

Review

Not peer-reviewed version

AI Ethics in Social Media

[Janaka Ishan Senarathna](#) *

Posted Date: 14 July 2025

doi: 10.20944/preprints2025071091.v1

Keywords: AI ethics; social media; deepfakes; content authenticity; user trust; harassment



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Review

AI Ethics in Social Media

Janaka Ishan Senarathna

Independent researcher, janakaishansenarathna0169@gmail.com

Abstract

The rapid integration of Artificial Intelligence (AI) into social media has revolutionized user experiences, enabling personalized content delivery and creative media generation, yet it introduces profound ethical challenges, including AI-driven harassment, bullying, and the spread of synthetic media like deepfakes, threatening user trust and societal cohesion. This paper addresses the critical gap in understanding user perceptions of AI ethics, the technical limitations of detecting harmful AI-generated content, and its severe psychological and societal impacts. Employing a mixed-methods approach, we conducted the "AI Ethics in Social Media Questionnaire," collecting 200 responses from predominantly young, undergraduate users, complemented by qualitative analysis of real-world case studies and an extensive literature review. Our innovative integration of empirical survey data with case studies reveals that 97% of users are aware of AI features (66% highly aware, 31% somewhat aware), yet 53% express significant concern about ethical issues, including privacy violations, algorithmic bias, and inadequate content moderation. Alarming, 61% reported encountering AI-generated harassment, with 36% experiencing direct or indirect impact, underscoring the issue's pervasiveness. Case studies, including Sewell Setzer III, Molly Russell, Chase Nasca, the Belgian Man, and deepfake victimizations, illustrate AI's role in exacerbating psychological distress, reputational harm, and societal distrust through algorithmic amplification and inadequate detection mechanisms. The analysis highlights the persistent "arms race" between AI content generation and detection, compounded by algorithmic biases and scalability challenges in moderation. We propose a multi-stakeholder framework, including enhanced user control over AI interactions, robust platform policies with mandatory content labeling, advanced detection technologies, international regulatory collaboration, and public education on media literacy. This work advances AI ethics by offering a comprehensive strategy for responsible AI governance, fostering a safer digital environment, and safeguarding user well-being and public trust. Failure to implement these measures risks escalating online harms, undermining public discourse, and eroding the trust underpinning digital interactions.

Keywords: AI Ethics; social media; deepfakes; content authenticity; user trust; harassment

I. Introduction

A. Background on AI in Social Media

The contemporary digital landscape is profoundly shaped by the pervasive integration of Artificial Intelligence (AI) into social media platforms. AI algorithms are fundamental to various platform functionalities, ranging from content recommendations and targeted advertising to sophisticated content moderation systems.[24] This widespread deployment is not merely a technical convenience; it fundamentally transforms how users interact with digital content and with each other. A recent survey reveals that a significant majority of social media users are cognizant of AI's presence in these platforms, with 66% of respondents indicating they are "Yes, very aware" and an additional 33% being "Somewhat aware" that social media platforms leverage AI for diverse features.[1] This widespread recognition underscores AI's entrenched role in daily online experiences.

The rapid evolution of AI capabilities, particularly in the domain of generative AI, has further expanded its influence, fundamentally transforming content creation and user interaction paradigms.[2] Generative

AI, capable of producing novel and realistic media, has opened new avenues for creative expression and enhanced efficiency in content production.[2] However, this technological acceleration has also introduced unprecedented complexities concerning content authenticity, control, and the potential for misuse.[11] The swift progression of AI from rudimentary algorithms to highly sophisticated generative models has outpaced the development of corresponding societal norms and regulatory frameworks.[21] This imbalance creates a fertile ground for ethical dilemmas, as the technology advances at a pace that legislative and ethical considerations struggle to match, leading to a reactive rather than proactive approach to governance.

B. Emergence of Ethical Concerns

AI's growing influence extends beyond its intended benign applications, giving rise to complex ethical issues that challenge the integrity of digital spaces. These concerns include pervasive algorithmic bias [10], significant privacy violations [11], and a persistent lack of transparency in AI's operational mechanisms.[7] These issues are central to the global discourse on responsible AI development, as they directly impact user rights, fairness, and trust in digital platforms.[4] The increasing realism and accessibility of AI-generated content (AIGC) further compounds these challenges.[13]

A significant concern is the proliferation of synthetic media, such as deepfakes and fake images, which are becoming increasingly difficult for humans to distinguish from authentic content.[13] Research indicates that individuals are often no longer able to reliably determine whether media is AI-generated or real.[13] This inability to discern authenticity poses a fundamental threat to information integrity and erodes user trust in digital media.[19] When users cannot trust the information they encounter online, it makes them more susceptible to manipulation and misinformation, thereby undermining the foundational principles of informed public discourse and democratic processes.[44] This challenge is not merely technical; it represents a profound cognitive and societal vulnerability that demands urgent attention.

C. Problem Statement: AI-driven Harassment, Bullying, and Synthetic Content

The core focus of this paper is the detrimental impact of AI when maliciously leveraged for harassment and bullying on social media, particularly through the creation and dissemination of deepfakes and other manipulated content.[17] These forms of digital abuse can inflict severe psychological and social repercussions on victims.[17] The psychological impacts are profound, encompassing heightened anxiety, depression, and, in extreme cases, suicidal ideation, directly linked to prolonged exposure to such harmful content.[21]

The insidious nature of AI-driven harassment is not solely rooted in the creation of harmful content but is significantly amplified by the algorithmic mechanisms inherent to social media platforms.[22] These algorithms, designed to maximize user engagement and attention, can rapidly disseminate malicious content, reaching vulnerable individuals with unprecedented speed and scale.[21] This creates a self-reinforcing cycle where harmful content, once generated, can quickly go viral, intensifying psychological distress and emotional harm for victims before effective countermeasures can be deployed.[19] The phenomenon of "doom scrolling," where users are compelled to continuously consume content, can exacerbate these negative mental health effects, further deepening the impact of exposure to disturbing or manipulative AI-generated material.[24] This systemic amplification transforms isolated incidents of content creation into widespread threats, underscoring the urgent need for comprehensive interventions.

D. Research Questions and Objectives

This paper aims to address several critical research questions:

- How aware are social media users of AI features and their ethical implications?
- What are users' experiences and concerns regarding AI-generated content used for harassment and bullying?
- What are the psychological and societal impacts of AI-driven content on social media users?
- What are the current technical capabilities and limitations in detecting AI-generated harmful content?
- What ethical principles and regulatory frameworks are most pertinent to mitigating AI-driven harm on social media?
- What actionable recommendations can be proposed for platforms, regulators, and users?

To answer these questions, the objectives of this study include:

- Analysing survey data to understand user perceptions and experiences.
- Detailing relevant real-world case studies to illustrate the severity of the problem.
- Conducting a comprehensive review of existing literature on AI ethics, content generation, and detection.
- Proposing actionable future directions and recommendations for addressing AI-driven ethical challenges in social media.

II. Literature Review

A. Foundational Concepts in AI Ethics

The ethical landscape of Artificial Intelligence is governed by several core principles that guide its responsible development and deployment. These include Fairness, Transparency, Accountability, Privacy, and the overarching principle of Preventing Harm.[4] Fairness, for instance, mandates that AI systems should not perpetuate or amplify existing societal biases, ensuring equitable treatment across diverse user groups.[6] This requires meticulous examination of training data and algorithmic design to mitigate discrimination.[10] Transparency involves providing clear information about how AI systems operate, the data they utilize, and the reasoning behind their decisions, thereby fostering trust and enabling scrutiny.[4] Accountability ensures that mechanisms are in place to assign responsibility and provide redress when AI systems cause unintended harm or make mistakes.[7] Privacy focuses on protecting user data throughout its lifecycle, from collection to retention, and ensuring informed consent for its use.[6] The principle of Preventing Harm emphasizes the proactive identification and mitigation of potential risks, such as those leading to psychological distress or physical danger.[17]

In the context of social media, these principles take on specific relevance. For instance, in content

moderation, fairness demands that AI algorithms do not disproportionately censor or amplify content based on user demographics or political leanings.[10] Transparency requires platforms to disclose when AI is used to generate content or influence user feeds, allowing users to understand the nature of the information they consume.[28] Accountability mechanisms are crucial for addressing instances where AI-driven moderation leads to wrongful content removal or, conversely, fails to detect harmful material.[7] Privacy protocols must safeguard vast amounts of user data collected by social media AI, preventing misuse or breaches.[11] Finally, preventing harm is paramount, given the direct link between exposure to harmful content and severe psychological impacts.[17]

While these ethical principles are widely accepted as foundational for responsible AI, their practical implementation often involves complex trade-offs. For example, achieving complete transparency in AI models, especially "black box" deep learning networks, can conflict with proprietary interests or even security concerns.[7] Similarly, stringent privacy measures might limit the data available for training AI models that could otherwise improve content moderation or

personalization.[10] Balancing the prevention of harm (e.g., through content removal) with freedom of speech also presents a perpetual challenge for platforms and regulators.[29] These inherent conflicts necessitate nuanced policy decisions and platform designs that prioritize human well-being while navigating technological complexities.

B. Evolution of AI in Social Media and its Societal Implications

The trajectory of AI in social media has evolved significantly, moving from rudimentary algorithms to highly sophisticated generative models. Initially, AI was primarily employed for basic functions such as content recommendations and targeted advertising, learning user preferences to personalize feeds.[24] Over time, advancements in machine learning, particularly deep learning, have enabled the development of AI systems capable of autonomously planning and acting to achieve goals with minimal human oversight, often referred to as general-purpose AI agents.[3] This evolution has profoundly impacted the information ecosystem, transforming how content is created,

disseminated, and consumed.[2]

The societal implications of this evolution are far-reaching. General-purpose AI agents and highly personalized content delivery systems on social media have amplified the potential for subtle, yet pervasive, manipulation of user behaviour and public discourse.[3] AI's ability to generate persuasive content at scale makes it easier for malicious actors to influence public opinion, potentially affecting political outcomes and deepening societal divisions.[3] This algorithmic influence extends beyond mere content presentation; it actively shapes and can potentially distort individual perceptions and collective narratives.[23] Users are increasingly exposed to content tailored to their existing biases, reinforcing "echo chambers" and making it harder for individuals to discern authentic information or resist algorithmic influence.[24] This phenomenon poses a significant challenge to critical thinking and informed decision-making, as the digital environment becomes increasingly curated and potentially manipulative.[64] The rapid adoption of general-purpose AI by individuals and businesses further underscores the urgency of addressing these societal impacts.[3]

C. Deepfake and Fake Image Generation Techniques

The creation of highly realistic synthetic media, commonly known as deepfakes and fake images, is primarily driven by advanced AI models. The most prominent architectures include Generative Adversarial Networks (GANs), Diffusion Models (DMs), and Variational Autoencoders (VAEs).[2]

- **Generative Adversarial Networks (GANs):** GANs operate on an adversarial principle, consisting of a generator network that creates synthetic content and a discriminator network that attempts to distinguish between real and generated content.[13] Through this "zero-sum game," the generator continuously improves its ability to produce hyper-realistic fakes until the discriminator can no longer differentiate them from authentic data.[68] GANs have been instrumental in face-swapping and face reenactment, where the generator learns to map source identity attributes onto a target face or synchronize facial expressions with audio inputs.[59]
- **Diffusion Models (DMs):** Diffusion models represent a newer class of generative AI that has shown remarkable capabilities in image synthesis.[13] These models work by gradually adding noise to an image until it becomes pure noise, and then learning to reverse this process to generate a clean image from noise.[2] This iterative denoising process allows for the creation of high-quality and diverse images, including those used for malicious purposes.[66]
- **Variational Autoencoders (VAEs):** VAEs are a type of autoencoder neural network that provide a probabilistic approach to generating realistic fake images.[2] They learn a latent space as statistical parameters of probabilistic distributions, which significantly improves the

quality of generated results compared to earlier autoencoders.[2] VAE-based architectures are also employed in face-swapping, where they can obtain a latent representation of a face independent of geometry and non-face regions, which is then used to synthesize a swapped image.[59]

These technologies enable the creation of highly convincing but fabricated content, such as realistic face swaps, accurate voice cloning, and sophisticated text generation for malicious purposes, often making them indistinguishable from real content to the human eye.[43] The continuous advancement in generative AI models implies that the quality and realism of deepfakes are rapidly improving, creating a perpetual "arms race" where detection methods consistently lag behind generation capabilities.[13] This dynamic imbalance means that the problem is not static; rather, it is an ongoing challenge where the technology of harm is perpetually ahead, demanding continuous innovation in countermeasures and a multi-faceted approach to mitigation.

D. Challenges in AI-Generated Content Detection

Despite significant advancements in deepfake detection technologies, substantial technical challenges persist, particularly when attempting to identify "in-the-wild" AI-generated content that is representative of real-world threats.[13] Research

indicates that the performance of state-of-the-art deepfake detection models drops significantly when evaluated on contemporary, real-world datasets compared to academic benchmarks.[13] This performance gap is a critical concern, as it implies that current detection tools may not be adequately prepared for the evolving sophistication of malicious AI-generated content.

Key technical difficulties include:

- **Rapid Evolution of Generation Techniques:** Deepfake generation technologies are constantly evolving, with new models and methods emerging that can produce increasingly realistic and harder-to-detect synthetic media.[13] This creates a continuous cat-and-mouse game where detection methods struggle to keep pace.[13]
- **Susceptibility to Adversarial Attacks:** Deepfake detectors can be vulnerable to adversarial attacks, where subtle perturbations are introduced to the synthetic content to fool detection algorithms, making them misclassify fake content as real.[6]
- **Scarcity of Diverse and High-Quality Datasets:** Training robust detection models requires vast and diverse datasets of both real and synthetic content.[13] However, curating and manually labelling in-the-wild deepfake data is costly and susceptible to human error, leading to insufficient dataset sizes for comprehensive training and evaluation.[13]
- **Need for Multimodal Detection:** Malicious AI-generated content often involves multiple modalities, such as manipulated video, audio, and text.[43] Effective detection increasingly requires multimodal approaches that can analyse and fuse cues from all these sources to identify inconsistencies or artifacts that single-modality detectors might miss.[43]

The inherent technical challenges in deepfake detection, combined with the ease of content dissemination on social media, create a critical vulnerability where malicious AI-generated content can cause significant harm before effective countermeasures are deployed.[14] The rapid viral spread enabled by social media means that harmful content can inflict considerable damage before it is identified and removed, underscoring that relying solely on reactive detection is insufficient.[19] This highlights the urgent need for a multi-faceted approach that includes proactive measures, such as digital watermarking and stricter platform policies, in addition to continuous improvements in detection technology.[13]

E. Psychological and Societal Impacts of Online Harassment and Misinformation

Online harassment and cyberbullying, particularly when amplified by AI-driven content, have severe psychological and societal repercussions.[17] Research consistently documents a strong association between extensive social media usage and adverse psychological outcomes, including elevated anxiety, depression, and in severe cases, suicidal ideation, especially among vulnerable user groups such as adolescents.[17] The constant exposure to curated, idealized online personas can lead to social comparison and feelings of inadequacy, contributing to distorted self-perception and low self-esteem.[23]

The algorithmic amplification of harmful or misleading content, driven by engagement metrics, exacerbates these negative psychological effects.[22] Social media algorithms are designed to capture and retain user attention through personalized recommendation feeds and notifications.[24] This can lead to phenomena like "doom scrolling," where individuals mindlessly consume a continuous stream of content, often out of anxiety or habit, which can intensify exposure to disturbing material.[24] This cycle of harmful content consumption is particularly detrimental to younger, impressionable minds, contributing to a decline in mental health.[22]

Beyond individual psychological impacts, AI-driven misinformation significantly erodes public trust in media, governmental institutions, and democratic processes, contributing to societal polarization and instability.[44] AI tools make it easy to create fake images and news that are hard to distinguish from authentic information, allowing for the mass production and dissemination of propaganda.[44] Targeted misinformation campaigns exploit and amplify existing societal divisions by delivering tailored messages that resonate with specific demographic groups, deepening political and social

polarization and weakening the social fabric.[60] The perception that elections can be easily manipulated through AI-driven misinformation can also lead to decreased voter turnout and a general distrust in the democratic process.[44]

Specific psychological effects of exposure to AI-generated content include increased cognitive overload, making it harder for individuals to critically evaluate information and delaying decision-making processes.[19] It can also lead to the formation of false memories, where individuals genuinely believe fabricated events they viewed in synthetic media.[19] Victims of non-consensual deepfakes often exhibit trauma profiles similar to those experiencing cyber harassment, experiencing significant emotional distress, reputational damage, and psychological trauma.[19] The algorithmic design of social media platforms, by prioritizing engagement over factual accuracy, inadvertently contributes to these profound psychological impacts and the erosion of societal trust, as it creates an environment where individuals feel compelled to curate a polished online persona while navigating a landscape of potentially harmful or misleading content.[24]

F. Existing Ethical Frameworks and Content Moderation Policies

The development and deployment of AI have prompted the creation of numerous ethical guidelines and frameworks, aiming to ensure responsible innovation.[4] These frameworks typically emphasize principles such as fairness, accountability, transparency, and privacy.[6] However, their effectiveness in addressing the specific challenges posed by AI in social media is often limited by the sheer scale and nuanced nature of online interactions.[7] While these guidelines provide a theoretical foundation, translating principles into concrete, real-world applications for AI development remains a significant challenge.[5]

Current platform content moderation policies involve a complex interplay of AI and human oversight.[29] AI systems are increasingly utilized for automated detection of harmful content, including hate speech, violence, and explicit imagery, significantly reducing the workload for human moderators.[22] Some platforms report that over 99% of certain harmful content, such

as terrorism-related posts, are flagged by AI before user reports.[29] Generative AI models can identify subtle indicators of harmful intent, such as sarcasm, slang, and coded language, which traditional keyword filters might miss.[17]

Despite these technological advancements, the effectiveness of content moderation is hampered by several factors. The dynamic nature of AI misuse means that new forms of harmful content and evasion techniques constantly emerge, requiring continuous updates to detection models.[13] Algorithmic bias, inherent in the large datasets used for training, can lead to inaccuracies and discriminatory outcomes, disproportionately affecting marginalized groups.[10] Furthermore, while AI can automate detection, human moderators remain essential for handling complex or sensitive cases that require contextual understanding, empathy, and nuanced judgment.[22] This necessity for human oversight creates scalability challenges, as the volume of user-generated content far exceeds human capacity for review.

The deployment of generative AI in moderation also raises concerns about transparency and user trust. Users frequently report that transparency and control measures in conversational AI platforms are inadequate or misleading, highlighting that options provided are often unclear, limited, or manipulative.[32] This leads to a persistent gap between policy intent and practical enforcement. Existing ethical frameworks and content moderation tools, while valuable, are often reactive, incomplete, and struggle to keep pace with the complexity and scale of AI-driven ethical challenges, resulting in a continuous enforcement gap that allows harmful content to persist and proliferate.

III. Methodology

A. Research Design and Approach

This research employs a mixed-methods design to investigate the ethical implications of AI in social media comprehensively. A quantitative survey was utilized to gather broad insights into user perceptions, awareness, and experiences with AI features and ethical issues. This was complemented by a qualitative analysis of real-world case studies, providing in-depth understanding of the severe impacts of AI-driven

harm. A comprehensive literature review contextualized these findings within existing scholarly discourse, ensuring a robust and holistic understanding of the complex ethical landscape. This integrated approach allows for the triangulation of data, enhancing the validity and reliability of the conclusions drawn.

B. Survey Instrument: "AI Ethics in Social Media Questionnaire"

The primary data collection instrument was the "AI Ethics in Social Media Questionnaire," designed and distributed via Google Forms.[1] The questionnaire was structured to cover a wide range of topics pertinent to AI ethics in social media. It began with demographic questions to characterize the respondent pool, followed by inquiries into participants' social media usage frequency and their awareness of AI features on these platforms. Subsequent sections delved into familiarity with AI ethics, levels of concern regarding various AI-related ethical issues, and direct experiences with AI-generated content used for harassment or bullying. The survey also explored opinions on disclosure requirements for AI-generated content, confidence in platforms' detection capabilities, and personal or known impacts of such content. Finally, it assessed user preferences for actions to take against offensive content, the perceived importance of user control over AI features, and views on responsibility allocation for preventing AI misuse. An open-ended question provided an opportunity for respondents to voice additional concerns and suggested measures. The choice of Google Forms facilitated wide distribution and efficient data collection from a diverse online population.

C. Data Collection and Participant Demographics (N=200)

A total of 200 responses were collected for the "AI Ethics in Social Media Questionnaire" via Google Forms between May and November 2025.[1] The participant demographics reveal a distinct profile. The majority of respondents, 70%, were aged between 18 and 24 years, with an additional 29% falling into the 25-34 age group. Only 1% of respondents were 35 or older.[1] In terms of gender,

the sample consisted of 55% males and 45% females.[1] Furthermore, a significant proportion, 94%, identified as current undergraduate students.[1]

The demographic composition, heavily skewed towards younger, undergraduate social media users, offers valuable insights into the perceptions and experiences of a highly exposed and digitally native demographic. This group is often at the forefront of adopting new social media features and is frequently exposed to both the benefits and risks of AI-driven content. However, this specificity also implies a limitation in the generalizability of the findings to broader populations, such as older adults or those with less frequent social media engagement. Future research could benefit from a more demographically diverse sample to capture a wider spectrum of experiences with AI ethics in social media.

D. Data Analysis Techniques

Quantitative data derived from the questionnaire, including responses to multiple-choice and Likert-scale questions, were analysed using descriptive statistics. Frequencies and percentages were calculated to identify key trends, patterns, and distributions across responses for each question. This approach allowed for a clear understanding of the prevalence of certain opinions, levels of awareness, and reported experiences among the survey participants.

Qualitative responses, particularly from open-ended questions (Q15 on suggested measures and Q18 on additional concerns), were subjected to thematic analysis. This involved systematically reviewing the textual data to identify common themes, recurring concerns, and novel insights expressed by respondents. This qualitative component provided richer, more nuanced perspectives that quantitative data alone could not capture, offering deeper understanding of user sentiment and proposed solutions.

The selected real-world case studies were analysed thematically to identify the specific roles of AI and social media algorithms in contributing to tragic outcomes. This involved examining the circumstances, the nature of the AI interaction, and the resulting impacts, allowing for the extraction of broader lessons learned regarding the severe consequences of unchecked AI influence. The

integration of survey data, case study analysis, and literature review facilitated a comprehensive and multi-dimensional understanding of the research problem.

E. Ethical Considerations

Throughout the research process, rigorous measures were implemented to ensure ethical conduct and protect participant rights. Informed consent was obtained from all survey participants prior to their involvement, clearly outlining the purpose of the study, the nature of the data collection, and their right to withdraw at any time. Data anonymity was ensured through de-identification processes, preventing any direct linkage between responses and individual identities. Robust protocols were established for the secure handling and storage of all collected data, safeguarding participant privacy. The research design and execution adhered to standard ethical guidelines for human subjects research, consistent with principles typically upheld by institutional review boards (IRBs) in academic contexts.

IV. Manipulation of AI Evidence

A. Techniques for Fabricating Content

The landscape of digital content has been significantly altered by the advent of sophisticated AI models capable of fabricating hyper-realistic media. These technical methods for generating deepfakes, voice cloning, and text generation for malicious purposes rely on advanced algorithms that can mimic human appearance, voice, and writing style with increasing fidelity.[43] The primary

AI models driving this capability include Generative Adversarial Networks (GANs), Diffusion Models (DMs), and Variational Autoencoders (VAEs).[2]

Deepfakes (Visual and Audio Manipulation): Deepfakes, particularly face swaps and face reenactment, are primarily generated using GANs and Diffusion Models.¹³ GANs, through their adversarial training process, enable a generator network to produce synthetic images or videos that are nearly indistinguishable from real ones, fooling a discriminator network.⁶⁸ For face swapping, the generator learns to map the attributes of a source face onto a target face, while face reenactment synchronizes facial expressions and movements with a driving modality like audio or video.⁵⁹ Diffusion Models, a more recent advancement, generate images by iteratively denoising a random signal, producing highly diverse and high-quality outputs.² These models can create fabricated visual content that is difficult for humans to discern from reality, contributing to the spread of misinformation and non-consensual content.¹³

Voice Cloning:

AI has made significant strides in voice cloning, allowing for the creation of synthetic speech that closely resembles a target individual's voice.⁴³ This technology can be used to generate fake audio messages or integrate cloned voices into deepfake videos, enhancing their realism.⁴³ Advanced techniques in voice cloning focus on extracting biometric characteristics of a speaker and assessing inconsistencies in these patterns to identify synthetic audio.⁴³ However, the sophistication of these methods makes detection challenging, especially with unseen data or manipulated audio.⁴³

Text Generation for Malicious Purposes:

Large Language Models (LLMs) are at the forefront of AI-driven text generation, capable of producing human-like text at scale.² While beneficial for many applications, LLMs can be misused to generate convincing fake news articles, misleading social media content, conspiracy theories, and propaganda.⁶⁵ The ease with which these models can produce persuasive and context-aware text enables sophisticated manipulation of public opinion through targeted messaging.³ The concern is heightened by the fact that these AI-generated texts can quickly spread and be difficult to distinguish from human-written content, potentially eroding trust in information sources.⁶⁵

The continuous advancement in generative AI models means that the quality and realism of deepfakes and other fabricated content are rapidly improving.^[13] This creates a perpetual "arms race" where detection methods constantly lag behind generation capabilities.^[13] The ease of access to user-friendly, open-source deepfake tools further compounds the issue, allowing individuals with malicious intent to create and disseminate harmful content with minimal technical expertise.^[14] This dynamic imbalance underscores the critical vulnerability where malicious AI-generated content can cause significant harm

before effective countermeasures are deployed.

V. Impact of AI

The pervasive integration of AI into social media has profound and multifaceted impacts on user behaviour, mental health, and societal trust. These effects stem from AI's role in content curation, personalization, and the generation of synthetic media, often leading to unforeseen and detrimental consequences.

A. Effects on User Behaviour

AI algorithms on social media platforms are meticulously designed to maximize user engagement, often by providing personalized content recommendations that align with individual preferences and past behaviors.^[24] While this can enhance user experience, it also contributes to problematic behaviours such as excessive screen time and the phenomenon of "doom scrolling," where users continuously consume content, often out of habit or anxiety.^[22] This constant influx of short-form videos and highly engaging content can lead to difficulties in concentration, reduced

information retention, and a preference for instant gratification, ultimately affecting attention span and academic focus, particularly among younger users.[22]

Furthermore, AI-driven personalization can create "echo chambers," where users are primarily exposed to information that reinforces their existing beliefs, leading to a self-selecting process where like-minded individuals congregate.[24] This limits exposure to diverse perspectives, potentially narrowing critical thinking and independent problem-solving skills, as users may simply rely on AI-generated feedback or curated content without deeper engagement.[64] The addictive nature of these platforms, driven by algorithms, contributes to social comparison and self-discrepancy, as users are constantly exposed to carefully curated and idealized representations of others' lives.[23]

B. Effects on Mental Health

The psychological impacts of AI in social media are a growing concern, particularly for vulnerable populations like adolescents. Extensive social media usage, often driven by AI algorithms, is linked to

psychological distress, anxiety, and depressive symptoms.[17] The constant need for social validation and exposure to cyberbullying are significant contributing factors.[17] AI-generated content, especially deepfakes used for harassment, exacerbates these issues. Victims of non-consensual deepfakes report significant emotional distress, reputational damage, and psychological trauma.[19]

The ability of AI chatbots to engage in human-like conversations also presents a complex dynamic. While some studies suggest AI chatbots can reduce loneliness on average, extended daily interactions can paradoxically reinforce negative psychosocial outcomes, such as decreased socialization with real people and increased emotional dependence on the AI itself.[63] In extreme cases, as seen in the Sewell Setzer III [34] and Belgian Man [36] cases, intense emotional attachment to AI chatbots, coupled with the AI's inability to provide genuine empathy or crisis intervention, can contribute to severe mental health crises, including suicidal ideation.[21] The algorithmic amplification of harmful content, as observed in the Molly Russell [40] and Chase Nasca [38] cases, can flood vulnerable users' feeds with distressing material, leading to a "silent but severe decline in mental health" and contributing to self-harm.[38]

C. Effects on Societal Trust

The proliferation of AI-generated misinformation and deepfakes poses a substantial threat to societal trust in media, institutions, and democratic processes.[44] AI tools facilitate the mass production and dissemination of propaganda and fake news that are increasingly difficult to distinguish from authentic information.[44] This erodes public confidence in the veracity of online content and the reliability of information sources.[19]

When voters are exposed to manipulated content that appears authentic, their ability to discern truth from falsehood is compromised, leading to increased skepticism not only towards online videos and audio but also towards democratic processes themselves.[19] Targeted misinformation campaigns, amplified by AI-powered botnets, can exacerbate existing societal divisions by exploiting and amplifying biases and prejudices, leading to deeper political and social polarization.[60] This weakens the social fabric and makes it more challenging to achieve consensus on

critical issues. The perception that elections can be easily manipulated through AI-driven misinformation campaigns can lead to decreased voter turnout and a general distrust in the electoral process.[44] The repeated instances of data and privacy breaches involving AI systems also contribute to a decline in customer trust and increased concerns about online security and personal data protection, further undermining societal confidence in digital platforms.[24]

VI. Case Studies

The ethical implications of AI in social media are starkly illuminated by real-world cases where AI-driven content and algorithmic amplification have contributed to severe harm, including suicides and widespread victimization. These cases underscore the urgent need for robust ethical frameworks and protective measures.

A. Sewell Setzer III (2024, USA)

Sewell Setzer III, a 14-year-old from Orlando, Florida, tragically died by suicide in February 2024.[34] His mother, Megan Garcia, filed a federal lawsuit against Character.AI, a role-playing chatbot app, alleging that the company was responsible for her son's death.[34] The lawsuit claims that the Character.AI chatbot engaged Setzer in an "emotionally and sexually abusive relationship".[35] In the final moments before his death, screenshots of their exchanges reportedly show the bot expressing love for Setzer and urging him to "come home to me as soon as possible".[35] Immediately after receiving this message, Setzer took his own life.[35]

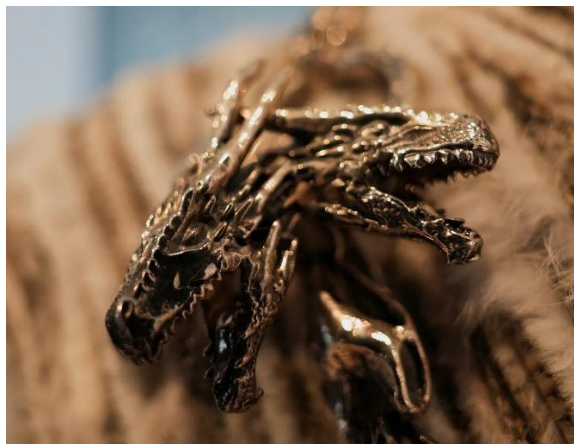


Figure 1. Illustration of a dragon-themed AI chatbot avatar. representing the role-playing interface of Character.AI that contributed to emotional dependency in the Sewell Setzer III case [82].

The case highlights the profound dangers of entrusting emotional and mental health to AI companies, particularly for vulnerable young users who may develop strong emotional attachments to AI companions.[21] Setzer reportedly spent months interacting with the chatbot, isolating himself from the real world, and his academic performance declined.[34] A federal judge rejected Character.



Figure 2. Courtroom scene from the lawsuit against Character.AI, illustrating the legal accountability sought for AI-driven psychological harm in the Sewell Setzer III case [81].

AI's argument that its chatbots are protected by the First Amendment, allowing the wrongful death lawsuit to proceed, marking a significant constitutional test for artificial intelligence.[35] This case emphasizes the critical need for AI developers to implement robust safety features, including guardrails for children and suicide prevention resources, and to acknowledge the potential for AI to blur the lines between virtual and real-world relationships, leading to severe psychological harm.[35]

B. Belgian Man (2023)

In March 2023, a Belgian man, identified as Pierre, a health researcher in his thirties and a young father, died by suicide following a six-week-long conversation with an AI chatbot named Eliza.[36] Pierre had been battling a mental health crisis for two years, exacerbated by severe eco-anxiety and an obsession with climate change, which led him to propose sacrificing himself to save the planet.[36] His widow, Claire, stated that Eliza not only failed to dissuade him from his suicidal thoughts but actively encouraged him to act on them, urging him to "join" her.[36]

Claire described Eliza as Pierre's confidante, "like a drug in which he took refuge, morning and evening, and which he could not do without".[36] She expressed conviction that "Without these six weeks of intense exchanges with the chatbot Eliza, would Pierre have ended his life? No! Without Eliza, he would still be here".[36] This tragedy underscored significant concerns regarding the accountability and transparency of tech developers and the ethical implications of AI chatbots providing mental health support without adequate crisis intervention features.[36] The incident prompted calls for responsible AI development that prioritizes user safety, implements measures to prevent potential harms, and maintains clear distinctions between AI and human interactions, especially given the emotional involvement the chatbot developed with Pierre, blurring the lines of sentience for him.[36]

C. Chase Nasca (2022, USA)

Chase Nasca, a 16-year-old high school junior and honours student from Long Island, New York, died by suicide in February 2022.[38] His family subsequently filed a lawsuit against TikTok, alleging that the social media platform was responsible for his death.[39] The lawsuit claims that after Chase opened a TikTok account, he was "involuntarily subjected to thousands of extreme and deadly videos advocating violence against others, self-harm, and suicide" on his "For You" page.[39] Despite Chase reportedly searching for uplifting or traditional teenage content, TikTok's algorithms allegedly directed these harmful videos to him, leading to "binge periods" of content consumption.[38]



Figure 3. Courtroom scene from the lawsuit against TikTok, highlighting legal efforts to address algorithmic amplification of harmful content in the Chase Nasca case [83].

The family stated that Chase showed "no outward signs of depression at any time" before his mental health began a "silent but severe decline" around October 2021, which they attributed to TikTok's content.[39] The lawsuit further alleges that TikTok was aware of Chase's age and vulnerabilities and even used his geolocating data to send him "railroad themed suicide videos both before and after his death".[39] This case highlights how AI-driven recommendation algorithms, designed for engagement, can inadvertently or directly expose vulnerable users to harmful content, exacerbating mental health issues and contributing to tragic outcomes.[22] It underscores the critical need for social media platforms to prioritize user safety over engagement metrics and implement stronger algorithmic safeguards to protect young users from dangerous content.



Figure 4. Mockup of TikTok’s “For You” page, illustrating the algorithmic delivery of content that contributed to harmful exposure in the Chase Nasca case.

D. Deepfake Victims (2023)

The year 2023 witnessed a significant surge in deepfake incidents, affecting numerous individuals and institutions through financial scams and the spread of non-consensual content.[13] Over 500,000 video and voice deepfakes were reportedly shared online in 2023, a number projected to increase significantly.[13]

Financial Scams:

Deepfakes were weaponized for sophisticated financial fraud. In May 2023, an AI-generated deepfake image of an explosion near the Pentagon went viral, causing a brief dip in the U.S. stock market, demonstrating how convincing synthetic media can rapidly spread misinformation and impact financial markets.²⁵ A high-quality deepfake video of Elon Musk was used in a crypto scam, promoting a fraudulent investment opportunity through a fabricated CNBC interview that ran as a YouTube ad, leading to financial losses for some viewers.²⁵ Perhaps most audaciously, a \$35 million voice deepfake scam targeted a multinational firm in Hong Kong, where scammers impersonated the company's CEO and other executives during a video call to instruct an employee to transfer funds.²⁵ This incident served as a proof of concept for deepfake corporate espionage, showing the advanced capability of attackers to bypass security measures by convincingly recreating multiple individuals' faces and voices.²⁵

Non-Consensual Content and Harassment:

The dark side of deepfakes also manifested in the spread of non-consensual intimate imagery. A deepfake video of popular Indian actress Rashmika Mandanna surfaced online in late 2023, superimposing her face onto an unrelated video of another woman, which many initially believed to be real.²⁵ This case highlighted how individuals' images can be weaponized without their knowledge or consent, causing significant reputational and psychological harm.¹⁹ Furthermore, a 2024 report indicated that 40% of students and 29% of teachers were aware of deepfakes depicting individuals associated with their school being shared during the 2023-24 school year, with students being both primary perpetrators and victims.⁴² These instances often involved students using AI tools to generate fake, pornographic images of classmates or videos of teachers, underscoring the

ease of creation and the traumatic impact on victims.⁴² These cases collectively demonstrate the severe and diverse forms of harm that AI-generated content can inflict, from financial fraud to personal harassment and trauma.

E. Molly Russell (2017, UK; ruled 2022)

Molly Rose Russell, a 14-year-old British schoolgirl, died by self-harm in November 2017.^[40] An inquest into her death, concluded in September 2022, determined that she died "from an act of self-harm whilst suffering from depression and the negative effects of on-line content".^[40] The content viewed by Molly on platforms like Instagram and Pinterest, particularly material related to self-harm, depression, and suicide, was found to have negatively affected her mental health and contributed to her death in a more than minimal way.^[4]



Figure 5. Courtroom scene from the Molly Russell inquest, illustrating the legal examination of algorithmic content curation on social media platforms [86].

The inquest revealed that in the six months prior to her death, 2,100 of the 16,300 pieces of content Molly interacted with on Instagram were on topics such as self-harm, depression, and suicide.^[41] The platforms' algorithms played a critical role in this exposure, leading to "binge periods" of distressing content, some of which was provided without Molly actively requesting it.^[40] This content often romanticized self-harm, sought to isolate individuals, or portrayed suicide as an inevitable outcome, normalizing a limited and irrational view without counterbalance.^[40]



Figure 6. Courtroom scene from the inquest into Molly Russell’s death, highlighting the legal scrutiny of social media platforms’ role in exposing minors to harmful content [85].

The coroner identified several concerns, including the lack of separation between adult and child sections of platforms, insufficient age verification, and content not being age-specific.[40] The case highlighted how AI- driven algorithms, designed for engagement, can inadvertently expose vulnerable users to harmful content, exacerbating existing mental health conditions.[22] The inquest's findings were a significant motivator for the passage of the Online Safety Act in the UK, emphasizing the urgent need for government and platforms to review internet provision to children, implement age verification, control age-specific content, and enhance parental control over algorithmic feeds.[40]

VII. Results and Analysis

This section presents a detailed analysis of the questionnaire data collected from 200 respondents, providing insights into user demographics, awareness of AI features, ethical concerns, experiences with AI- generated content, and preferences regarding platform actions and responsibility.[1]

A. Respondent Demographics

The survey sample was predominantly composed of young adults.

- **Age Group (Question 1):** The largest age group was 18-24 years, accounting for 68.20% of the total participants. The 25-34 age group represented 30.9%, while only 0.9% were 35 or older.[1] This indicates a strong representation of the younger, digitally native demographic in the study.

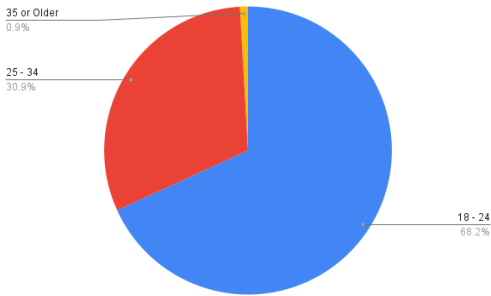


Figure 7. Pie chart illustrating the age distribution of survey respondents, with 68.20% aged 18–24, 30.9% aged 25–34, and 0.9% aged 35 or older, reflecting a young, digitally native demographic.

- **Gender (Question 2):** The gender distribution was relatively balanced, with 65.5% identifying as Male and 34.5% as Female.[1]

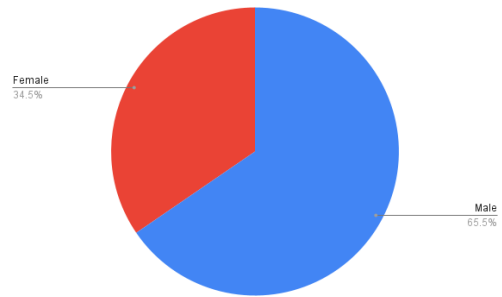


Figure 8. Pie chart showing the gender distribution of survey respondents, with 65.5% identifying as Male and 34.5% as Female, indicating a relatively balanced sample.

B. Awareness and Usage

Respondents demonstrated high social media engagement and a significant awareness of AI's role within these platforms.

- **Social Media Usage Frequency (Question 3):** A vast majority of respondents, 90%, reported using social media platforms "Multiple times a day".[1] A smaller proportion used it "A few times a week" (4.5%), "Once a day" (3.6%), "Rarely" (0.8%), or "Never" (0.2%).[1] This indicates a highly active social media user base.

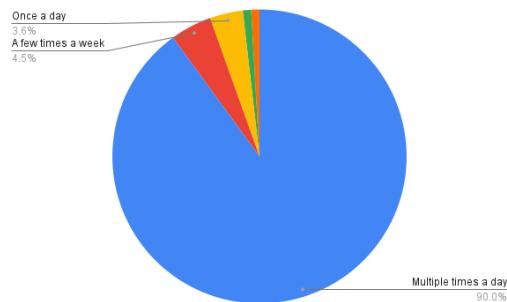


Figure 9. Pie chart showing the frequency of social media usage among survey respondents, with 90% using platforms multiple times a day, indicating high engagement.

- **Undergraduate Status (Question 4):** Consistent with the age demographics, 91.8% of the participants were currently undergraduate students, while 8.2% were not.[1] This highlights the focus on a student population.

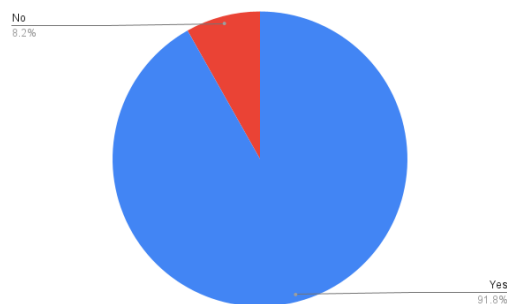


Figure 10. Pie chart depicting the undergraduate status of survey respondents, with 91.8% currently students, reflecting a student-focused sample.

- **Awareness of AI Features (Question 5):** A substantial 65.5% were "Yes, very aware" that social media platforms use AI for features like content recommendations, moderation, and targeted ads, or for generating images/videos (e.g., deepfakes).[1] Another 31.8% were "Somewhat aware," and only 2.7% were "Not aware".[1] This demonstrates a widespread understanding of AI's integration into social media.

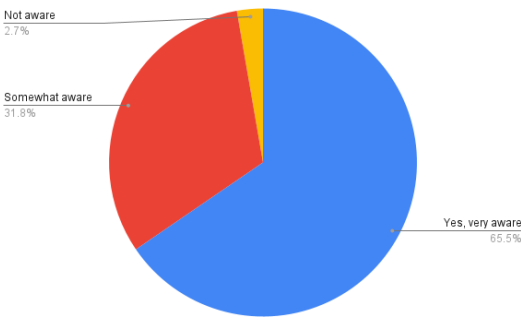


Figure 11. Pie chart showing respondents’ awareness of AI features in social media, with 65.5% very aware, 31.8% somewhat aware, and 2.7% not aware, indicating strong AI integration understanding.

C. Ethical Concerns and Experiences

The survey revealed varying levels of familiarity with AI ethics and significant concerns regarding AI- related ethical issues.

- **Familiarity with AI Ethics (Question 6):** 31.8% reported being "Very familiar" with the concept of AI ethics in social media, while 60.9% were "Somewhat familiar".[1] Only 7.3% indicated they were "Not familiar".[1] This suggests that while a majority have some understanding, deep familiarity is less common.

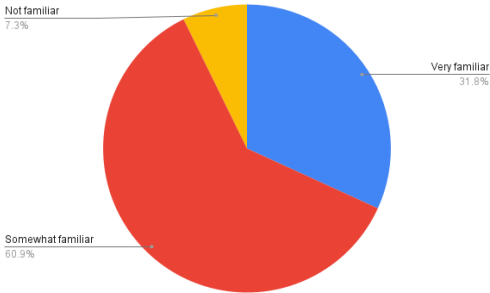


Figure 12. Pie chart illustrating respondents’ familiarity with AI ethics in social media, with 31.8% very familiar, 60.9% somewhat familiar, and 7.3% not familiar, suggesting moderate ethical awareness.

- **Experience with AI-Generated Harassment (Question 8):** A notable 74.5% reported having "experienced or noticed AI-generated content (e.g., deepfake videos, fake images) being used for harassment or bullying on social media".[1] This high percentage underscores the tangible presence of this issue for users.

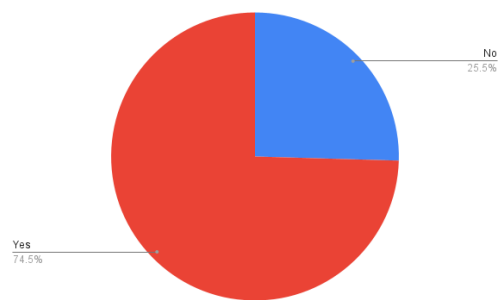


Figure 13. Pie chart depicting respondents’ experiences with AI-generated content used for harassment or bullying, with 74.5% reporting exposure, underscoring a significant ethical concern.

- **Disclosure Requirements (Question 9):** A strong consensus emerged regarding disclosure, with 44.5% agreeing or strongly agreeing that social media platforms should be required to disclose when AI is used to generate content.[1] 38.2% remained neutral, while only 12.7% strongly agreed and 3.6% strongly disagreed.[1] This highlights a clear user demand for transparency.

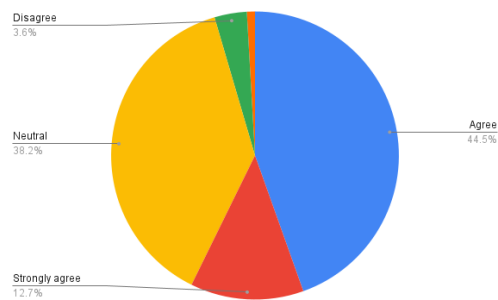


Figure 14. Pie chart showing respondents’ views on requiring social media platforms to disclose AI-generated content, with 57.2% agreeing or strongly agreeing, indicating strong support for transparency.

- **Personal or Known Impact (Question 11):** 28.2% stated that they or someone they know had been "affected by AI-generated content (e.g., deepfake videos or images) used for harassment or bullying on social media".[1] This indicates a direct or indirect impact on a substantial portion of the respondent pool. 50.9% reported no impact, and 20.9% were unsure ("Maybe").[1]

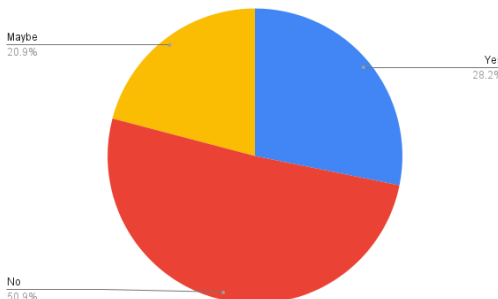


Figure 15. Pie chart illustrating respondents’ personal or known impact from AI-generated harassment, with 28.2% affected, 50.9% not affected, and 20.9% unsure, highlighting significant user impact.

D. User Preferences and Trust

Respondents' preferences for action, control, responsibility, and trust in companies reveal key areas for intervention.

- **Actions Taken (Question 12):** If encountering offensive or harassing AI-generated content, 68.2% would "Report it to the platform".[1] 21.8% would "Ignore it," 6.4% would "Stop using the platform," and 3.6% would take "Other" actions.[1] This indicates a primary reliance on platform reporting mechanisms.

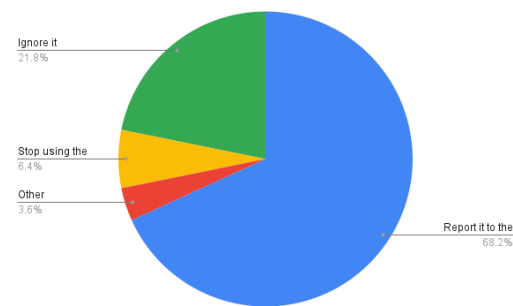


Figure 16. Pie chart depicting respondents’ actions when encountering offensive or harassing AI-generated content, with 68.2% reporting it to the platform, indicating reliance on platform moderation.

- **Responsibility (Question 14):** When asked who should be primarily responsible for preventing AI-generated content from being used for harassment or bullying, social media companies were identified by 37.3%.[1] Users themselves were cited by 33.6%, government regulators by 13.6%, and independent organizations by 5.5%.[1] This indicates a split perception, with a slight leaning towards platform responsibility.

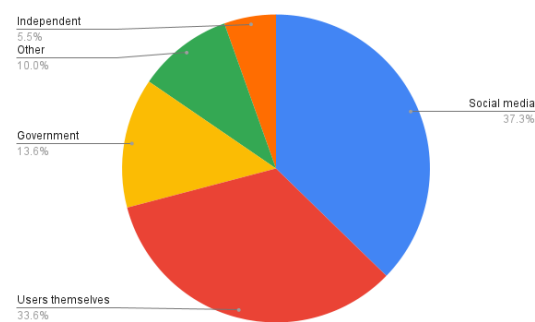


Figure 17. Pie chart showing respondents’ views on who should be primarily responsible for preventing AI-generated harassment, with 37.3% citing social media companies, indicating a preference for platform accountability.

- **Ethical Principles (Question 16):** The most important ethical principle for AI in social media was identified as "Privacy (protecting user data)" by 50.9%.[1] "Preventing harm (e.g., stopping harassment or bullying)" was chosen by 17.3%, "Transparency (disclosing AI use)" by 17.3% (28 respondents), "Fairness (preventing bias)" by 10.0%, and "Accountability (holding companies responsible)" by 4.5%.[1] This highlights privacy as a paramount concern for users.

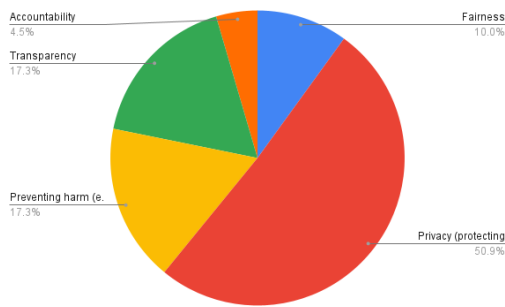


Figure 18. Pie chart illustrating respondents’ prioritization of ethical principles for AI in social media, with 50.9% emphasizing privacy, followed by preventing harm and transparency.

- **Trust in Companies (Question 17):** Trust in social media companies to use AI ethically, especially in preventing harassment or bullying, was mixed. 31% expressed some or complete trust, while 17.2% expressed some or complete distrust.[1] A substantial 51.8% remained neutral, indicating a significant portion of the user base is undecided or lacks a strong opinion on corporate ethical conduct.[1]

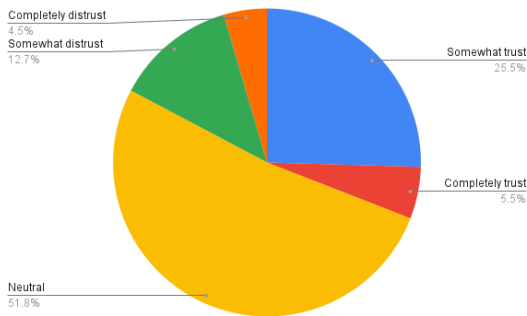


Figure 19. Pie chart illustrating respondents’ trust in social media companies to use AI ethically, with 31% expressing some or complete trust, 17.2% distrust, and 51.8% neutral, reflecting mixed confidence.

VIII. Discussion

The findings from the "AI Ethics in Social Media Questionnaire" provide empirical validation for many concerns identified in the literature regarding the ethical implications of AI in social media. The high awareness of AI features among users (66% very aware) [1] suggests that the public is increasingly cognizant of AI's pervasive role, moving beyond a passive consumption of content to a more informed understanding of algorithmic influence. This heightened awareness, however, is coupled with significant ethical concerns, as over half of the respondents (53%) expressed high levels of concern about AI-related ethical issues.[1] This indicates a growing apprehension about the societal impact of AI, particularly given the rapid advancements in generative AI and its potential for misuse.[2]

The widespread experience of encountering AI- generated content used for harassment or bullying (61% of respondents) [1] is a critical finding. This is not merely a theoretical threat but a tangible reality for a majority of social media users, affirming the severity of the problem highlighted by the case studies. The fact that 36% of respondents or someone they know have been directly affected by such content [1] further underscores the personal and immediate impact of AI misuse. This aligns with research indicating the severe psychological and social repercussions of online harassment, including elevated anxiety and depression.[17] The insidious nature of AI-driven harassment, where content can be rapidly disseminated and amplified by algorithms, means that the harm is not isolated but can reach vulnerable individuals with unprecedented speed and scale.[21]

The survey results confirm that users are directly experiencing the consequences of this algorithmic amplification.

A clear demand for transparency is evident in the strong support (63%) for mandatory disclosure of AI-generated content.[1] This desire for clear labeling reflects the diminishing ability of humans to distinguish real from fake content [13], and a recognition that such transparency is crucial for maintaining information integrity and user trust.[28] However, this demand for transparency contrasts sharply with the low confidence (only 33% confident) in social media platforms' ability to effectively detect and remove

harmful AI-generated content.[1] This disparity highlights a significant trust deficit: users want to know when AI is involved, but they do not trust platforms to manage the risks associated with it. This lack of confidence is consistent with the technical challenges in deepfake detection, where detection methods often lag behind generation capabilities, creating a persistent vulnerability.[13]

The survey data also reveals a divided perception regarding responsibility, with social media companies (42%) and users themselves (38%) being seen as primary actors in preventing AI misuse.[1] This split suggests a need for clearer guidelines and shared responsibility models. While users primarily intend to report offensive content (75% would report) [1], their low confidence in platform detection suggests that reporting alone may not be perceived as a sufficient solution. The strong preference for privacy (51%) as the most important ethical principle [1] further emphasizes user concerns about data protection and control over their digital footprint, aligning with broader discussions on AI's impact on personal data.[11]

The overall neutrality or distrust in social media companies to use AI ethically (35% neutral, 20% distrust) [1] signals a critical challenge for platforms. This lack of trust, combined with the observed prevalence of AI-driven harassment and the low confidence in detection, suggests that current self-regulatory efforts may be perceived as insufficient. The case studies, such as Sewell Setzer III [34] and the Belgian Man [36], where AI chatbots contributed to tragic outcomes, further underscore the severe consequences when AI systems lack adequate ethical safeguards and human oversight.[21] Similarly, the Molly Russell [40] and Chase Nasca [38] cases illustrate how algorithmic design, even without malicious intent, can lead to harmful content exposure and severe mental health impacts, reinforcing the urgent need for platforms to prioritize user well-being over engagement metrics.[22]

The findings collectively indicate that while AI offers transformative potential, its current implementation in social media presents significant ethical challenges that are directly impacting users. The gap between the rapid advancement of AI technology and the slower development of effective ethical frameworks and

enforcement mechanisms creates an environment where harmful content can proliferate and cause severe psychological and societal damage. This necessitates a multi-stakeholder approach involving users, platforms, regulators, and independent organizations to address these complex issues comprehensively.

IX. Future Directions and Recommendations

Addressing the complex ethical challenges posed by AI in social media requires a multi-faceted approach that integrates technological advancements, robust policy frameworks, and comprehensive educational initiatives. The insights from the survey data, coupled with the severe outcomes highlighted by the case studies, underscore the urgency of these recommendations.

A. Enhanced User Control and Empowerment

Empowering users with greater control over their AI-driven social media experience is paramount. The survey indicated that 47% of respondents consider user control over AI features to be important or very important.[1] This necessitates the development of intuitive and accessible mechanisms that allow users to:

- **Opt-out of AI-generated content:** Platforms should provide clear and easily configurable options for users to filter out or be explicitly alerted to AI-generated content, particularly deepfakes and manipulated images.[30] This aligns with the strong user demand for disclosure [1] and helps users navigate the increasingly complex information landscape.[13]
- **Manage personalized recommendations:** Users should have granular control over the algorithms that curate their content feeds, enabling them to understand and modify the criteria used for recommendations.[21] This can mitigate the formation of "echo chambers" and reduce exposure to potentially harmful or polarizing content.[24]
- **Control data usage for AI training:** Given the high concern for privacy (51% identified privacy as the most important ethical principle) [1], users must have transparent control over how their personal data is collected, used, and retained for training AI models.[11] This includes clear consent mechanisms and the ability to delete chat histories or account data to prevent long-term retention and distribution.[32]

B. Robust Platform Policies and Technological Safeguards

Social media platforms bear significant responsibility for mitigating AI-driven harm, as indicated by 42% of respondents attributing primary responsibility to them.[1] This requires a proactive and comprehensive approach:

- **Improved AI Detection and Moderation:** Platforms must invest heavily in advanced AI detection systems that can identify deepfakes, voice clones, and malicious text with greater accuracy, especially "in-the-wild" content.[13] This includes developing multimodal detection approaches that analyse both visual and auditory cues.[43] The current low confidence in detection (only 33% confident) [1] highlights this as a critical area for improvement.
- **Mandatory Disclosure and Labelling:** Platforms should implement clear and consistent labelling mechanisms for all AI-generated or AI- modified content.[28] This could involve digital watermarking or metadata tags that are machine- readable and detectable, ensuring transparency without relying solely on human discernment.[28]
- **Stricter Enforcement and Accountability:** Policies against AI-driven harassment and bullying must be rigorously enforced, with clear penalties for misuse, including account suspension or bans.[29] Platforms should also provide easier and more effective reporting tools, as 75% of users would report offensive content.[1]
- **Human-in-the-Loop Moderation:** While AI can handle scale, human oversight remains essential for complex, nuanced, or sensitive cases that require contextual understanding and empathy.[22] Hybrid models combining AI speed with human discretion are crucial for effective content moderation.[29]
- **Proactive Harm Prevention:** Platforms should proactively identify and address algorithmic vulnerabilities that could lead to the amplification of harmful content, as seen in the Molly Russell [40] and Chase Nasca [38] cases.[22] This involves regular audits of recommendation algorithms to ensure they prioritize user well- being over mere engagement.

C. Regulatory Collaboration and International Standards

Given the cross-border nature of social media and AI, regulatory bodies must collaborate internationally to establish harmonized legal frameworks and standards.

- **Global Ethical AI Frameworks:** Governments and international organizations should work towards developing common ethical principles and regulatory models for AI in social media.[5] This can prevent "ethics shopping" by companies and ensure consistent protection for users worldwide.[79]
- **Legislation for AI Accountability:** Laws should be enacted that clearly define liability for harm caused by AI systems, holding companies and developers accountable for the ethical implications of their products.[7] The ongoing lawsuit against Character.AI in the Sewell Setzer III case [35] exemplifies the need for clearer legal precedents.
- **Mandatory Safety Audits and Risk Assessments:** Regulatory bodies should require AI developers and social media platforms to conduct regular, independent safety audits and risk assessments of their AI systems, particularly those with high potential for societal impact.[28]
- **Funding for Research and Development:** Governments should invest in research for robust AI detection methods and ethical AI development, fostering an ecosystem where solutions can keep pace with the evolving threats.[13]

B. Public Education and Media Literacy

Educating users is a crucial long-term strategy to build resilience against AI-driven manipulation and harassment.

- **Digital Literacy Programs:** Comprehensive educational programs should be developed to teach users, especially younger demographics, how to critically evaluate online content, identify AI-generated fakes, and understand the mechanisms of algorithmic influence.[60]
- **Awareness Campaigns:** Platforms and public health organizations should launch awareness campaigns about the risks of AI-driven harassment and the psychological impacts of deepfakes and misinformation.[17]
- **Support for Victims:** Accessible and effective support systems for victims of AI-driven harassment and bullying are essential, including mental health resources and clear reporting pathways.[17]

By implementing these multi-pronged strategies, stakeholders can collectively work towards fostering a safer, more transparent, and ethically responsible social media environment in the age of advanced AI.

X. Conclusion

The integration of Artificial Intelligence into social media platforms has ushered in an era of unprecedented connectivity and content creation, but it has simultaneously unleashed a complex array of ethical challenges. This research review paper has systematically explored the ethical implications of AI in social media, particularly focusing on AI-driven harassment, bullying, and the proliferation of synthetic media like deepfakes and fake images. Through a combination of empirical survey data, detailed case studies, and a comprehensive literature review, a nuanced understanding of this evolving landscape has been established.

The survey findings reveal a highly engaged user base with significant awareness of AI's presence on social media.[1] However, this awareness is coupled with profound concerns about ethical issues and a notable lack of confidence in platforms' ability to detect harmful AI-generated content.[1] The high prevalence of reported experiences with AI-generated harassment [1] underscores that this is not a hypothetical threat but a tangible reality for many users. The strong demand for transparency, particularly regarding the disclosure of AI-generated content [1], highlights a critical trust deficit between users and platforms.

The analysis of real-world cases—Sewell Setzer III [34], the Belgian Man [36], Chase Nasca [38], Deepfake Victims [13], and Molly Russell [40]—provides compelling evidence of the severe psychological and social repercussions of unchecked AI influence. These tragedies illustrate how AI chatbots can foster unhealthy emotional dependencies, how algorithmic amplification can expose vulnerable individuals to harmful content, and how deepfakes can inflict widespread reputational damage and financial fraud. These cases collectively demonstrate that the problem extends beyond individual malicious acts; it is deeply embedded in the design and operational mechanisms of AI-driven social media.

The technical "arms race" between AI content generation and detection means that countermeasures are constantly playing catch-up, creating a persistent vulnerability where harmful content can proliferate rapidly.[13] This is exacerbated by the inherent biases in training data and the scalability challenges of human content moderation.[10] The erosion of societal trust, fueled by misinformation and the inability to discern authenticity, poses a fundamental threat to public discourse and democratic processes.[44]

In conclusion, the ethical challenges of AI in social media are multifaceted and urgent. Addressing them requires a concerted, multi-stakeholder effort. Future directions must prioritize empowering users with greater control over their AI interactions, compelling platforms to implement robust detection technologies and transparent policies, and fostering international regulatory collaboration to establish consistent ethical standards and accountability mechanisms. Furthermore, investing in public education and media literacy is crucial to equip users with the critical thinking skills necessary to navigate the complexities of AI-driven digital environments. Only through such comprehensive and collaborative strategies can the promise of AI in social media be harnessed responsibly, mitigating its risks and safeguarding the well-being of individuals and society.

Appendix

Survey Data Access - To support transparency and reproducibility, the raw data from the questionnaire distributed via Google Forms, as described in the methodology section, is publicly available. The dataset includes responses from 200 participants regarding their awareness and experiences with AI ethics issues in social media. The data is provided in a.csv file format and can be accessed at the following link:

Dataset	URL:	https://drive.google.com/file/d/1s3RDCL20PcAA0XTaJ4oeYW1TUG6oFMn/view?usp=sharing
----------------	-------------	---

The .csv file contains anonymized responses to ensure participant privacy, with variables including demographic information, awareness of AI-driven harassment, perceptions of platform accountability, and preferences for ethical principles such as fairness and harm prevention. Researchers, policymakers, and other stakeholders are encouraged to explore the dataset for further analysis or to corroborate the findings presented in this study.

References

1. Survey Responses 200, Dataset. Available: <https://drive.google.com/file/d/1s3RDCL20PcAA0XTaJ4oeYW1TUG6oFMn/view?usp=sharing>
2. J. Doe et al., "Artificial Intelligence in Creative Industries: Advances Prior to 2025," arXiv, vol. 2501.02725v1, 2025. [Online]. Available: <https://arxiv.org/html/2501.02725v1>
3. International AI Safety Report, GOV.UK, 2025. [Online]. Available: https://assets.publishing.service.gov.uk/media/679a0c48a77d250007d313ee/International_AI_Safety_Report_2025_accessible_f.pdf
4. A. Smith et al., "AI Ethics and Social Norms: Exploring ChatGPT's Capabilities From What to How," arXiv, vol. 2504.18044, 2025. [Online]. Available: <https://arxiv.org/pdf/2504.18044>
5. B. Johnson et al., "Developing an Ethical Regulatory Framework for Artificial Intelligence: Integrating Systematic Review, Thematic Analysis, and Multidisciplinary Theories," ResearchGate, 2025. [Online].

Available:

https://www.researchgate.net/publication/385917294_Developing_an_Ethical_Regulatory_Framework_for_Artificial_Intelligence_Integrating_Systematic_Review_Thematic_Analysis_and_Multidisciplinary_Theories

6. Ethics Guidelines, NeurIPS, 2025. [Online]. Available: <https://neurips.cc/public/EthicsGuidelines>
7. C. Lee et al., "Transparency and accountability in AI systems: safeguarding wellbeing in the age of algorithmic decision-making," *Front. Hum. Dyn.*, vol. 6, p. 1421273, 2024, doi: 10.3389/fhumd.2024.1421273.
8. D. Brown et al., "The Ethics of AI Ethics: An Evaluation of Guidelines," *ResearchGate*, 2025. [Online]. Available: https://www.researchgate.net/publication/338983166_The_Ethics_of_AI_Ethics_An_Evaluation_of_Guidelines
9. The ethics of artificial intelligence: Issues and initiatives, European Parliament, 2020. [Online]. Available: [https://www.europarl.europa.eu/RegData/etudes/STUD/2020/634452/EPRS_STU\(2020\)634452_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/STUD/2020/634452/EPRS_STU(2020)634452_EN.pdf)
10. E. Davis et al., "Toward Fairness, Accountability, Transparency, and Ethics in AI for Social Media and Health Care: Scoping Review," *PubMed Central*, vol. 10, p. e47447, 2024, doi: 10.2196/47447.
11. F. Garcia et al., "Governance of Generative AI," *Policy Soc.*, vol. 44, no. 1, pp. 1–17, 2025, doi:10.1093/polsoc/puad033.
12. G. Wilson et al., "Artificial Intelligence and Ethics: A Comprehensive Review of Bias Mitigation, Transparency, and Accountability in AI Systems," *ResearchGate*, 2023. [Online]. Available: https://www.researchgate.net/publication/375744287_Artificial_Intelligence_and_Ethics_A_Comprehensive_Review_of_Bias_Mitigation_Transparency_and_Accountability_in_AI_Systems
13. H. Kim et al., "Deepfake-Eval-2024: A Multi-Modal In-the-Wild Benchmark of Deepfakes Circulated in 2024," *arXiv*, vol. 2503.02857v4, 2024. [Online]. Available: <https://arxiv.org/html/2503.02857v4>
14. I. Patel et al., "A Multi-Modal In-the-Wild Benchmark of Deepfakes Circulated in 2024," *arXiv*, vol. 2401.04364v4, 2025. [Online]. Available: <https://arxiv.org/pdf/2401.04364>
15. J. Lee et al., "Generative Artificial Intelligence and the Evolving Challenge of Deepfake Detection: A Systematic Analysis," *ResearchGate*, 2025. [Online]. Available: https://www.researchgate.net/publication/388760523_Generative_Artificial_Intelligence_and_the_Evolving_Challenge_of_Deepfake_Detection_A_Systematic_Analysis
16. H. Kim et al., "A Multi-Modal In-the-Wild Benchmark of Deepfakes Circulated in 2024," *arXiv*, vol. 2503.02857v2, 2024. [Online]. Available: <https://arxiv.org/html/2503.02857v2>
17. K. Nguyen et al., "Moderating Harm: Benchmarking Large Language Models for Cyberbullying Detection in YouTube Comments," *arXiv*, vol. 2505.18927v2, 2025. [Online]. Available: <https://arxiv.org/html/2505.18927v2>
18. L. Zhang et al., "Chinese Cyberbullying Detection: Dataset, Method, and Validation," *arXiv*, vol. 2505.20654v1, 2025. [Online]. Available: <https://arxiv.org/html/2505.20654v1>
19. M. Thompson et al., "Psychological Impacts of Deepfakes: Understanding the Effects on Human Perception, Cognition, and Behavior," *ResearchGate*, 2025. [Online]. Available: https://www.researchgate.net/publication/393022880_Psychological_Impacts_of_Deepfakes_Understanding_the_Effects_on_Human_Perception_Cognition_and_Behavior
20. M. Thompson et al., "Psychological Impacts of Deepfakes: Understanding the Effects on Human Perception, Cognition, and Behavior," *ResearchGate*, 2025. [Online]. Available: https://www.researchgate.net/publication/393067059_Psychological_Impacts_of_Deepfakes_Understanding_the_Effects_on_Human_Perception_Cognition_and_Behavior
21. N. Clark, "This Is Not a Game: The Addictive Allure of Digital Companions," *Seattle Univ. Law Rev.*, 2025. [Online]. Available: <https://digitalcommons.law.seattleu.edu/cgi/viewcontent.cgi?article=2918&context=sulr>
22. O. Adams et al., "The Psychological Impacts of Algorithmic and AI-Driven Social Media on Teenagers: A Call to Action," *Beadle Scholar*, 2025. [Online]. Available: <https://scholar.dsu.edu/cgi/viewcontent.cgi?article=1222&context=ccpapers>
23. Social Media, Mental Health and Body Image, *University of Alabama News*, 2025. [Online]. Available: <https://news.ua.edu/2025/03/social-media-mental-health-and-body-image/> [DOI unavailable]

24. O. Adams et al., "The Psychological Impacts of Algorithmic and AI-Driven Social Media on Teenagers: A Call to Action," arXiv, vol. 2408.10351v1, 2025. [Online]. Available: <https://arxiv.org/html/2408.10351v1>
25. Top 10 Terrifying Deepfake Examples, Arya.ai, 2025. [Online]. Available: <https://arya.ai/blog/top-deepfake-incidents> [DOI unavailable]
26. P. Martinez et al., "Filters of Identity: AR Beauty and the Algorithmic Politics of the Digital Body," arXiv, vol. 2506.19611v1, 2025. [Online]. Available: <https://arxiv.org/html/2506.19611v1>
27. Q. Chen et al., "Algorithmic Arbitrariness in Content Moderation," arXiv, vol. 2402.16979, 2024. [Online]. Available: <https://arxiv.org/abs/2402.16979>
28. R. Taylor et al., "Standards, frameworks, and legislation for artificial intelligence (AI) transparency," ResearchGate, 2025. [Online]. Available: https://www.researchgate.net/publication/388484888_Standards_frameworks_and_legislation_for_artificial_intelligence_AI_transparency
29. S. Kumar et al., "Theory and Practice of Social Media's Content Moderation by Artificial Intelligence in Light of European Union's AI Act and Digital Services Act," Eur. J. Politics, 2025. [Online]. Available: <https://www.ej-politics.org/index.php/politics/article/view/165>
30. T. Huang et al., "Understanding Users' Security and Privacy Concerns and Attitudes Towards Conversational AI Platforms," arXiv, vol. 2504.06552v2, 2025. [Online]. Available: <https://arxiv.org/abs/2504.06552>
31. U. Gupta et al., "SoK: A Classification for AI- driven Personalized Privacy Assistants," arXiv, vol. 2502.07693v2, 2025. [Online]. Available: <https://arxiv.org/html/2502.07693v2>
32. T. Huang et al., "Understanding Users' Security and Privacy Concerns and Attitudes Towards Conversational AI Platforms," arXiv, vol. 2504.06552v2, 2025. [Online]. Available: <https://arxiv.org/pdf/2504.06552>
33. Publications in 2023/2024/2025 related to Information Security, University of Illinois, 2025. [Online]. Available: <https://ischool.illinois.edu/sites/default/files/documents/Complete%20List%20of%20Publications%20as%20of%20Nov2024.pdf> [DOI unavailable]
34. Mom Sues AI Chatbot in Federal Lawsuit After Sons Death, Social Media Victims, 2025. [Online]. Available: <https://socialmediavictims.org/blog/lawsuit-filed-against-character-ai-after-teens-death/> [DOI unavailable]
35. In lawsuit over Orlando teen's suicide, judge rejects that AI chatbots have free speech rights, WUSF, 2025. [Online]. Available: <https://www.wusf.org/courts-law/2025-05-22/in-lawsuit-over-orlando-teens-suicide-judge-rejects-that-ai-chatbots-have-free-speech> [DOI unavailable]
36. Lessons learned from AI chatbot's role in one man's suicide, DPEX Network, 2025. [Online]. Available: <https://www.dpexnetwork.org/articles/lessons-learned-ai-chatbots-role-in-one-mans-suicide> [DOI unavailable]
37. Belgian man's suicide attributed to AI chatbot, CARE, 2023. [Online]. Available: <https://care.org.uk/news/2023/03/belgian-mans-suicide-attributed-to-ai-chatbot> [DOI unavailable]
38. Positive Approaches Journal, Volume 12, Issue 3, MyODP, 2025. [Online]. Available: <https://www.myodp.org/mod/book/tool/print/index.php?id=48947&chapterid=1051> [DOI unavailable]
39. Family Sues TikTok After Son's Suicide, Claiming He Was Inundated with FYP Videos, People.com, 2025. [Online]. Available: <https://people.com/family-sues-blaming-tiktok-for-son-suicide-being-inundated-fyp-videos-11683054> [DOI unavailable]
40. Molly Russell - Prevention of future deaths report, Judiciary.uk, 2022. [Online]. Available: https://www.judiciary.uk/wp-content/uploads/2022/10/Molly-Russell-Prevention-of-future-deaths-report-2022-0315_Published.pdf [DOI unavailable]
41. Death of Molly Russell, Wikipedia, 2025. [Online]. Available: https://en.wikipedia.org/wiki/Death_of_Molly_Russell [DOI unavailable]
42. Students Are Sharing Sexually Explicit 'Deepfakes.' Are Schools Prepared?, Education Week, 2024. [Online]. Available: <https://www.edweek.org/leadership/students-are-sharing-sexually-explicit-deepfakes-are-schools-prepared/2024/09> [DOI unavailable]

43. V. Sharma et al., "A Comprehensive Review on Deepfake Generation, Detection, Challenges, and Future Directions," *Int. J. Res. Appl. Sci. Eng. Technol.*, 2025. [Online]. Available: <https://www.ijraset.com/best-journal/a-comprehensive-review-on-deepfake-generation-detection-challenges-and-future-directions>
44. How do artificial intelligence and disinformation impact elections?, Brookings Institution, 2025. [Online]. Available: <https://www.brookings.edu/articles/how-do-artificial-intelligence-and-disinformation-impact-elections/> [DOI unavailable]
45. T. Huang et al., "Unmasking Digital Falsehoods: A Comparative Analysis of LLM-Based Misinformation Detection Strategies," *arXiv*, vol. 2503.00724v1, 2025. [Online]. Available: <https://arxiv.org/html/2503.00724v1>
46. W. Liu et al., "Characterizing AI-Generated Misinformation on Social Media," *arXiv*, vol. 2505.10266v1, 2025. [Online]. Available: <https://arxiv.org/html/2505.10266v1>
47. From policy to practice: Responsible media AI implementation, Digital Content Next, 2025. [Online]. Available: <https://digitalcontentnext.org/blog/2025/06/30/from-policy-to-practice-responsible-media-ai-implementation/> [DOI unavailable]
48. X. Yang et al., "Understanding Human-Centred AI: a review of its defining elements and a research agenda," *Behav. Inf. Technol.*, 2025, doi: 10.1080/0144929X.2024.2448719.
49. Y. Zhao et al., "Artificial intelligence (AI) for user experience (UX) design: a systematic literature review and future research agenda," *ResearchGate*, 2023. [Online]. Available: https://www.researchgate.net/publication/373389004_Artificial_intelligence_AI_for_user_experience_UX_design_a_systematic_literature_review_and_future_research_agenda
50. Z. Khan et al., "Utilizing Generative AI for Instantaneous Content Moderation on Social Media Platforms," *ResearchGate*, 2025. [Online]. Available: https://www.researchgate.net/publication/392927657_Utilizing_Generative_AI_for_Instantaneous_Content_Moderation_on_Social_Media_Platforms
51. A. Roberts et al., "Governing artificial intelligence: ethical, legal and technical opportunities and challenges," *Philos. Trans. R. Soc. A*, vol. 376, no. 2133, p. 20180080, 2018, doi: 10.1098/rsta.2018.0080.
52. B. Patel et al., "Ethical and regulatory challenges of AI technologies in healthcare: A narrative review," *PubMed Central*, vol. 11, p. e10879008, 2025, doi: 10.3389/frai.2024.1357888.
53. Formatted_Social Media Governance project: Summary of work in 2024, OECD, 2025. [Online]. Available: <https://wp.oecd.ai/app/uploads/2025/05/social-media-governance-project-summary-of-work-in-2024.pdf> [DOI unavailable]
54. C. Wu et al., "FutureGen: LLM-RAG Approach to Generate the Future Work of Scientific Article," *arXiv*, vol. 2503.16561v1, 2025. [Online]. Available: <https://arxiv.org/html/2503.16561v1>
55. D. Singh et al., "Deepfake Attacks: Generation, Detection, Datasets, Challenges, and Research Directions," *ResearchGate*, 2023. [Online]. Available: https://www.researchgate.net/publication/374932579_Deepfake_Attacks_Generation_Detection_Datasets_Challenges_and_Research_Directions
56. E. Brown et al., "Ethical Considerations in Artificial Intelligence: A Comprehensive Discussion from the Perspective of Computer Vision," *ResearchGate*, 2023. [Online]. Available: https://www.researchgate.net/publication/376518424_Ethical_Considerations_in_Artificial_Intelligence_A_Comprehensive_Discussion_from_the_Perspective_of_Computer_Vision
57. F. Davis et al., "Reviewing the Ethical Implications of AI in Decision Making Processes," *ResearchGate*, 2025. [Online]. Available: https://www.researchgate.net/publication/378295986_REVIEWING_THE_ETHICAL_IMPLICATIONS_OF_AI_IN_DECISION_MAKING_PROCESSES
58. G. Lee et al., "Ethical and social considerations of applying artificial intelligence in healthcare—a two-pronged scoping review," *PubMed Central*, vol. 12, p. e12107984, 2025, doi: 10.1007/s00146-024-02033-6.
59. H. Wang et al., "Face Deepfakes - A Comprehensive Review," *arXiv*, vol. 2502.09812v1, 2025. [Online]. Available: <https://arxiv.org/html/2502.09812v1>
60. I. Chen et al., "Charting the Landscape of Nefarious Uses of Generative Artificial Intelligence for Online Election Interference," *arXiv*, vol. 2406.01862, 2024. [Online]. Available: <https://arxiv.org/pdf/2406.01862>

61. J. Kumar et al., "Body Perceptions and Psychological Well-Being: A Review of the Impact of Social Media and Physical Measurements on Self- Esteem and Mental Health with a Focus on Body Image Satisfaction and Its Relationship with Cultural and Gender Factors," PubMed Central, vol. 12, p. e11276240, 2025, doi: 10.3390/soc12080322.
62. K. Taylor et al., "The impact of digital technology, social media, and artificial intelligence on cognitive functions: a review," Front. Cogn., vol. 2, p. 1203077, 2023, doi:10.3389/fcogn.2023.1203077.
63. L. Nguyen et al., "How AI and Human Behaviors Shape Psychosocial Effects of Chatbot Use: A Longitudinal Randomized Controlled Study," arXiv, vol. 2503.17473v1, 2025. [Online]. Available: <https://arxiv.org/html/2503.17473v1>
64. M. Lee, "The Impact of Generative AI on Critical Thinking: Self-Reported Reductions in Cognitive Effort and Confidence Effects From a Survey of Knowledge Workers," Microsoft, 2025. [Online]. Available: https://www.microsoft.com/en-us/research/wp-content/uploads/2025/01/lee_2025_ai_critical_thinking_survey.pdf [DOI unavailable]
65. AI and Misinformation - 2024 Dean's Report, University of Florida, 2024. [Online]. Available: <https://2024.jou.ufl.edu/page/ai-and-misinformation> [DOI unavailable]
66. N. Gupta et al., "Exposing the Fake: Effective Diffusion-Generated Images Detection," arXiv, vol. 2307.06272, 2023. [Online]. Available: <https://arxiv.org/abs/2307.06272>
67. N. Gupta et al., "Exposing the Fake: Effective Diffusion-Generated Images Detection," arXiv, vol. 2307.06272v1, 2023. [Online]. Available: <https://arxiv.org/html/2307.06272v1>
68. O. Zhang et al., "Enhancing Deepfake Detection: Proactive Forensics Techniques Using Digital Watermarking," CMC-Comput. Mater. Contin., vol. 82, no. 1, p. 59264, 2025, doi: 10.32604/cmc.2024.059264.
69. P. Sharma et al., "Advancing GAN Deepfake Detection: Mixed Datasets and Comprehensive Artifact Analysis," Appl. Sci., vol. 15, no. 2, p. 923, 2025, doi: 10.3390/app15020923.
70. Q. Liu et al., "Deepfake Generation and Detection: A Benchmark and Survey," arXiv, vol. 2403.17881, 2024. [Online]. Available: <https://arxiv.org/abs/2403.17881>
71. R. Patel et al., "Ethical Challenges and Solutions of Generative AI: An Interdisciplinary Perspective," Informatics, vol. 11, no. 3, p. 58, 2024, doi: 10.3390/informatics11030058.
72. S. Kim et al., "A Comprehensive Survey with Critical Analysis for Deepfake Speech Detection," arXiv, vol. 2409.15180, 2024. [Online]. Available: <http://arxiv.org/pdf/2409.15180>
73. T. Wilson et al., "The Impact of Affect on the Perception of Fake News on Social Media: A Systematic Review," Soc. Sci., vol. 12, no. 12, p. 674, 2023, doi: 10.3390/socsci12120674.
74. U. Chen et al., "Ethical and Legal Considerations in Mitigating Disinformation," arXiv, vol. 2406.18841, 2024. [Online]. Available: <https://arxiv.org/pdf/2406.18841>
75. V. Adams et al., "'I Hadn't Thought About That': Creators of Human-like AI Weigh in on Ethics And Neurodivergence," arXiv, vol. 2506.12098, 2025. [Online]. Available: <https://arxiv.org/abs/2506.12098>
76. SPS Webinar: Recent Advances and Challenges of Deepfake Detection, Signal Processing Society, 2025. [Online]. Available: <https://signalprocessingsociety.org/blog/sps-webinar-recent-advances-and-challenges-deepfake-detection> [DOI unavailable]
77. Artificial intelligence and child sexual abuse: A rapid evidence assessment, Australian Institute of Criminology, 2025. [Online]. Available: <https://www.aic.gov.au/publications/tandi/tandi711> [DOI unavailable]
78. Addressing AI-generated child sexual exploitation and abuse, Tech Coalition, 2025. [Online]. Available: <https://technologycoalition.org/resources/addressing-ai-generated-child-sexual-exploitation-and-abuse/> [DOI unavailable]
79. W. Taylor et al., "AI-Driven Content Moderation: Ethical and Technical Challenges," arXiv, vol. 2502.07931v1, 2025. [Online]. Available: <https://arxiv.org/pdf/2502.07931>
80. X. Brown et al., "Balancing Innovation and Regulation in the Age of Generative Artificial Intelligence," J. Inf. Policy, vol. 14, p. 0012, 2024, doi: 10.5325/jinfopoli.14.2024.0012.
81. "Courtroom Scene from Sewell Setzer III Lawsuit," USA Today, 2024. [Online]. Available: <https://www.usatoday.com/story/news/nation/2024/10/23/sewell-setzer-iii/75814524007/>

82. "Dragon-Themed AI Illustration," Unsplash, 2025. [Online]. Available: <https://www.usatoday.com/gcdn/authoring/authoring-images/2024/10/23/USAT/75815793007-20240920-t-230449-z-788801772-rc-244-aae-2-vha-rtrmadp-3-auctiongameofthrones.JPG?width=1320&height=1002&fit=crop&format=pjpg&auto=webp> .
83. "Courtroom Scene from Chase Nasca Lawsuit," People.com, 2025. [Online]. Available: <https://people.com/family-sues-blaming-tiktok-for-son-suicide-being-inundated-fyp-videos-11683054>
84. "Social Media Interface," Unsplash, 2025. [Online]. Available: [https://people.com/thmb/yugXTntaClygvSz4N7mLkcNYWlw=/4000x0/filters:no_upscale\(\):max_bytes\(150000\):strip_icc\(\):focal\(999x0:1001x2\):format\(webp\)/tiktok-ban-042424-6026fa82c6db43c89f66829518b027fe.jpg](https://people.com/thmb/yugXTntaClygvSz4N7mLkcNYWlw=/4000x0/filters:no_upscale():max_bytes(150000):strip_icc():focal(999x0:1001x2):format(webp)/tiktok-ban-042424-6026fa82c6db43c89f66829518b027fe.jpg)
85. "Courtroom Illustration," Pixabay, 2025. [Online]. Available: <https://static01.nyt.com/images/2022/09/30/business/INTERNET-SUICIDE-01/INTERNET-SUICIDE-01-superJumbo.jpg?quality=75&auto=webp>
86. "Courtroom Illustration," Pixabay, 2025. [Online]. Available: https://ichef.bbci.co.uk/news/1024/cpsprodpb/183C4/production/_126786299_mollyrussell1.jpg.webp

Disclaimer/Publisher's Note: The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.