# The Impact of Deepfakes on Society And The Role of AI In Mitigating Harm

**Huzaifa Iqbal**

B.Tech (CS), First Year Student (NIET) , INDIA

## ABSTRACT:

Deepfake technology has emerged as one of the most concerning developments in the field of artificial intelligence. By manipulating images, audio, and videos, deepfakes can challenge the boundary between authentic and synthetic media . While this technology has some positive applications in entertainment, education, and accessibility, its misuse poses serious risks to society. From spreading misinformation and damaging reputations to influencing politics and increasing cybercrime, deepfakes present challenges that demand urgent attention. This paper explores both the harmful and beneficial impacts of deepfakes, highlighting the increasing concerns regarding media reliability and authenticity. It further examines how artificial intelligence itself can play a crucial role in addressing these issues, particularly through detection tools, authentication systems, and ethical frameworks. The discussion emphasizes the need for awareness, regulation, and technological solutions to ensure that innovation in AI does not come at the cost of societal harm. This paper contributes by providing a structured review of risks, applications, detection methods, and by presenting original survey results (n=200) on public awareness and trust in AI for deepfake detection.

**Key words:** Deepfakes, Artificial Intelligence (AI), Misinformation, Digital Media Authenticity, Cybercrime, Deepfake Detection, Ethical AI

---

## 1. Introduction

The rapid growth of artificial intelligence has given rise to a variety of technologies that are reshaping human communication and interaction. Among these, deepfakes represent both a technological advancement and a potential threat . A deepfake is an audio, image, or video created using AI techniques to convincingly mimic real individuals, often making it difficult to distinguish between genuine and manipulated content. While this technology has beneficial uses in areas such as film production, education, and accessibility, its misuse has far-reaching consequences. Instances of identity theft, fake news, political manipulation, and online harassment demonstrate the risks associated with deepfakes. These challenges have raised significant concern among researchers, policymakers, and the wider public.

At the same time, artificial intelligence—the very technology that enables deepfakes—also offers solutions to mitigate their harmful effects. AI-driven detection tools, watermarking methods, and authentication systems are being developed to counteract manipulation and restore trust in digital media. This paper examines the dual nature of deepfakes by discussing their social impacts and analyzing the role of AI in reducing the risks they pose.
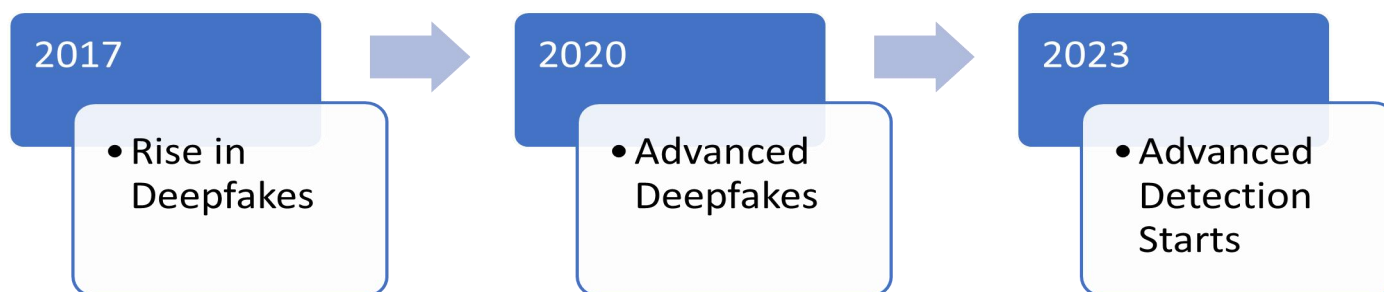


**Figure 1: Timeline tracing Advancement of Deepfake Technology and Detection**

## 2. The Dangers of Deepfake Technology.

With AI technology advancing rapidly, it is now possible to generate highly realistic yet fabricated videos and images. Although deepfakes may have entertainment value, they also pose serious risks. The following points show some of the ways deepfakes can cause harm

## 2.1 Fraud:

Deepfakes not only enable several forms of financial crimes, but also amplify them. For instance, in 2019, the CEO of a British energy company was asked by fraudsters to transfer AC 220,000 to a bank account in Hungary using audio deepfake technology to imitate the voice of the CEO of its parent company. Deepfakes can be utilized in several types of fraudulent activity, like romance fraud and grandma scams. They often serve as an entry point or launching pad for other types of crime and are frequently associated with identity theft. This type of offence increases challenges in identity verification and Management, which in turn affects society and national security more broadly. For instance, deepfakes could be used to bypass biometrics in cases of identity theft the case of fraud associated with identity theft illustrates the cascading effects of deepfake crimes. It highlights the need for more effective identity verification and management strategies, as, as well as more robust policies and practices to combat the growing threat of deepfakes.

## 2.2 Pornography and Abuse:

Most deepfake technology is used to generate non-consensual pornography, according to Sloot and Wagensveld (European Parliament).
In fact, a 2019 Deeptrace report revealed that 96% of the 14,600 online deepfake videos were used to forge non-consensual pornographic material, with a total of 134,364,438 views across the top four deepfake pornography websites. Shockingly, over 90% of deepfakes online are pornographic depictions of women.
Apart from being used for non-consensual pornography, deepfake technology poses a significant risk for victims of domestic violence as it can be used to threaten
and abuse them. EPRS (Ruff, 2022), Hayward & Maas (Ali et al., 2022), Ferreira
et al. (Ajder et al., 2019) and Lucas (Deepfake & Fight,2022) report that the most common uses of deepfakes for abuse include producing non-consensual content such as revenge, child and fantasy pornography, sextortion (cyber blackmail), and perpetrated abuse.

## 2.3 Information Manipulation and Forgery Misinformation:

Moreover, deepfakes can trigger armed conflicts that threaten national security on both sides of the Atlantic, as was the case with the purported deepfake of Gabonese President Ali Bongo Ondimba in 2019. The mere knowledge of deepfakes can undermine the credibility of any video, making it difficult to determine its authenticity, thereby leading to irreparable consequences, whether real or fake.

## 3. Positive Applications of Deepfake Technology

While deepfakes often get attention for their risks, the technology behind them also has promising and positive uses. From entertainment to education, deepfakes can open up new creative possibilities and make certain tasks easier or more engaging. Some of the key benefits include :

## 3.1 Film Industry:

Deepfake technology is significantly transforming the film industry. It can recreate the voices of actors who can no longer speak due to illness or age, preserving their performances for future audiences. Filmmakers can update or improve existing footage instead of reshooting entire scenes, which saves a lot of time and effort. Classic scenes can be brought back exactly as they were—or even improved—while actors who passed away long ago can appear in new films, letting audiences see them on screen again.

Deepfakes also make dubbing movies into different languages much easier and more realistic, allowing global audiences to access films in their preferred languages and educational content in their own language. Overall, it's a tool that's making movie-making more efficient, innovative, and widely accessible.

## 3.2 Social and Medical Benefits of Deepfakes:

Deepfake technology is starting to make a real difference in social and medical life. It can help people say goodbye to lost loved ones through virtual representations and assist Alzheimer's patients by recreating familiar, younger faces, making interactions more comforting and personal.

In medicine, deepfakes can simulate patient scenarios so doctors and nurses can practice safely, and they can create controlled virtual environments for mental health therapy, helping patients face and manage fears or trauma in a supportive way.

## 4. Types of Deepfakes:

Deepfakes are created in different ways, and each type serves a unique purpose depending on how the technology is applied. Understanding these categories helps in recognizing how deepfakes are used across various fields.

## 4.1 Photo Shopping:

Photo-editing software such as Photoshop can be regarded as a precursor to modern deepfake method, as it allows images to be altered in ways that change reality. While its edits are usually limited to still pictures, it laid the foundation for how digital tools can manipulate appearances and create convincing but misleading visuals.

## 4.2 Face Swapping:

Much of the activity in deepfake-focused communities has revolved around face-swapping technologies such as **FakeApp** and **FaceSwap**. These software tools allow users to map one actor's face onto another actor's performance, producing results that can appear highly realistic. What once started as an experiment within small online groups quickly spread, as these tools became more user-friendly and accessible to the general public. This accessibility has made face-swapping one of the most common and recognizable applications of deepfake technology.

## 4.3 Lip-Syncing and Voice Synthesis:

Another strand of deepfake production focuses not on the dynamic virtual performance of face-swapping, but on generating new speech content for recognizable figures. A well-known example is the fake video of Mark Zuckerberg, created using proprietary software. This bespoke video-faking tool combined audio

and video files from Zuckerberg's Senate testimony in April 2018 to produce a new virtual performance. The video gained widespread attention as it demonstrated how convincingly public figures could be manipulated to appear as though they expressed statements they never made. Such examples highlight both the creative possibilities and the potential dangers of voice synthesis and lip-syncing technologies.
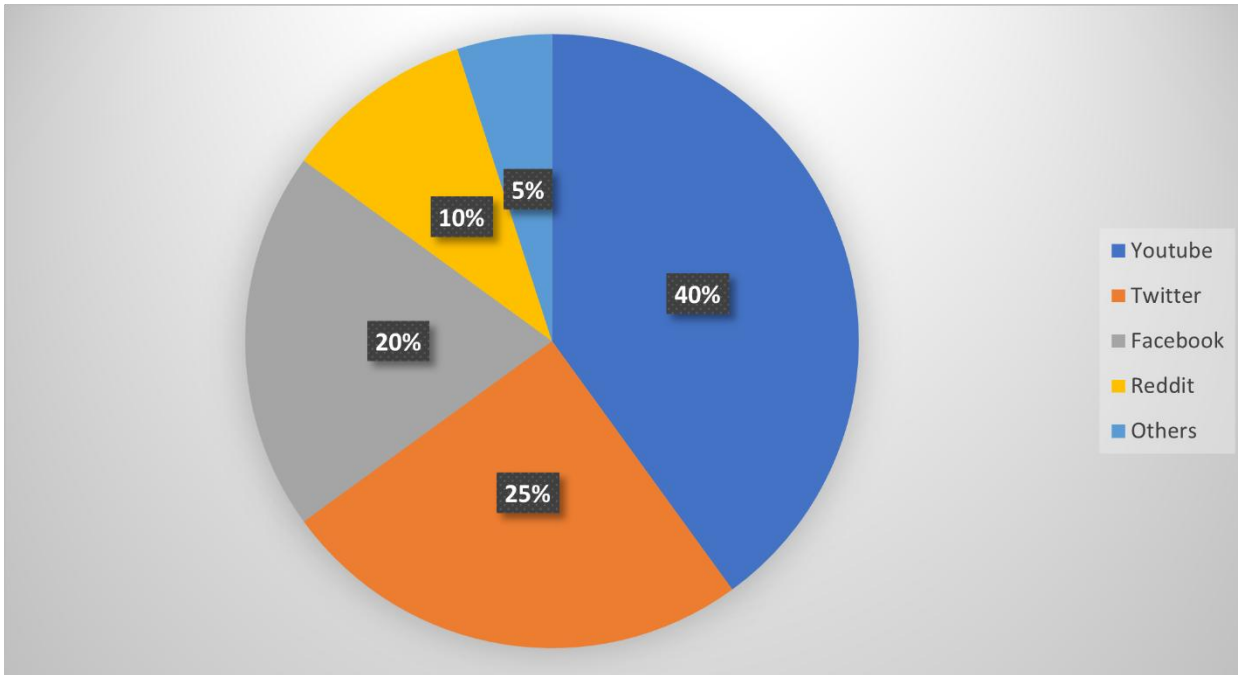


**Figure 2: Distribution of Deepfakes Across Social Media**

# 5. Preventing Misuse of Deepfakes and the Challenges Involved:

As deepfakes become more realistic and widespread, preventing their misuse has become a major concern. Harmful deepfakes can spread misinformation, damage reputations, and even threaten privacy or security. Understanding the challenges in detecting and controlling these videos is essential for developing effective prevention strategies.

## 5.1 Challenges of an Outright Ban:

Currently, no criminal or civil liability framework specifically bans the creation or distribution of deepfakes, and implementing a blanket ban is not straightforward. While deepfakes can cause significant harm, this is not true in every case. Outright prohibition could also restrict routine modifications that improve the clarity or quality of content. Moreover, deepfakes touch on issues of freedom of expression, even when they involve intentionally false statements. Since deepfakes cannot and arguably should not be banned across the board, the question remains whether their creators and distributors should be held civilly liable for any harm they cause. This makes the issue more about accountability and responsible use rather than a simple question of legality. Without clear guidelines, the boundary between creative application and harmful misuse becomes ambiguous, leaving both users and victims in a difficult position.

## 5.2 Legal Challenges in Suing Platforms Over Deepfakes:

It is challenging to achieve individual accountability for harmful deepfakes, but creators are not the only ones responsible. Online platforms already have some incentive to screen content, influenced by

moral pressure, market dynamics, and political scrutiny. However, they currently do not face any significant civil liability risk for user-generated content. This gives them the power to either address the spread of damaging deepfakes or simply ignore it, and too often commercial interests may take precedence over safety concerns. At the same time, it is also true that platforms deal with enormous volumes of content every day, making perfect monitoring extremely difficult. While a few platforms may act responsibly without legal pressure, relying solely on morals is not enough. In most cases, these forces remain insufficient to prevent real harm, which is why the debate over their role and responsibility continues.

# 6. Deepfake Detection Techniques:

Deepfake detection refers to the process of recognizing manipulated or artificially generated media content, such as videos or images, which have been generated through the utilization of deep learning methodologies. The detection of deepfakes plays a pivotal role in upholding trust in media and mitigating the dissemination of inaccurate or deceptive content. Deepfake detection involves collecting and preparing a dataset of real and potentially fake media content, ensuring format consistency, and extracting relevant features such as facial landmarks, audio spectrograms, and temporal data . The dataset is then divided into training, validation, and test sets for evaluation and training purposes. Then, feature engineering is conducted to differentiate between authentic and deepfake content.

# 6.1 Face and Body Analysis:

Face and body analysis play a crucial role in the detection of deepfakes, which are videos or images created using artificial intelligence techniques that have been altered or fabricated. Face and body analysis can help identify anomalies and inconsistencies that may indicate the presence of a deepfake. Here are some essential facial and body analysis techniques and considerations for deepfake detection:

• **Facial Landmark Detection**: This technique identifies and tracks specific points on a person's face, such as the eyes, nose, mouth, and other facial features . Deepfake detectors analyze the alignment and mobility of these landmarks over time using this information. These landmarks may not move organically or consistently in deepfake videos, indicating that the video has been manipulated.

• **Blink Analysis:** Blink analysis focuses on identifying blinking patterns in videos that are not natural. Deepfake detectors investigate how often and when a person blinks in a video. Blink patterns that exhibit anomalies or inconsistencies may indicate that a video is a deepfake.

• **Lip Synchronization Detection:** Lip synchronization analysis determines whether the audio and facial movements in a video are in sync. In deepfake videos, the vocal movements may not precisely correspond to the spoken words, which can be an indication of manipulation. Deepfake detectors examine this synchronization in order to identify potential inconsistencies.

# 6.2 Generative model analysis:

Generative models, particularly Generative Adversarial Networks (GANs), have played a major role in the creation of deepfake content. However, they can also be used to detect deepfakes, albeit in a different manner. Here is how generative models can be utilized to detect deepfakes:

• **Model Artifacts:** Deepfake generation models, such as Generative Adversarial Networks (GANs), frequently introduce particular artifacts into the generated content. These artifacts are unusual patterns or distortions that do not exist in actual images or videos. The detection of these artifacts may be indicative of a deepfake.

• **Detection of GAN Noise:** GANs, which are frequently used to generate deepfake images, incorporate noise patterns into deepfake images. Analyzing these noise patterns can be a useful method for detecting GAN generated content.

# 7. Survey on Awareness and Perception of AI in Deepfake Detection:

## 7.1 Introduction:

To explore how people perceive the role of AI in addressing the problem of deepfakes, a survey was conducted using Google Forms. This method made it possible to reach a diverse group of participants efficiently while maintaining respondent anonymity . Using both structured questions and open-ended ones allowed the study to capture measurable data as well as personal insights, providing a comprehensive overview of public opinion.

## 7.2 Methods:

To understand people's views in a clear and organized way, this study used a survey as the main method of data collection. The following parts explain how the survey was designed, who took part in it, and how the responses were gathered and analyzed.

### • Survey Design:

A structured online survey was conducted on **September 19, 2025** to collect primary data on public perception and awareness of AI in combating deepfakes. The survey included students and professionals from diverse backgrounds and consisted of over 10+ questions, incorporating multiple-choice, Likert scale, and open-ended formats, allowing participants to express their opinions in various ways. The aim was to understand not only how familiar people are with deepfakes and AI, but also how much they trust AI to detect manipulated content and their overall perspective on its effectiveness. The survey also explored participants' concerns about the misuse of deepfakes and their expectations for AI's role in mitigating such risks. This helped provide a more nuanced understanding of public attitudes toward emerging detection technologies.

### • Participants:

A total of over 200 respondents participated in the survey. Participants were selected using a combination of conventional sampling, random sampling, and targeted outreach, ensuring a diverse mix of ages, educational backgrounds, and occupations. This approach allowed the survey to capture a wide range of perspectives on public awareness of deepfakes and trust in AI for detecting manipulated content.

### •Procedure:

Responses were collected via Google Forms from participants who voluntarily took part in the survey and remained anonymous throughout. Quantitative data were analyzed using percentages and mean scores, while open-ended responses were carefully examined and categorized into recurring themes. Graphs and charts were then generated to visualize the results, providing a clear and comprehensive view of public awareness, perceptions, and trust in AI for detecting deepfakes.

## 7.3 Results:

The survey included over 200+ participants from diverse backgrounds, providing both quantitative and qualitative insights into public awareness and trust in AI for detecting deepfakes. The following statistics and

visualizations summarize key findings from the responses, highlighting trends in familiarity, trust, and perceptions of AI effectiveness.

## • Awareness and Understanding of Deepfakes:

The survey revealed that 73% of respondents had heard about deepfakes. Out of this group, 62% were able to correctly or partially explain what deepfakes are, while the remaining participants admitted they had only heard the term without fully understanding it. Students represented the largest group of respondents demonstrating familiarity with deepfakes, with the majority coming from the 19–29 age group. Interestingly, a significant number of respondents from the 40+ age group also reported awareness, showing that knowledge of deepfakes is not limited to younger generations.

## • Trust in AI Detection Tools:

When asked about their willingness to rely on AI tools to detect deepfakes, 61% of respondents expressed trust in such systems. While this shows a promising level of acceptance, many participants also voiced important concerns. Respondents frequently pointed to three issues: a lack of strong action from authorities in regulating or responding to deepfakes, limitations in the current capability of AI systems, and doubts over whether AI itself could be manipulated or misused. These perspectives highlight that, although people are open to AI as a solution, confidence in its reliability is still developing.

## • Challenges in Combating Deepfakes:

The survey further explored what participants saw as the biggest obstacles in fighting deepfakes with AI. The most common responses included the lack of motivation from authorities, distrust in the technology, and the perception that AI is not yet advanced enough to deal with rapidly evolving deepfake techniques. Together, these concerns suggest that both technological and institutional improvements are needed before AI can become a fully trusted safeguard.

## • Experience and Exposure:

In terms of direct exposure, 36% of respondents reported having encountered a deepfake in real life. Among them, 78% stated that they were able to identify the content as fake right away, 9% said they recognized it only after closer inspection, while the rest admitted that they could not identify it until someone else pointed it out. This finding indicate that while some individuals are developing the ability to recognize deepfakes, a large portion of the public remains vulnerable to being misled. Most respondents agreed that AI currently seems most capable of detecting image- and video-based deepfakes, with less confidence expressed about its ability to detect audio or text-based manipulations such as fake news or voice clones.

## • Responsibility for Preventing Deepfakes:

When asked who should be primarily responsible for preventing the spread of deepfakes, responses showed that participants overwhelmingly believe accountability should be shared. A strong majority (94%) of respondents selected the option 'all of the above' emphasizing the need for a collective effort. Breaking it down further, 87% placed responsibility on the creators of deepfakes, 68% pointed to social media platforms as key actors in controlling their spread, and 57% felt government and regulators must play a stronger role. These results highlight a broad recognition that no single group can address the issue alone, and that effective prevention will require cooperation between technology developers, platforms, and policymakers.
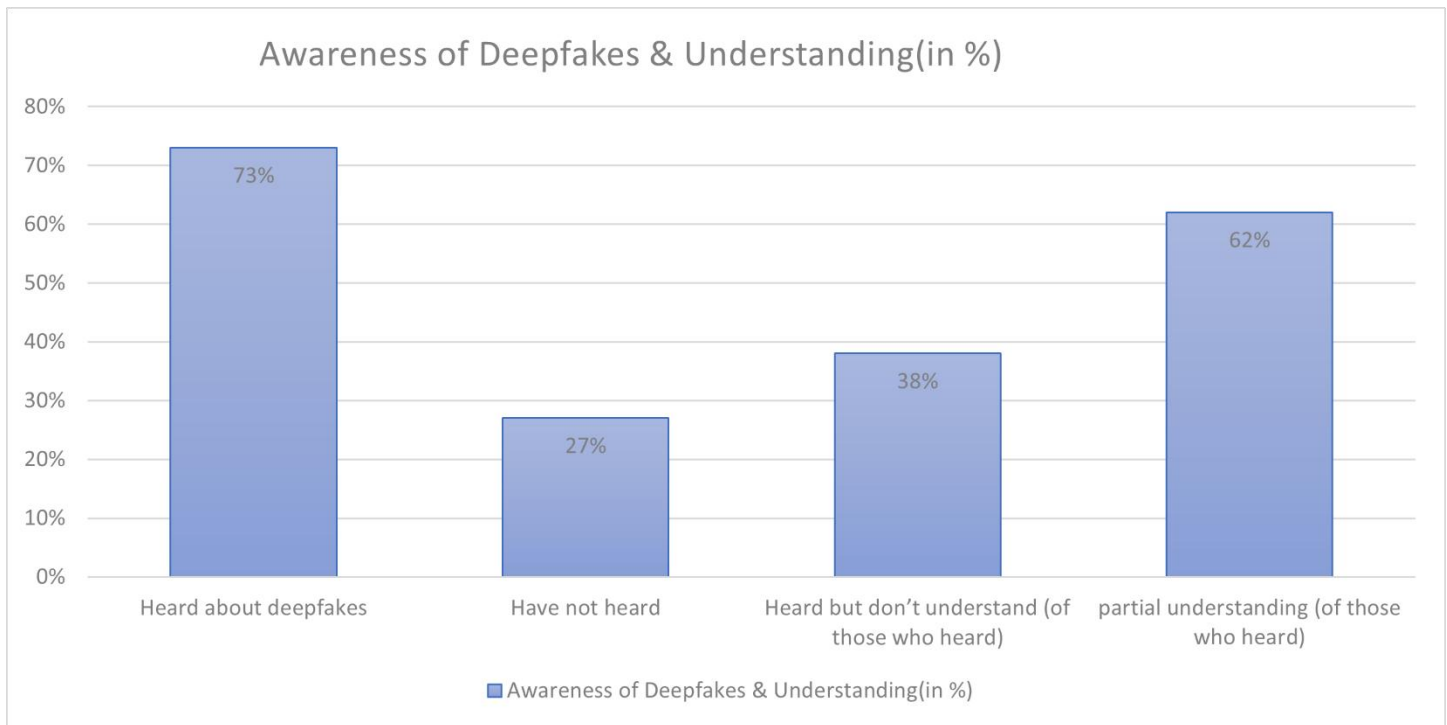
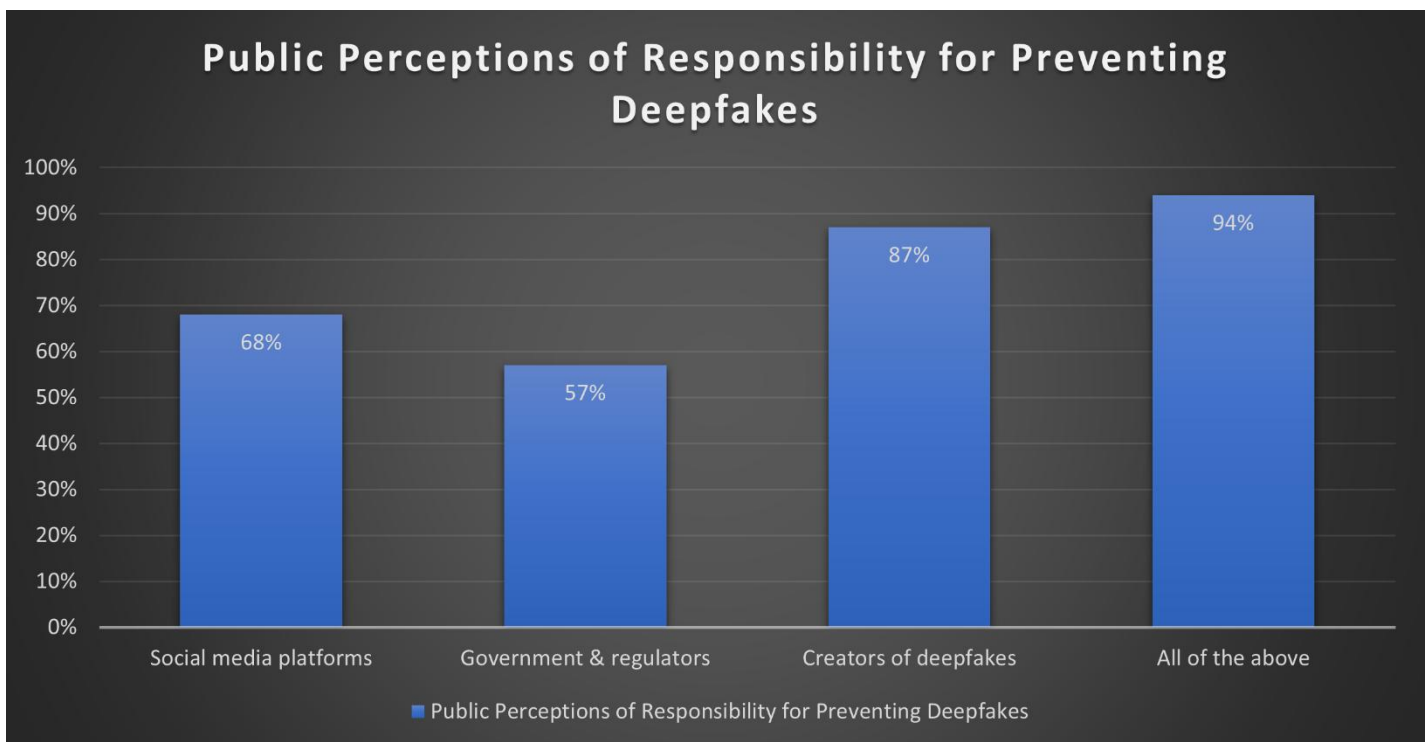**Figure 3: Awareness & Understanding of Deepfakes Among Respondents**



**Figure 4: Public Perceptions of Responsibility in Preventing Deepfakes**

# 8. Conclusion:

Deepfakes are an innovative yet challenging advancement in artificial intelligence. While they open doors to creativity in film, education, and medicine, their misuse—through fraud, abuse, and misinformation—poses serious risks to individuals and society. This research highlights that, although many people are aware of deepfakes, understanding and trust in AI detection tools vary, with concerns about reliability and manipulation still present.

The findings also show that stopping harmful deepfakes is not the responsibility of a single group. Creators, social media platforms, and governments all need to work together, alongside informed and vigilant users. With ethical AI, ongoing technological improvements, and public awareness, society can balance the benefits of deepfakes with the need to minimize harm, making sure this powerful technology is used responsibly.

Future research should focus on improving AI-based detection of audio deepfakes and on balanced innovation with safeguards against misuse.

# 9. References:

[1]. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3213954
[2]. https://peerj.com/articles/cs-2037/
[3]. https://www.sciencedirect.com/science/article/pii/S240584402500653X
[4]. https://datasociety.net/library/deepfakes-and-cheap-fakes/
[5]. https://www.tandfonline.com/doi/full/10.1080/23742917.2023.2192888#d1e159
[6]. https://ijarsct.co.in/Paper15308.pdf
[7]. https://www.researchgate.net/publication/351300442_Deep_Insights_of_Deepfake_Technology_A_Review
[8]. https://link.springer.com/article/10.1007/s42454-024-00054-8
[9]. https://arxiv.org/pdf/1909.11573