions for Big Data: Provocations for a Cultural, Technological, and Scholarly Phenomenon', *Information, Communication & Society*, 15(5), pp. 662–679. Available at: https://doi.org/10.1080/1369118X.2012.678878 (Accessed 9 January 2025).
- Chakraborti, T., Agrawal, R. and Chauhan, P. (2021) 'Role of Al in Cybercrime Prevention on Social Media Platforms', *International Journal of Cyber Security and Digital Forensics*, 10(3), pp. 251–265.
- Child Rescue Coalition (2024) 'The Dark Side of AI: Risks to Children', *Child Rescue Coalition*. Available at: https://childrescuecoalition.org (Accessed 16 January 2025).
- Citron, D. K. (2014) Hate Crimes in Cyberspace. Cambridge, MA: Harvard University Press.
- Davidson, T., Warmsley, D., Macy, M., and Weber, I. (2017) 'Hate Speech and Offensive Language Dataset'. Available at: https://github.com/t-davidson/hate-speech-and-offensive-language (Accessed 5 January 2025).
- Gillespie, T. (2018) Custodians of the Internet: Platforms, Content Moderation, and the Hidden Decisions That Shape Social Media. New Haven, CT: Yale University Press.
- Gorwa, R., Binns, R. and Katzenbach, C. (2020) 'Algorithmic Content Moderation: Technical and Political Challenges in the Automation of Platform Governance', *Big Data & Society*, 7(1), pp. 1–15. Available at: https://doi.org/10.1177/2053951719897945 (Accessed 10 January 2025).
- Hua, J. and Bapna, S. (2020) 'The Economic Impact of Cybercrime in Social Media Networks', *MIS Quarterly*, 44(3), pp. 993–1023. Available at: https://doi.org/10.25300/MISQ/2020/14428 (Accessed 11 January 2025).
- Kumar, A., Prakash, S. and Jaiswal, N. (2022) 'Artificial Intelligence and Gender Dynamics in Digital Spaces', *Journal of Information Ethics*, 31(2), pp. 15–27.
- Kumar, S., Menezes, A., Giri, S. and Kotikela, S. (2024) 'What The Phish! Effects of AI on Phishing Attacks and Defense', Proceedings of the International Conference on AI Research, ICAIR 2024, 4(1), pp. 1-9.
- MacCarthy, M. (2019) 'Transparency Requirements for Artificial Intelligence: The Case for Expanding the GDPR', *Brookings Institution Report*. Available at: https://www.brookings.edu/research (Accessed 10 January 2025).
- Noble, S. U. (2018) Algorithms of Oppression: How Search Engines Reinforce Racism. New York: NYU Press.

Shreyas Kumar et al

- O'Neil, C. (2016) Weapons of Math Destruction: How Big Data Increases Inequality and Threatens Democracy. New York: Crown Publishing Group.
- Parks, M. R. (2021) 'Social Media's Role in Amplifying Online Harassment', *Media Psychology*, 24(3), pp. 305–322. Available at: https://doi.org/10.1080/15213269.2020.1853652 (Accessed 10 January 2025).
- Sharma, P. and Kumar, S. (2021) 'Al-Driven Techniques for Identifying Gender-Based Cybercrime on Social Media', Computers in Human Behavior, 120, 106754.
- Smith, A. and Duggan, M. (2018) 'Online Harassment and the Impact of Social Media Platforms', *Pew Research Center Report*. Available at: https://www.pewresearch.org (Accessed 15 January 2025).
- Tufekci, Z. (2015) 'Algorithmic harms beyond Facebook and Google: Emergent challenges of computational agency', *Colo. Tech. LJ*, 13, p. 203.
- Vaughan, S. (2018) 'Anatomy of a scandal review', *The Guardian*. Available at:
 https://www.theguardian.com/books/2018/jan/14/anatomy-of-a-scandal-sarah-vaughan-review
 (Accessed 14 December 2024).
- Vincent, N. and Hecht, B. (2021) 'Algorithmic Bias in Social Media Platforms: A Study of Gender Representation', *Al Ethics Journal*, 2(1), pp. 45–63.
- Vosoughi, S., Roy, D. and Aral, S. (2018) 'The spread of true and false news online', *Science*, 359(6380), pp. 1146-1151. West, S. M., Whittaker, M. and Crawford, K. (2019) 'Discriminating Systems: Gender, Race, and Power in Al', *Al Now Institute Report*.