# Code of silence: cyber security strategies for combating deepfake disinformation

RAJEEV  KUMAR

SUHEL AHMAD  KHAN

NAWAF  ALHARBE

RAEES AHMAD  KHAN

**Abstract:** At a time when deepfake technology is extensively employed, this study explores the serious problems produced by deepfake misinformation and recommends a comprehensive cyber security architecture to lessen its effects. Deepfakes are a serious threat to the authenticity of digital information because they use sophisticated artificial intelligence algorithms that have outlived their original novelty. This study examines the psychological effects of manipulated media, sheds light on the intricacies of creating deepfakes, and examines actual incidents that highlight the pervasive effects of deepfake misinformation. Our suggested cyber security plans include cutting-edge detection algorithms, blockchain technologies and extensive outreach programmes. Through promoting a shared dedication to openness and truth, symbolised by the figurative 'code of silence', this study seeks to strengthen the digital environment against the ubiquitous impact of misleading content, thereby adding to the larger conversation about preserving authenticity and trust in an increasingly digital age.

Deepfake technology has emerged as a new and severe threat to the integrity of our digital landscape in an era where digital communication and information exchange are paramount. Artificial intelligence algorithms have created synthetic media that have quickly progressed from novelty to sophisticated instruments that can create audio-visual content. This appears genuine but is completely fake. Though this technology has great potential for the creative and entertainment industries, this revolutionary power has also given rise to a dangerous kind of misinformation.

The emergence of deepfake disinformation presents a serious risk to people, institutions and the foundation of public confidence. In an era where determining the legitimacy of digital content is becoming more difficult, there is an increased need for efficient cyber security techniques to stop deepfake misinformation. This article examines the complexities of deepfake technology, examines the difficulties that arise from its harmful application, and suggests inventive cyber security protocols to prevent the spread of misleading material.

As we undertake this investigation, it becomes clear that the code of silence we speak of is not one of confidentiality but rather a shared resolve to shatter the hush of deceit. We hope to create a climate in which the code of silence surrounding deepfake disinformation is broken down and replaced by a strong defence that preserves the values of truth, transparency and digital authenticity through the creation and application of cutting-edge cyber security techniques.

## The deepfake menace

Cutting-edge machine learning algorithms are used in deepfakes, along with other deceptive modifications of visual and audio information, making it harder and harder to distinguish between real and fake media.[1] Using this technology, bad actors can create false narratives, harm people's reputations, influence public opinion and even manipulate financial markets. The once-accepted 'code of silence' around cyber security needs to be broken down in order to protect against these pernicious dangers as the frequency of deepfakes increases.

It's important to understand that deepfake technology is not limited to well-known targets, in order to fully appreciate the seriousness of the situation.[2] The harmful impacts of distorted media can affect everyone, including celebrities, public figures and regular people.[3] As a result, in the face of a growing digital arms race, it becomes imperative to strengthen cyber security safeguards.

## Issues and challenges

**Rapid evolution of deepfake technology:** One of the main obstacles is the deepfake technology's quick development. Cyber security methods find it difficult to keep up with deepfake algorithms as they continue to surpass current detection mechanisms with increasing sophistication. Because these developments are dynamic, defensive measures must be updated and improved on a regular basis.

**Multimodal complexity:** Deepfakes are not only about manipulating one mode of data – they frequently incorporate text, video and sometimes audio. For cyber security tactics, addressing the multimodal complexity of deepfakes is a big challenge. It is still difficult to create detection and authentication systems that are both comprehensive and capable of efficiently analysing and validating content in a variety of modalities.

**Adversarial training and countermeasures:** Malicious actors are using adversarial training techniques more frequently to make deepfakes more resistant to detection systems. A game of cat and mouse is thereby introduced between those who are producing deepfakes and those who are constructing defence mechanisms. In the field of deepfake cyber security, there is a continuing need for countermeasures to keep up with adversary approaches.

**Ethical implications and privacy concerns:** The application of aggressive deepfake detection technologies raises moral concerns, particularly those related to individual privacy. Finding a balance between the necessity to identify and combat false information and the protection of privacy is a challenging task. In order to preserve public confidence, cyber security strategies must adhere to legal and ethical guidelines.

**Limited standardisation and regulation:** Cyber security tactics are less effective when there are no regulated frameworks or established methods in place to counteract deepfake material. It is difficult to mount a united front against this global threat because of the fragmented picture created by inconsistent legal definitions, jurisdictional issues and different levels of commitment among nations.

**Educational gaps and awareness:** One major problem is that the general public does not know or comprehend the existence of phoney technology or its repercussions. To enable people to assess online content critically and support the robustness of the digital ecosystem, educational gaps must be closed in order to foster digital literacy and scepticism.

**Resource intensiveness:** Robust cyber security tactics against deepfake deception require significant resources, including money for technology purchases and highly qualified labour. Many businesses, especially smaller ones with tighter budgets, would find it difficult to implement and maintain state-of-the-art defences, which would leave them open to sophisticated deepfake attacks.

**Cross-border collaboration and information-sharing:** Since deepfake dangers frequently cross national boundaries, effective international cooperation is needed to counter them. The establishment of a unified worldwide defence against deepfake disinformation is hampered by the resistance to or difficulties with cross-border information exchange and cooperative efforts among governments, business organisations and cyber security specialists. Resolving these issues is essential to the success of the suggested 'code of silence'.

## Significance

These issues are extremely important in the current digital environment since deepfake disinformation is becoming a serious concern. Deepfake technology is developing quickly, and since it can change audio-visual content with previously unheard-of realism, there are reasons to be seriously concerned. The suggested cyber security tactics described in this study are crucial for a number of reasons.

First, there is the matter of preserving trust in digital communications. In a time when digital communication is essential to the spread of knowledge, maintaining trust in systems and services is critical. There is a pressing need to protect the legitimacy of digital content by providing tactics that, when put into practice, can lessen the damage that fake misinformation does to people's trust.

We also need to protect individuals and organisations from harm. Disinformation that is deeply fabricated has the ability to do a great deal of harm, including threats to national security, financial losses and reputational injury. The tactics suggested in this article seek to strengthen people and organisations against these dangers, offering a defence system that is essential for safeguarding individual and business interests.

Mitigating societal instability is another issue as deepfake misinformation has societal repercussions that go beyond its effects on individuals and organisations. Deepfake-fuelled misinformation operations in politics have the potential to sow division, taint elections and cause social unrest. This study hopes to add to the larger endeavour to lessen these dangers and preserve the stability

of democratic processes.

Encouraging responsible technology development is crucial. This study highlights the significance of conscientious AI development and application by tackling the obstacles presented by counterfeit technology. To guarantee that scientific improvements in AI benefit society without enabling malevolent intent requires striking a balance between technological innovation and ethical issues.

Alongside this is the matter of fostering digital literacy and awareness. It's important to emphasise how crucial awareness and education are in reducing the harmful effects of fake news. The recommended techniques seek to foster a more resilient and educated digital society by empowering individuals to distinguish between genuine and manipulated content and by increasing digital literacy.

The biggest achievements will be made by those involved contributing to global cyber security efforts. In light of the global reach of deepfake risks, we need to promote international cooperation and the creation of standardised cyber security methodologies. It is imperative to present a united front on a worldwide basis against deepfake deception, and the suggested tactics in this study add to the existing discussion about implementing efficient cyber security measures globally.

Finally, there is the business of innovating technological solutions. The aim of this study is to advance the continuous development of cyber security technology. The suggested tactics make use of state-of-the-art developments such as blockchain technology and deepfake detection algorithms, demonstrating how technological solutions are always evolving to counter new dangers in the digital sphere.

## Integrated framework

The comprehensive framework we are calling the 'code of silence' aims to provide a cohesive and multifaceted strategy for addressing the growing threat that deepfake disinformation poses. To develop a strong defence against the misuse of digital content, this framework combines cutting-edge technological solutions, educational programmes, ethical considerations and international collaboration. The following are some of this integrated framework's essential elements.

**Advanced detection algorithms:** Apply cutting-edge deepfake detection techniques that make use of neural networks and machine learning. These algorithms examine audio and visual material for irregularities, paying particular attention to minute discrepancies that can point to the existence of deepfake manipulation. This involves collaboration with the research community and frequent updates guarantee flexibility in changing deepfake approaches. Specific elements include:

> Research and development collaboration: Collaborate with the research community to stay abreast of the latest advancements in deepfake detection techniques and establish partnerships with academia and industry experts for continuous knowledge exchange.

> Neural network and machine-learning implementation: Implement cutting-edge neural networks and machine-learning algorithms for deepfake detection and train algorithms using diverse datasets to enhance accuracy and adaptability.

> Audio material examination: Develop algorithms specifically designed to examine audio material for irregularities indicative of deepfake manipulation; utilise signal-processing techniques to identify discrepancies in voice patterns and content.

> Visual material examination: Create algorithms tailored for the examination of visual material, focusing on minute discrepancies in facial expressions, movements, or other visual cues; leverage image-processing techniques to detect anomalies associated with deepfake content.

> Flexibility through frequent updates: Establish a system for regular updates to the deepfake detection algorithms; adapt the algorithms to evolving deepfake techniques through continuous learning and improvements.

> Agile response framework: Develop a responsive framework that allows quick adjustments to detection algorithms based on emerging deepfake threats; implement a feedback loop for information exchange with the research community to ensure agility.

**Multimodal authentication systems:** We need to create and implement multimodal authentication systems that include

behavioural biometrics, speech recognition and facial recognition. By combining these several authentication techniques, identity verification becomes more accurate and dependable, increasing the difficulty for malevolent actors to create fake media. Key elements include:

Define authentication components: Identify key components, including behavioural biometrics, speech recognition and facial recognition, for multimodal authentication.

Research and development: Conduct research on state-of-the-art technologies and methodologies related to behavioural biometrics, speech recognition and facial recognition; develop or adopt cutting-edge algorithms and models for each authentication component.

Integration planning: Plan the integration of behavioural biometrics, speech recognition and facial recognition into a cohesive multimodal authentication system.

System architecture design: Design the overall architecture of the multimodal authentication system, considering the seamless integration of individual authentication components.

Data collection and training: Gather diverse datasets for each authentication component to train the respective models effectively; train models on collected data to enhance accuracy and reliability.

Implement multimodal system: Implement the multimodal authentication system, combining behavioural biometrics, speech recognition and facial recognition.

Testing and validation: Conduct rigorous testing to ensure the accuracy and dependability of the multimodal authentication system; validate the system's effectiveness in preventing fake media creation.

Optimisation: Fine-tune algorithms and parameters for optimal performance; address any identified weaknesses or limitations through optimisation.

User interface integration: Integrate the multimodal authentication system seamlessly into user interfaces to ensure user-friendly experiences.

Documentation and training: Create comprehensive documentation for system usage and maintenance; provide training for end users and administrators on utilising the multimodal authentication system.

Security measures: Implement additional security measures to protect the multimodal authentication system from potential attacks or manipulation.

Continuous improvement: Establish mechanisms for ongoing monitoring and improvement, incorporating user feedback and technological advancements.

**Blockchain technology for content verification:** We can make use of blockchain technology to create a decentralised, impenetrable digital content record. This blockchain-based technology offers an unchangeable ledger for media files, ensuring the accuracy of the data. By automating the verification process, smart contracts enable users to safely and independently verify the legitimacy of transmitted content. The following are the key elements:

Blockchain implementation: Develop and deploy a blockchain network for content verification; utilise decentralised nodes to create a distributed ledger that records transactions.

Digital content record: Implement the blockchain as a secure and immutable digital content record; ensure that every piece of media file data is time-stamped and appended to the blockchain.

Data accuracy assurance: Leverage the blockchain's characteristics to guarantee the accuracy and integrity of the stored data; employ cryptographic hashing to create unique identifiers for media files.

Smart contract integration: Develop smart contracts to automate the verification process; define rules within smart contracts to validate the legitimacy of transmitted content.

User verification process: Enable users to independently verify the authenticity of content through smart contracts; establish user-friendly interfaces for interacting with smart contracts and accessing blockchain

records.

Security measures: Implement robust security measures to protect the blockchain network from tampering or unauthorised access; utilise encryption techniques to safeguard the privacy and security of the stored content.4

Decentralisation principles: Adhere to decentralisation principles to ensure the resilience of the blockchain against single points of failure; promote a distributed network to enhance the overall security and trustworthiness of the digital content record.

Automation of verification: Enable automated verification processes triggered by smart contracts to streamline and expedite the content verification workflow; ensure that the system can efficiently handle a large volume of verification requests.

**Ethical considerations and privacy preservation:** When developing and implementing deepfake detection tools, we must take ethics into account. Put user privacy first by using privacy-preserving methods and making sure that cyber security tactics are implemented in accordance with the law. Retaining public confidence requires striking a balance between the necessity for security and the right to personal privacy. There are multiple aspects to this:

Ethical considerations in tool development: Consider ethical implications during the development of deepfake detection tools; ensure that the tools adhere to ethical standards and principles.

User privacy priority: Prioritise user privacy as a fundamental principle in tool design and implementation; establish privacy-centric guidelines for the entire development lifecycle.

Privacy-preserving methods: Integrate privacy-preserving methods into the design of deepfake detection tools; employ techniques that minimise the exposure of sensitive user data.

Legal compliance in cyber security tactics: Ensure that the cyber security tactics employed align with legal frameworks and regulations; regularly update tactics to stay compliant with evolving laws.

Balancing security and privacy: Strike a balance between the necessity for security measures and the protection of personal privacy; implement measures that effectively counter deepfake threats without compromising individual privacy rights.

Building public confidence: Prioritise actions that contribute to building and maintaining public confidence; communicate transparently about the ethical considerations and privacy safeguards implemented.

**Comprehensive education and awareness campaigns:** We need to initiate educational programmes to increase knowledge of deepfake disinformation's existence and possible effects. We also need to encourage digital literacy so that people can assess Internet content critically, as well as foster a culture of informed digital citizenship by promoting scepticism and appropriate sharing habits, actively enhancing the resilience of the digital ecosystem.[5] There are some important steps to achieving these goals:

Initiate educational programmes: Develop and implement structured educational programmes aimed at raising awareness about the existence and potential impacts of deepfake disinformation; design curriculum content that provides insights into the techniques used in creating deepfakes and their potential consequences.

Encourage digital literacy: Promote digital literacy initiatives to enhance individuals' ability to critically assess content on the Internet, including recognising potential signs of deepfakes; provide resources and training materials to improve users' understanding of the digital landscape and emerging threats.

Foster a culture of informed digital citizenship: Launch campaigns that advocate for responsible digital behaviour, emphasising the importance of being informed and vigilant online; encourage users to verify information before sharing; promote responsible content-sharing habits.

Promote skepticism: Instil a sense of skepticism regarding online content through awareness campaigns,

emphasising the need to question the authenticity of media; provide tools and resources to help users identify and evaluate potential deepfake content.

Enhance resilience of the digital ecosystem: Advocate for collaborative efforts to build a resilient digital ecosystem that can withstand the impact of deepfakes; promote the adoption of secure communication practices and technologies to mitigate the spread of disinformation.

**Policy development and international collaboration:** We should encourage the design and implementation of laws that deal with the production and distribution of false information, and work together with foreign partners to create an all-encompassing legal system that cuts across national borders. We need to encourage information exchange and teamwork to counteract cyberthreats worldwide. Key components of this element include:

Policy development: Encourage the design and development of laws specifically addressing the production and distribution of false information; collaborate with legal experts, policymakers and relevant stakeholders to create robust legal frameworks tailored to counteract the threat of deepfake disinformation.

International collaboration: Foster collaboration with foreign partners, including governments, international organisations and law-enforcement agencies; establish mechanisms for sharing information and intelligence on emerging deepfake threats; and explore the creation of bilateral or multilateral agreements to facilitate coordinated responses across national borders.

All-encompassing legal system: Work towards the creation of a comprehensive legal system that transcends national borders, ensuring a unified approach to combatting deepfake-related offences; address jurisdictional challenges by promoting standardised legal norms applicable globally.

Information exchange: Encourage the exchange of critical information related to deepfake threats among participating countries; establish secure channels and protocols for the timely sharing of intelligence on potential risks and vulnerabilities.

Teamwork for cyberthreat mitigation: Promote collaborative efforts among international partners to enhance cyber security capabilities; facilitate joint initiatives, such as joint investigations and operations, to counteract deepfake-related cyberthreats effectively.

**Continuous monitoring and adaptation:** To achieve our goals, we need to install a system of ongoing observation to evaluate cyber security tactics' efficacy instantly. This would include frequently updating detection algorithms to take advantage of new threats and tricks used by adversaries. And it would involve encouraging cooperation between government organisations, tech developers and cyber security specialists to stay ahead of the ever-changing deepfake issues. Important elements of this are:

Install an ongoing observation system: Establish a continuous monitoring system to assess the effectiveness of cyber security tactics in real time; implement tools and technologies for constant surveillance of digital content and potential deepfake threats.

Evaluate cyber security tactics: Regularly assess the performance of existing cyber security tactics in detecting and mitigating deepfake threats; analyse the results of ongoing monitoring to identify areas for improvement in the overall security posture.

Update detection algorithms: Develop a mechanism for frequent updates to deepfake detection algorithms; stay abreast of the latest advancements in neural networks, machine learning and other technologies relevant to deepfake detection.

Address new threats and tricks: Monitor emerging trends and tactics employed by adversaries in creating deepfake content; incorporate insights from threat intelligence to adapt detection algorithms and strategies accordingly.

Encourage cooperation: Foster collaboration between government organisations, tech developers and cyber security specialists; establish forums, working groups, or collaborative platforms to facilitate information exchange and joint efforts.

Stay ahead of deepfake issues: Promote a proactive approach to deepfake mitigation by anticipating potential threats; encourage continuous learning and knowledge-sharing to ensure that the collective defence against deepfakes evolves with the threat landscape.

## Future emerging strategies

Anticipating the dynamic terrain of deepfake deception, the code-of-silence paradigm recognises the necessity of ongoing attention to detail, and flexibility. It will be necessary to incorporate future emerging approaches into this cyber security plan in order to keep ahead of persistently improving bad actors. The following factors should be taken into account when updating the framework to handle new threats.[6,7]

**Deepfake generative adversarial network (GAN) countermeasures:** Future deepfakes might use more sophisticated GANs, which would make it harder to identify fake material. The primary goal of research and development should be to develop counter-GAN technologies that can recognise and neutralise the improved generative capabilities of upcoming deepfake models.

**Zero-day vulnerability response:** The system has to incorporate techniques that resolve zero-day vulnerabilities unique to deepfake detection algorithms as cyberthreats change. In order to generate fast fixes and upgrades to fight evolving threats, it will be essential to establish rapid reaction teams and collaborate with the cyber security community.

**Explainability in detection algorithms:** To increase user trust and comprehension, we need to make deepfake detection algorithms more explainable. Transparency should be given priority in future iterations of the system, with detailed explanations of how computers detect fraudulent content. This can enable consumers to make well-informed decisions and aid in the continuous improvement of detection techniques.

**Decentralised content verification platforms:** Examine the possibilities of blockchain-based decentralised content verification systems. Prospective approaches could encompass the creation of decentralised applications (DApps) that empower users to autonomously confirm the legitimacy of media material without depending on centralised authority, hence augmenting the robustness of the framework.

**Quantum-safe cryptography:** Be prepared for the arrival of quantum computing, which may threaten current cryptography techniques. In order to guarantee the ongoing security and integrity of blockchain-based verification systems as well as other cryptographic components, the code-of-silence framework should include quantum-safe cryptographic algorithms.

**Enhanced behavioural biometrics:** In the future, deepfakes may try to imitate not just voice and face features but also subtle behavioural clues. To strengthen the resistance against complex deepfake attacks, we must integrate advanced behavioural biometrics into multimodal authentication systems, such as mouse movements and typing patterns.

**Dynamic policies and regulations:** We must acknowledge the ever-changing nature of technology and the necessity of flexible laws and rules. Agile policymaking that can react quickly to new dangers while maintaining moral principles and individual privacy should be a key component of future strategies.

**Collaborative AI platforms:** We need to create cooperative artificial intelligence (AI) systems that combine data and knowledge from various sources. The framework can gain collective insights by establishing networked AI systems, which makes it possible to respond more coordinatedly to the hostile actors' shifting strategies in the deepfake arena.

**Immersive technologies and 3D models:** Expect to see 3D models and immersive technologies incorporated with deepfake material. Future tactics should be able to identify manipulations in virtual reality (VR) and augmented reality (AR) settings, where conventional detection techniques can encounter particular difficulties.

**Global cyber security standards:** We must encourage the creation of international guidelines for the mitigation and detection of deepfakes. In order to ensure a coordinated and successful response to cross-border deepfake threats, cooperation with governments and international groups will be crucial.

## Conclusions

As our increasingly digital culture continues to seek authenticity and truth, the code-of-silence paradigm provides a comprehensive and flexible response to the growing threat of deepfake deception. The tactics described in this framework represent a shared commitment to strengthening our digital ecosystem against the pernicious influence of synthetic media as we negotiate the complex terrain of AI-powered manipulations.

One cannot over-stress the importance of battling deepfake deception. The potential ramifications, which include harm to one's reputation, instability in society and a decline in trust, highlight the necessity of taking proactive cyber security measures. The integrated framework presented here makes use of blockchain technology, multimodal authentication systems, sophisticated detection algorithms and ethical considerations to produce a comprehensive defence mechanism.

The code of silence is more than just a figure of speech – it's a commitment to ending the silence of dishonesty. This framework aims to establish a digital environment where truth wins out over manipulation by promoting international collaboration, protecting user privacy, enhancing detection technology and encouraging digital literacy.

The architecture recognises that new risks will inevitably arise in the future and prepares for strategies such as counter-GAN technology, quantum-safe cryptography and cooperative AI platforms. The code of silence aims to be at the forefront of the ongoing fight against deepfake disinformation by adopting agility in policymaking, integrating immersive technology detection and advocating for global standards.

To sum up, the code of silence is a living, breathing example of our dedication to digital authenticity. Technology must progress, and our defences must too. Through the application of the tactics delineated in this framework, we strive to maintain the authenticity of digital communication, cultivate a responsible culture and preserve the trust that underpins our globalised society. We are starting a journey towards a future where the digital world is characterised by resilience, openness and an uncompromising dedication to truth by shattering the silence of falsehood.

## About the authors

*Dr Rajeev Kumar is an assistant professor in the Centre for Innovation and Technology at the Administrative Staff College of India, Hyderabad, Telangana, India.*

*Dr Suhel Ahmad Khan is currently working as an assistant professor in the Department of Computer Science at the Indira Gandhi National Tribal University (A Central University), Amarkantak, Madhya Pradesh, India.*

*Dr Nawaf Alharbe is an associate professor in the Applied College, Taibah University, Medina, Saudi Arabia.*

*Prof Raees Ahmad Khan is a professor in the Department of Information Technology at the Babasaheb Bhimrao Ambedkar University (A Central University), Lucknow, Uttar Pradesh, India.*

## References:

1.Riehle, Cornelia. 'Europol report criminal use of deepfake technology'. EUcrim, 9 May 2022. Accessed Apr 2024. https://eucrim.eu/news/europol-report-criminal-use-of-deepfake-technology/.

2.Dahiya, Yash. 'The Rise Of Deepfake Technology: A Threat To Evidence In Arbitration?'. Live Law, 22 Nov 2023. Accessed Apr 2024. www.livelaw.in/articles/the-rise-of-deepfake-technology-a-threat-to-evidence-in-arbitration-242718, www.livelaw.in/articles/the-rise-of-deepfake-technology-a-threat-to-evidence-in-arbitration-242718.

3.Shankar, Vasundhara. 'Deepfakes call for stronger laws'. The Hindu Business Line, 16 Jul 2023. Accessed Apr 2024. www.thehindubusinessline.com/business-laws/deepfakes-call-for-stronger-laws/article67077019.ece.

4.Ahmad Khan, S; Kumar, R; Ahmad Khan, R. 'Software Security: Concepts & Practices'. Chapman and Hall/CRC, 2023. Accessed Apr 2024. https://doi.org/10.1201/9781003330516.

5.Kumar, R; Ahmad Khan, S; Ahmad Khan, R. 'Software Durability: Concepts and Practices'. CRC Press, 2023. Accessed Apr 2024. https://doi.org/10.1201/9781003322351.

6.Sivaraman, R. 'Tamil Nadu cybercrime police issue advisory on deepfake scams'. The Hindu, 8 Aug 2023. Accessed Apr 2024. www.thehindu.com/news/national/tamil-nadu/tn-cybercrime-police-issue-advisory-on-deepfake-scams/article67171553.ece.

7.Varma, Vishnu. 'Kerala deepfake fraud case: Efforts on to nab accused from Gujarat, say cops'. Hindustan Times, 15 Aug 2023. Accessed Apr 2024. www.hindustantimes.com/india-news/aibased-deepfake-scammer-identified-accused-of-cheating-elderly-man-in-kerala-police-launch-manl