# Deepfake videos and images detection

**Pratik Prakash Kanojiya**

**Ritick Ramsagar Rai**

**Guide: Asst. Prof. Gauri Ansurkar**

Keraleeya Samajam's Model College, Khambalpada Road, Thakurli,
Dombivli(East),Maharashtra.

## Abstract

With the rapid proliferation of deepfake technology, the ability to generate highly convincing fake media content has raised significant concerns regarding misinformation and potential threats to individual privacy. This research aims to explore the current landscape of deepfake detection methodologies, focusing on advancements, limitations, and emerging challenges.

The study begins by providing a comprehensive overview of deepfake generation techniques, emphasizing the complexity of the problem and the need for advanced detection mechanisms. Subsequently, it reviews existing detection approaches, including traditional methods and state-of-the-art deep learning-based models. The strengths and weaknesses of these techniques are critically evaluated, shedding light on the evolving nature of deepfake attacks and the corresponding arms race in detection strategies.

In addressing the challenges associated with the dynamic nature of deepfake technology, the research investigates the feasibility of cross-modal detection, combining insights from multiple modalities such as audio, video, and textual content. This holistic approach aims to improve detection accuracy and resilience against adversarial attacks.

The study also considers the ethical implications of deepfake detection, emphasizing the importance of balancing security measures with the preservation of user privacy. It explores potential regulatory frameworks and ethical guidelines to govern the deployment of deepfake detection systems in various domains.

**Keyword**:  Deepfake creation, Deepfake Detetcion, Faceswap,
Generative Adversarial Networks

## Introduction

In today's digital age, where technology continues to advance at a rapid pace, a new form of digital deception has emerged, known as deepfakes. Deepfakes are sophisticated manipulations of audio and video content that use artificial intelligence to replace or superimpose faces and voices onto existing footage, creating convincing but entirely fabricated content.

As these deepfakes become more prevalent, the need for reliable detection methods becomes increasingly crucial. Imagine a world where anyone could be portrayed saying or doing anything without their knowledge or consent. The potential consequences are vast, ranging from misinformation and fake news to more insidious uses.

This research aims to explore and develop a simple yet effective approach to detecting deepfakes. By understanding the basic principles behind deepfake creation and leveraging accessible techniques,This tendency has made it easier and more accurate for non-technical people to edit films and manipulate data by changing expressions, introducing dialogue, and swapping out faces. Such manipulations have a significant impact on political campaigns, public opinion, and people's reputations, but they can be

very hard to detect. To generate realistic pictures and videos, deep learning models frequently require large amounts of photo and video data. Public personalities, like politicians and celebrities, are often targeted because of their vast libraries of online media content. The ability to edit these photos and videos allows for the creation of deepfakes that change these people's appearances or place them in fake situations. This technology's impacts extend beyond entertainment and may provide serious risks.

## Background on Deepfakes

### Where did deepfakes originate?

Naturally, the modification of photos and video has existed for nearly as long as the history of photography and film. Throughout the 20th century, advances in technology and methods for manipulating images and sounds were made; in the early 1990s, academic institutions began to create deepfake technology. A further development of deepfake technology was carried out independently by amateurs; more recently, deepfakes have also been employed by businesses and other organizations. On Reddit, the word "deepfake" first appeared. On the subreddit deepfakes, users posted their own deepfake videos. Most of these videos featured adult actresses in explicit clips with the faces of celebrities switched out. However, other online groups continue to distribute explicit material on platforms that have not prohibited deepfake explicit material. On February 7, 2018, the subreddit deepfakes was closed for "involuntary explicit material" due to the nonconsensual nature of the media. Communities that share deepfakes on Reddit still exist, but they show scenes that aren't obscene.
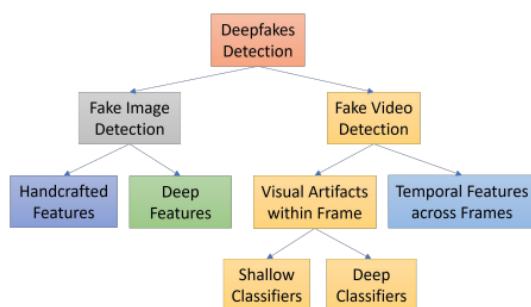


Fig 1:Types of Deepfakes

## Significance of Deepfake Detection

The significance of deepfake detection cannot be overstated in the contemporary digital landscape. With the rapid advancement of artificial intelligence and deep learning technologies, the creation and dissemination of highly convincing fake multimedia content, known as deepfakes, have become prevalent. These manipulated videos, images, and audio recordings pose serious threats to individuals, businesses, and society at large.

One primary concern is the potential for deepfakes to be used in malicious activities, such as spreading misinformation, defaming public figures, or manipulating public opinion. The ability to convincingly fabricate content that appears genuine makes it challenging for individuals to discern between real and manipulated media. This not only jeopardizes the trustworthiness of information but also has the potential to undermine the foundations of democratic societies where informed decision-making relies on accurate and truthful data.In addition to the societal impact, there are significant personal and privacy implications. Deepfakes can be leveraged for identity theft, cyberbullying, or other forms of online harassment. Individuals may find themselves targeted by false narratives created through manipulated content, leading to reputational damage and emotional distress.Moreover, the rise of deepfakes poses challenges to the authenticity of digital evidence in various sectors, including law enforcement and the judiciary. The potential for fabricated videos or audio recordings to be presented as genuine evidence in legal proceedings raises questions about the reliability of digital media in the modern era.The development of robust deepfake detection methods is crucial to address these multifaceted challenges. It not only safeguards the integrity of information but also protects individuals from the harmful consequences of malicious deepfake use. By advancing our capabilities to identify and mitigate the impact of deepfakes, we contribute to the preservation of trust, privacy, and the responsible use of digital media in our increasingly interconnected world. The significance of deepfake detection extends beyond technological advancements; it is a cornerstone in upholding the ethical and moral standards essential for a secure and trustworthy digital environment.

## Purpose and Scope of the Research

The purpose of this research is to investigate and contribute to the evolving field of deepfake detection, aiming to develop robust methods capable of identifying manipulated multimedia content with high accuracy. In the digital age, where the creation

and dissemination of deepfakes have become increasingly sophisticated, there is a pressing need for advanced detection techniques to mitigate the potential societal, political, and personal consequences of deceptive media.The primary objective is to enhance our understanding of the challenges posed by deepfakes and to propose a novel detection framework that surpasses current limitations. By delving into the intricacies of deepfake technology, the research seeks to identify key features and patterns that distinguish manipulated content from authentic media. This involves a comprehensive exploration of existing literature, methodologies, and datasets associated with deepfake detection.The scope of this research extends to various domains, including but not limited to cybersecurity, media forensics, and digital ethics. The proposed detection framework aims to be versatile and applicable across different types of multimedia content, such as videos, images, and audio recordings. Additionally, the research explores the potential integration of machine learning models, feature extraction techniques, and other innovative approaches to bolster the effectiveness of the detection system.Furthermore, the research has broader implications for the development of countermeasures against the malicious use of deepfakes. As deepfake technology evolves, so must our detection capabilities, and this study strives to contribute meaningful insights that can be translated into practical applications for safeguarding individuals, institutions, and society at large.In summary, the purpose of this research is to advance the state of deepfake detection through a comprehensive examination of existing methodologies and the development of an innovative framework. The scope encompasses a multidisciplinary approach, aiming to address the complex challenges posed by deepfakes and contribute to the ongoing efforts to maintain the integrity and trustworthiness of digital media in the contemporary landscape.

## Face swap

Face swapping is one of the most common subcategories of face modification, which has grown in popularity in the modern day. In order to create a new image that combines the physical traits of both people, the facial features of one person are digitally replaced with those of another. The original attempt at face switching was called FakeApp, and it was made by a Reddit user. During the Deepfake development stage, an autoencoder-decoder pairing

structure is frequently used. After the autoencoder has first extracted the latent features from the facial images, the decoder reconstructs and decodes the images. To allow for facial substitution between the source and destination images, a dual encoder-decoder architecture is needed. The two pairs share the encoder network settings, and each pair trains on a different set of images. Because of the common encoder network, both couples are able to consistently extract facial traits.

The shared encoder is able to recognize and understand the similarities between two different sets of facial photos by using this method. Because the placement of the mouth, nose, and eyes are examples of relative traits seen in facial structures, this exercise is not too difficult. demonstrates a methodical process for producing Deepfakes, which involve combining the feature set of a source face A with decoder B to create an altered version of face B that bears similarities to the original face A. This technique is used in a number of projects, including DeepFaceLab and Deep-Fake tensorflow.
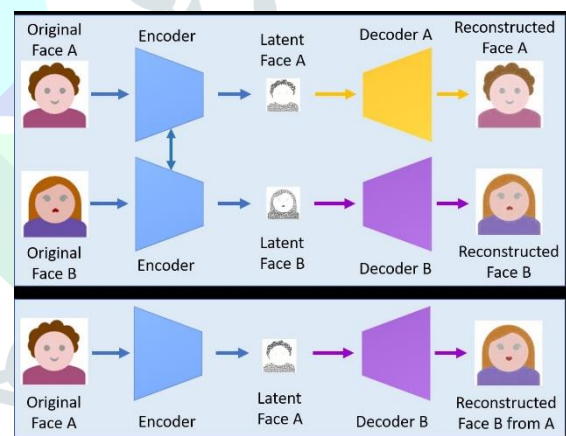


Fig 2: A deepfake creation model using two encoder-decoder pairs

## Generative Adversarial Networks (GAN)

The dual neural network architecture known as a GAN (Generative Adversarial Network) consists of a generator and a discriminator. While the discriminator distinguishes between real and fake content, the generator creates fake content by employing a random vector. GANs can generate nonexistent real faces and produce realistic deepfakes. STYLEGAN and VGGFace are popular GAN-based methods for producing deepfakes, including the "this person does not exist" website.The adversarial loss and perceptual loss layers are two extra layers included in the architecture of deepfake techniques such as GANs.

These layers improved the quality and realism of created synthetic images by employing an autoencoder-decoder technique to capture latent facial traits like eye movements. CycleGAN is a deepfake technique that applies unique characteristics from one picture to another by utilizing the GAN architecture. To find latent characteristics more quickly, this approach uses a cycle loss function. Unlike supervised approaches, the unsupervised CycleGAN method may perform image-to-image translation without requiring matched samples. Put differently, the model may learn attributes of several pictures from both the source and target domains, even if they have no relation to each other.
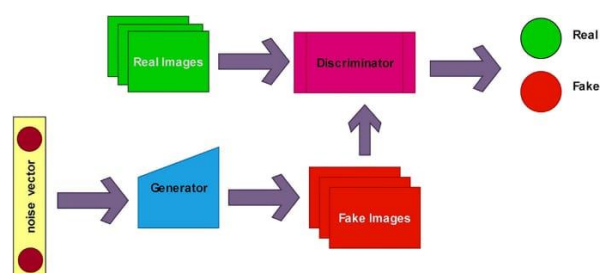


Fig 3: The GAN architecture consisting of a generator and a discriminator, and each can be implemented by a neural network.

Deepfake detection poses a critical challenge, often treated as a binary classification problem distinguishing authentic from tampered videos. The efficacy of such methods hinges on extensive datasets for training classifiers. While the number of fake videos is growing, creating a benchmark for validation remains limited. To address this gap, Korshunov and Marcel devised a substantial deepfake dataset using a GAN model, generating 620 videos via Faceswap-GAN from the VidTIMIT database. Mimicking facial expressions and movements, these videos serve as a valuable testing ground for various detection methods.

Results indicate that popular face recognition systems, like VGG and Facenet, struggle to detect deepfakes effectively. Lip-syncing approaches and image quality metrics with SVM also exhibit high error rates. This underscores the urgent need for robust methods capable of discerning deepfakes from genuine content.Surveying deepfake detection methods reveals two major categories: fake image detection and fake video detection. The latter is subdivided into single-frame-based methods and temporal features across frames-based methods. Temporal feature methods typically employ deep learning recurrent classification models, while visual artifacts within video frames can be addressed by both deep and shallow classifiers.

## Methodology

### How deepfake generation works steps:

Early detection methods relied on handcrafted features, while recent approaches, such as those leveraging deep learning, automatically extract salient and discriminative features to identify deepfakes. As deepfakes continue to pose threats to privacy, societal security, and democracy, ongoing research aims to develop more robust detection methods.Step of deepfake video & image preperation

### 1. Data Collection:
Collect an extensive dataset comprising images or videos featuring the target person, ensuring diversity and high quality for a more authentic deepfake.

### 2. Preprocessing:
Clean and preprocess the gathered data, maintaining consistency in lighting conditions, angles, and facial expressions.
Utilize face detection algorithms to identify and extract facial features from images or frames.

### 3. Feature Extraction:
Employ deep neural networks, such as Convolutional Neural Networks (CNNs), to extract intricate facial features and patterns from the preprocessed data.
Capture representations of facial expressions, head movements, and other nuanced details crucial for a realistic deepfake.

### 4. Generative Models:
Utilize generative models, with Generative Adversarial Networks (GANs) being a prominent choice. GANs consist of a generator responsible for creating synthetic content and a discriminator tasked with evaluating whether the content is real or fake.

### 5. Training the Model:
Train the generative model using the pre-processed dataset. The generator refines its ability to create realistic content by learning from actual examples in the dataset.
Simultaneously, the discriminator hones its capacity to distinguish between real and synthetic content.

### 6. Feedback Loop:
Establish a feedback loop during training where the generator endeavors to produce content indistinguishable from real data, while the discriminator continually improves its discrimination capabilities.

## 7. Fine-Tuning:

Fine-tune the model based on feedback and evaluation, adjusting parameters to enhance the quality and realism of the generated content.
Iterate this process to achieve a generator capable of producing visually convincing deepfakes.

## 8. Post-Processing:

Apply post-processing techniques to refine the visual quality of the generated content, involving adjustments to colors, smoothing transitions, and enhancing details.

## 9. Deployment:

Once the generative model is trained and fine-tuned, it can be utilized to generate deepfake videos or images by providing input data, such as a source video featuring a different individual.

It is important to acknowledge that while deepfake technology holds legitimate and creative potential, ethical concerns arise, especially in the context of misinformation and privacy. Advances in deepfake detection techniques are crucial to mitigate potential misuse of this technology.
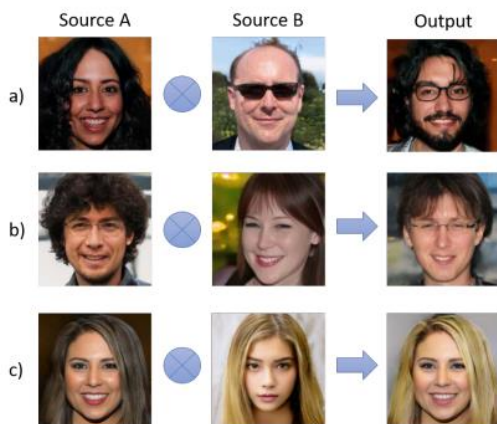


Fig 4: Examples of mixing styles using StyleGAN

## Deep Features-based Methods

Deep features-based methods represent a class of techniques in machine learning and computer vision that exploit features extracted by deep neural networks for diverse applications. These methods capitalize on the hierarchical representations learned by deep models like Convolutional Neural Networks (CNNs) to capture intricate features from raw input data.

### Feature Extraction:

In these methods, the initial layers of pre-trained deep neural networks serve as feature extractors. These layers capture hierarchical features, such as textures and shapes, progressively learning complex patterns. This is particularly relevant for tasks related to images.

### Transfer Learning:

One notable advantage of deep features-based methods lies in transfer learning. Models pre-trained on extensive datasets, like ImageNet, acquire generic features applicable to a broad spectrum of tasks. This enables the utilization of learned features for new tasks with limited labeled data, enhancing both performance and efficiency.

### Classification and Recognition:

Deep features are commonly applied in tasks like image classification and object recognition. The extracted representations encode discriminative information, facilitating accurate identification of objects or scenes. This proves valuable where traditional handcrafted features face challenges.

### Semantic Segmentation:

In semantic segmentation, deep features play a crucial role in understanding the context of individual pixels in an image. The learned representations aid in delineating object boundaries and assigning labels to regions, contributing to a comprehensive understanding of a scene.

### Anomaly Detection:

These methods find applications in anomaly detection by training on normal data. The model learns to represent typical patterns, and deviations from these patterns can indicate anomalies or outliers in the data.

### Image Retrieval:

Deep features are instrumental in content-based image retrieval. By mapping images into a shared feature space, the similarity between images can be measured, enabling efficient retrieval of semantically similar images.

### Limitations and Considerations:

Despite their effectiveness, deep features-based methods pose challenges. Ensuring interpretability of learned features, addressing biases in pre-trained models, and managing domain-specific intricacies are ongoing considerations in their application.

In summary, deep features-based methods showcase the adaptability and efficacy of leveraging pre-trained deep neural networks for diverse tasks. This underscores the transformative impact of deep learning, particularly in the realm of computer vision, while emphasizing the need for ongoing scrutiny of potential challenges and biases.

## Detection Models

**Review of Machine Learning Models:**
Traditional ML models, such as Support Vector Machines (SVMs), K-Nearest Neighbors (KNNs), and Random Forests, were initially employed for deepfake detection. These models typically rely on handcrafted features extracted from images or videos, such as texture, color, and facial landmarks.

However, handcrafted features may not capture the complex patterns inherent in deepfakes, leading to limited performance. Additionally, traditional ML models can be susceptible to variations in data and require extensive feature engineering, making them less generalizable and adaptable.

**Deep Learning Models**
Deep learning models have surpassed traditional ML models in recent years due to their ability to automatically learn complex features from data. Several deep learning architectures have been successfully applied to deepfake detection, including:

**Convolutional Neural Networks (CNNs):** CNNs excel at capturing spatial features in images and videos, making them well-suited for detecting deepfakes. They have achieved state-of-the-art performance in various deepfake detection tasks.
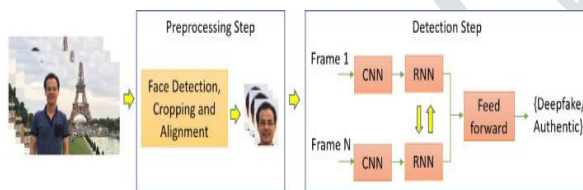


Fig 5: A two-step process for face manipulation detection where the preprocessing step aims to detect, crop and align faces on a sequence of frames and the second step distinguishes manipulated and authentic face images by combining convolutional neural network (CNN) and recurrent neural network (RNN)

**Recurrent Neural Networks (RNNs) and Long Short-Term Memory (LSTM) networks:**
RNNs and LSTMs are adept at learning temporal dependencies in sequences, making them suitable for analyzing video frames and detecting inconsistencies that may indicate deepfakes.
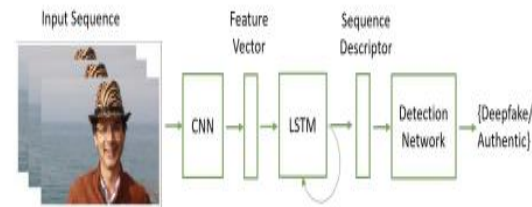


Fig 6: A deepfake detection method using convolutional neural network (CNN) and long short term memory (LSTM) to extract temporal features of a given video sequence, which are represented via the sequence descriptor}

**Generative Adversarial Networks (GANs):**
GANs can be trained to distinguish between real and fake data, making them effective for deepfake detection. They can learn the subtle differences between real and fake data, even when the deepfakes are highly realistic.

**Autoencoders:**
Autoencoders are trained to reconstruct the input data, and any reconstruction errors can be indicative of manipulated content. They can be effective for detecting deepfakes, especially when combined with other models.
Image of Convolutional Neural Network (CNN) architecture Convolutional Neural Network (CNN) architecture
These deep learning models have demonstrated significant improvements over traditional ML models. They can achieve higher accuracy in detecting deepfakes, even when dealing with subtle manipulations and variations in data.

**Hybrid Models:**
Recent research has explored hybrid models that combine traditional ML and deep learning techniques. These models leverage the strengths of both approaches:

Traditional ML models can be used for pre-processing the data and extracting initial features.
Deep learning models can then be used to learn more complex features and perform the final classification. Hybrid models have shown promise in achieving even better performance than purely deep learning models, especially for specific deepfake detection tasks.
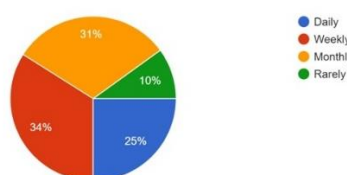
## Public Survey
We first conducted a poll of people through Google form creator and data collection service to acquire information regarding people's awareness
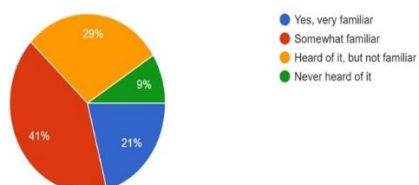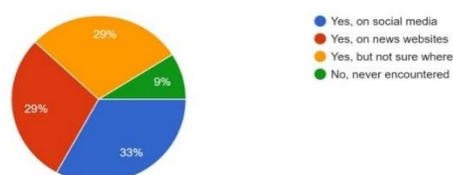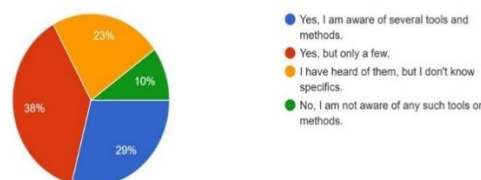
## Questionnaire

- How often do you use social media platforms?
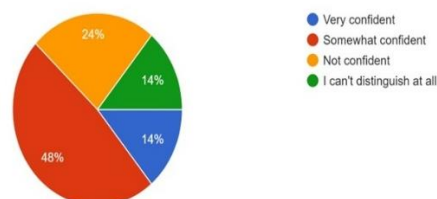- Are you familiar with the concept of deepfake technology?
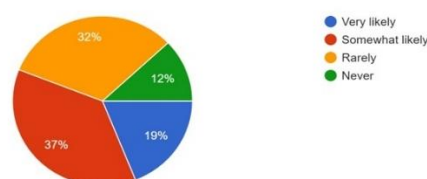- Have you ever encountered a deepfake video? If yes, where?
- Are you aware of any tools or methods to detect deepfake videos?
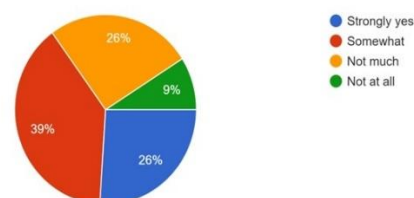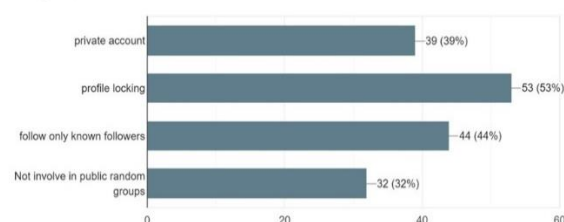- How confident are you in distinguishing real videos from deepfakes on social media?
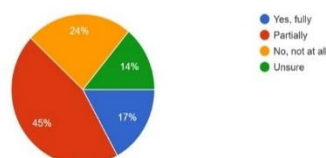- How likely are you to believe video content shared on social media is authentic?
- Do you think deepfakes on social media can influence political opinions?
- which security practice you do for social media choose which you used?
- Should social media companies be held responsible for deepfake content spread on their platforms?
- What measures do you think can be effective in combating the spread of deepfakes?

**Are you aware of any tools or methods to detect deepfake videos?**
100 responses

- Yes, I am aware of several tools and methods.
- Yes, but only a few.
- I have heard of them, but I don't know specifics.
- No, I am not aware of any such tools or methods.

23%, 10%, 38%, 29%

**How confident are you in distinguishing real videos from deepfakes on social media?**
100 responses

- Very confident
- Somewhat confident
- Not confident
- I can't distinguish at all

24%, 14%, 48%, 14%

**How likely are you to believe video content shared on social media is authentic?**
100 responses

- Very likely
- Somewhat likely
- Rarely
- Never

32%, 12%, 37%, 19%

**Do you think deepfakes on social media can influence political opinions?**
100 responses

- Strongly yes
- Somewhat
- Not much
- Not at all

26%, 9%, 39%, 26%

## Results

**How often do you use social media platforms?**
100 responses

- Daily
- Weekly
- Monthly
- Rarely

31%, 10%, 34%, 25%

**Are you familiar with the concept of deepfake technology?**
100 responses

- Yes, very familiar
- Somewhat familiar
- Heard of it, but not familiar
- Never heard of it

29%, 9%, 41%, 21%

**Have you ever encountered a deepfake video? If yes, where?**
100 responses

- Yes, on social media
- Yes, on news websites
- Yes, but not sure where
- No, never encountered

29%, 9%, 29%, 33%

**which security practice you do for social media choose which you used?**
100 responses

| Practice | Value |
| --- | --- |
| private account | 39 (39%) |
| profile locking | 53 (53%) |
| follow only known followers | 44 (44%) |
| Not involve in public random groups | 32 (32%) |

**Should social media companies be held responsible for deepfake content spread on their platforms?**
100 responses

- Yes, fully
- Partially
- No, not at all
- Unsure

24%, 14%, 45%, 17%

**What measures do you think can be effective in combating the spread of deepfakes?**
100 responses

| Measure | Value |
| --- | --- |
| Stricter laws and regulations | 31 (31%) |
| Better technology for detection | 46 (46%) |
| Public awareness campaigns | 42 (42%) |
| All of the above | 39 (39%) |

## Descriptive Analysis

Descriptive statistics is a means of describing features of a data set by generating summaries about data samples

| Are you aware of any tools or methods to detect deepfake videos? | |
|---|---|
| Mean | 2.14 |
| Standard Error | 0.095367445 |
| Median | 2 |
| Mode | 2 |
| Standard Deviation | 0.953674446 |
| Sample Variance | 0.909494949 |
| Kurtosis | -0.741305354 |
| Skewness | 0.426716252 |
| Range | 3 |
| Minimum | 1 |
| Maximum | 4 |
| Sum | 214 |
| Count | 100 |
| Largest(1) | 4 |
| Smallest(1) | 1 |
| Confidence Level(95.0%) | 0.1892297 |

| How likely are you to believe video content shared on social media is authentic? | |
|---|---|
| Mean | 2.37 |
| Standard Error | 0.092828722 |
| Median | 2 |
| Mode | 2 |
| Standard Deviation | 0.928287225 |
| Sample Variance | 0.861717172 |
| Kurtosis | -0.81853751 |
| Skewness | 0.119067101 |
| Range | 3 |
| Minimum | 1 |
| Maximum | 4 |
| Sum | 237 |
| Count | 100 |
| Largest(1) | 4 |
| Smallest(1) | 1 |
| Confidence Level(95.0%) | 0.184192325 |

| Have you ever encountered a deepfake video? If yes, where? | |
|---|---|
| Mean | 2.14 |
| Standard Error | 0.098494 |
| Median | 2 |
| Mode | 1 |
| Standard Deviation | 0.984937 |
| Sample Variance | 0.970101 |
| Kurtosis | -1.04804 |
| Skewness | 0.295452 |
| Range | 3 |
| Minimum | 1 |
| Maximum | 4 |
| Sum | 214 |

| | |
|---|---|
| Count | 100 |
| Largest(1) | 4 |
| Smallest(1) | 1 |
| Confidence Level(95.0%) | 0.195433 |

| Should social media companies be held responsible for deepfake content spread on their platforms? | |
|---|---|
| Mean | 2.35 |
| Standard Error | 0.092523543 |
| Median | 2 |
| Mode | 2 |
| Standard Deviation | 0.925235433 |
| Sample Variance | 0.856060606 |
| Kurtosis | -0.673901659 |
| Skewness | 0.334762918 |
| Range | 3 |
| Minimum | 1 |
| Maximum | 4 |
| Sum | 235 |
| Count | 100 |
| Largest(1) | 4 |
| Smallest(1) | 1 |
| Confidence Level(95.0%) | 0.183586783 |

| How confident are you in distinguishing real videos from deepfakes on social media? | |
|---|---|
| Mean | 2.38 |
| Standard Error | 0.089646 |
| Median | 2 |
| Mode | 2 |
| Standard Deviation | 0.896458 |
| Sample Variance | 0.803636 |
| Kurtosis | -0.57144 |
| Skewness | 0.369327 |
| Range | 3 |
| Minimum | 1 |
| Maximum | 4 |
| Sum | 238 |
| Count | 100 |
| Largest(1) | 4 |
| Smallest(1) | 1 |
| Confidence Level(95.0%) | 0.177877 |

## Findings:

- Most of the people not aware of the deepfake detection tools
- Somewhat likely believe media on social media is authentic
- Strongly deepfake on social media can influence political opinions
- Somewhat social media company also responsible for not detecting real content
- Better technology for detection can combating the spread of technology

## Challenges and Future Directions

Despite the significant progress made in ML-based deepfake detection, several challenges remain:
Evolving deepfake techniques: Deepfake creators are constantly improving their methods, making it difficult for detection models to keep pace.
Limited training data: Large amounts of high-quality training data are required for deep learning models to achieve robust performance.
Computational cost: Training and deploying deep learning models can be computationally expensive and resource-intensive.
Future research directions include:
Developing more robust and generalizable models: Models need to be able to adapt to evolving deepfake techniques and perform well across different types of deepfakes.
Exploring new data augmentation techniques: Generating more diverse and realistic training data can improve the performance of deepfake detection models. Investigating alternative model architectures: Exploring emerging deep learning architectures and techniques, such as transformers and attention mechanisms, may lead to further breakthroughs in deepfake detection.
In conclusion, ML models have proven to be a powerful tool for deepfake detection. With continued research and development, ML-based solutions have the potential to effectively combat the growing threat of deepfakes and ensure the integrity of online content.

## Conclusion

As deepfake content has become more prevalent, people's faith in media content has declined because seeing is no longer enough to believe. This issue has wide-ranging consequences, which could cause distress and negative impacts on those who are targeted, promote hate speech and disinformation, and potentially incite conflict over politics and violence. In the current context, the severity of this issue is more acute because deepfake content and altered videos can now be easily produced thanks to technology. Moreover, social media networks have the capacity to quickly spread this kind of content throughout the world. To address this issue, this paper provides a current overview of deepfake creation and detection techniques. This research paper looks at the problems and developments in the topic of deepfakes, providing valuable information to artificial intelligence researchers throughout the globe. It will assist them in creating practical methods for identifying deepfakes and protecting people from the possible harm they may bring.

## References

1] Oscar De Lima, Sean Franklin, Shreshtha Basu, Blake Karwoski, and AnnetGeorge. Deepfake detection using spatiotemporal convolutional networks.
arXiv preprint arXiv:2006.14749, 2020.

[2] Luca Guarnera, Oliver Giudice, and Sebastiano Battiato. Deepfake detection by analyzing convolutional traces. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops, pages 666–667, 2020.

[3] Siwei Lyu. Deepfake detection: Current challenges and next steps. In 2020 IEEE international conference on multimedia & expo workshops (ICMEW), pages 1–6. IEEE, 2020.

[4] Thanh Nguyen, Cuong M. Nguyen, Tien Nguyen, Thanh Duc, and Saeid Nahavandi. Deep learning for deepfakes creation and detection: A survey. 09 2019.

[5] Md Shohel Rana, Mohammad Nur Nobi, Beddhu Murali, and Andrew H Sung. Deepfake detection: A systematic literature review. IEEE access, 10:25494–25513, 2022.

[6] Hanqing Zhao, Wenbo Zhou, Dongdong Chen, Tianyi Wei, Weiming Zhang,and Nenghai Yu. Multi-attentional deepfake detection. In Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, pages 2185–2194, 2021.