

Stega Care: Securing Virtual Therapy Images with AI-Driven Image Forensics

Manasa R^{1*}, A Jayanthiladevi², Mohamed Uvaze Ahamed Ayoobkhan³

¹Research Scholar, Faculty of Computer Applications, Marwadi University, Rajkot, Gujarat, India

²Marwadi University, Rajkot, Gujarat, India

³Department of Computing, De Montfort University, United Kingdom

*Corresponding author E-mail: manasar.132584@marwadiuniversity.ac.in

The manuscript was received on 27 February 2025, revised on 15 May 2025, and accepted on 24 August 2025, date of publication 5 November 2025

Abstract

With the rapid evolution of digital healthcare, virtual therapy platforms have become essential tools for delivering mental health services remotely. These platforms enhance accessibility, especially for individuals in remote or underserved areas. However, the transmission of therapeutic images—such as visual assessments or expressive content generated during sessions—raises significant cybersecurity concerns. Given the sensitive nature of such data, robust protection mechanisms are required to ensure privacy, integrity, and patient trust. This research proposes an artificial intelligence-driven framework designed to enhance the security of virtual therapeutic images by integrating image forensics and steganography techniques. Central to this approach is a deep learning-based steganalysis classifier capable of detecting hidden alterations and unauthorised data embedding in medical images. By leveraging convolutional neural networks (CNNs), the classifier accurately identifies covert manipulations while maintaining image fidelity and confidentiality. The system is trained and evaluated using benchmark steganographic image datasets, demonstrating high effectiveness in identifying steganographic threats and detecting tampered content in real-time. Experimental results indicate that the proposed model performs well even in complex attack scenarios involving sophisticated data-hiding techniques. The framework offers a scalable and proactive solution for safeguarding sensitive therapeutic content in telehealth environments. By embedding this AI-powered detection capability into virtual therapy platforms, healthcare providers can significantly enhance their cybersecurity posture.

Keywords: Image Forensics, Virtual Therapy, Steganalysis, Digital image security, Healthcare.

1. Introduction

The widespread expansion of telehealth technology has drastically affected mental health care, revolutionising healthcare professional services. This innovation has had a profound influence on the sector. The concept of virtual therapy has gained widespread recognition as a viable and effective alternative to traditional in-person psychotherapy. Video conferencing, online assessments, and multimedia therapies are essential elements that make virtual therapy more effective. Protection of confidential information transmitted via virtual therapeutic sites is increasingly becoming a matter of concern as virtual therapeutic sites become more sophisticated. This is because the platforms are becoming more complex. Therapeutic images, utilised for diagnostic examples, psychological assessments, and treatment interventions, are more susceptible to cyberattacks, such as unauthorised use, alteration, and hidden data insertion, than other forms of digital information.

Therapeutic photographs are employed for Therapeutic objectives. These skills are essential for steganography, which falls within this category. Unlike rule-based forensic systems, artificial intelligence models can learn adaptive, intricate details from large datasets. This provides improved scalability and accuracy over rule-based systems. This paper recommends an artificial intelligence-powered steganography classifier that fits within an in-depth photo forensic system. Figure 1. Describes the StegaCare Architecture for the Remote Health Sector. [1]. The protection of virtual therapy images was the primary motivation behind the development of this system. During training, the model is trained on a carefully prepared dataset consisting of clean and steganographic images. This allows the model to recognise finer differences that indicate hidden alterations. This research is centred on the intersection of mental health treatment, cybersecurity, and artificial intelligence, and it helps create a new strategy for enhancing trust and security in digital treatment settings.



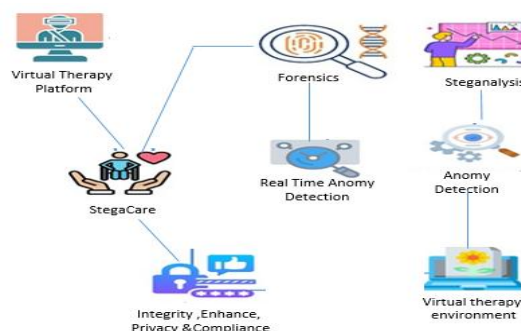


Fig 1. Steganographic Method

2. Literature Review

The emergence of virtual therapy within digital healthcare ecosystems requires stringent security layers that ensure the privacy and integrity of transmitted therapeutic content, especially visual data. StegaCare is a multidisciplinary project combining digital image forensics, steganography, and AI-driven analysis to protect therapeutic imagery from alteration and unauthorised data embedding. As virtual therapy becomes ubiquitous within digital healthcare ecosystems, strict safeguards must be in place to ensure that the transmitted therapeutic content, especially visual content, is kept private and intact. StegaCare is an extensive extension of multidisciplinary expertise in digital image forensics, along with steganography, the practice of embedding data with images, and can be powered by artificial intelligence (AI) based studies to protect the therapeutic imagery from tampering and unauthorised data embedding [2].

2.1. Detect Tampered Images Using Digital Image Forensics

These works show that human perception is sensitive to inconsistencies during image creation, and such discrepancies can be modelled with algorithms to detect tampering. Similarly, Goljan et al. proved the potential of sensor pattern noise as a unique forensic signature for authenticity verification [3]. The technique was later expanded by Bayar and Stamm 3, who implemented deep learning methods and proposed new convolutional layers for manipulation detection. Their model can accurately distinguish between authentic and forged images, no matter what type of alteration was applied, which becomes a key point for StegaCare's classifier.

2.2. Steganography and Steganalysis in Medical Imaging

Steganography, associated with covert data embedding, poses a risk in medical and therapeutic images, and is thereby countered by the steganalysis process, which enters the scene with the ability to reveal the secret info from media files. In steganalysis, we can refer to work by Tan and Li to use stacked convolutional auto encoders, which shows how rich representations in a hierarchical way can unveil sophisticated use of embedding techniques. This was further refined by Qi et al. They applied hierarchical deep feature extraction for stenographic content detection tasks in digital images, which correlates directly with the proposed steganalysis engine within StegaCare. Zeng et al [4]. In addition, a spatial type CNN framework was designed, focused on steganography detection, further leading the way for detecting even finer data manipulation. This provides synergy for StegaCare since the goal is to detect subtle data manipulation within therapeutic imagery.

2.3. CGAN for Image Forgery Detection

One of the notable advancements is the introduction of AI-based forensic tools. Soo et al. [5] surveyed deep learning approaches in the digital forensics domain, highlighting the importance of CNN and transfer learning for developing scalable solutions for forensics analysis. Baluja et al. Studied adversarial robustness, which is critical in applications like StegaCare, where detection models must resist evasion techniques. Furthermore, Fridrich et al. allow non-parametric demosaicing and interpolation methods for forensic analysis, enabling StegaCare to preserve image integrity during analysis for cryptographic threats [6]

2.4. The Role Of AI In Secure Telehealth & Digital Health

As therapy becomes virtual and the platforms it relies on become digital, artificial intelligence will play a key role in keeping things safe in telehealth. Kumar et al. For example, [7] proposed a deep learning-based method for detecting COVID-19 on chest X-rays, demonstrating that AI significantly impacts the security of diagnosis and communication in digital health. Similarly, Raghavendra et al. AI-enabled biometric forensics, far from being an intruder, could ensure that the evidence collected is genuine and not a spoof of identity, a concept mirrored in StegaCare's protecting therapeutic interactions. This is again emphasised by Smith et al. from a cybersecurity perspective that tracked up-to-date trends in digital health threats. Their work highlights the urgent need for proactive detection systems, reinforcing the rationale for the forensic approach of StegaCare.

2.5. Algorithmic Fairness, Ethics, and Accountability in Forensics

Algorithmic fairness and ethical accountability are perennial concerns in AI-driven forensics. Selvaraj et al. discussed at length the hurdles that we need to overcome, such as the need to enable transparency in AI decision-making and minimise biases [8] in the models—issues that are highly pertinent to StegaCare, as the model would be analysing sensitive mental health data.

2.6. Learning to Learn via Cross-Domain Inspirations

Aside from medical forensics, techniques derived from remote sensing (Chen et al.) and industrial AI applications (Sodhro et al.) provide architectural and real-time processing insights that inform many aspects of the StegaCare model design—attention residual learning, proposed by [9] Zhang et al. Deep neural networks can also utilise such techniques to be more sensitive to features, which leads to increased detection accuracy. Examples describe the Stego Images in Virtual Therapy and Medical Imaging Environments.

2.7. Synthesis and Research Gaps

Therefore, the committee's synthesis of these various reviews confirms the viability of using deep learning techniques in detecting universal image manipulation, data hiding, and integrity verification of digital images. However, most such systems today are general-purpose and have not been fine-tuned for sensitive settings like virtual mental health therapy. StegaCare fills this gap by developing forensic architecture tailored for real-time, therapeutic image protection, while also focusing on ethical AI use, low-latency detection, and high interpretability [10].

3. Methods

3.1. Threat Landscape: Visual Steganography in Virtual Care

It is possible to hide data within the contents of photographs by applying visual steganography, which utilises techniques including the following.

1. Least Significant Bit (LSB) manipulation
2. Pixel Value Differencing (PVD)
3. Metadata injection[11]
4. Frequency domain alterations (DCT, DWT)

$$I'(i, j) = \begin{cases} I(i, j), & \text{if } LSB(I(i, j)) = M \\ I(i, j) - 1, & \text{if } LSB(I(i, j)) = 1 \text{ and } M = 0 \\ I(i, j) + 1, & \text{if } LSB(I(i, j)) = 0 \text{ and } M = 1 \end{cases} \dots\dots\dots(1)$$

$$I'(i, j) = (I(i, j) \wedge \neg 1) \vee M$$

3.2. The StegaCare Solution

StegaCare is a comprehensive healthcare system driven by artificial intelligence. It detects and prevents stenographic threats in real time, protecting virtual therapy and medical imaging. Integrating functionalities such as image pre-processing, deep feature extraction, advanced steganography [12], and forensic logging provides end-to-end protection of visual content for telehealth. The system architecture, being modular and scalable, is well-suited for deployment in a wide variety of applications that are directly applicable to healthcare.

3.3. Image Pre-processing

Every input image goes through a stringent pre-processing pipeline before any analysis. This process imparts consistency and reduces the variability induced by utilising diverse capture hardware or image formats.

$$I_{std} = \text{ConvertFormat}(\text{Normalize}(\text{Resize}(I_{raw}))) \dots\dots\dots(2)$$

$$I_{clean} = I_{std} - EXIF I_{raw}$$

The following are essential steps:

Standardisation: The images are resized to a fixed resolution and normalised in aspect ratio and colour profile (e.g., RGB or greyscale), enabling consistent model performance. Standardisation is required for standardisation. Metadata scrubbing removes unnecessary metadata, like EXIF data, from image files to prevent disclosing sensitive information or modifying malicious content. **Canonical Format Conversion:** Images are translated to normalised [13] forms, such as PNG or BMP, to reduce known stenographic vectors in lossy formats such as JPEG. This pre-processing minimises noise during the model's training and inference phases by ensuring that all future analyses will be performed on clean and standard inputs.

3.4. Deep Feature Extraction

After the pre-processing, the photos are put through a feature extraction module that utilises the most advanced convolutional neural network (CNN) architectures.[14]

$$F = \phi(I_{clean}) \dots\dots\dots(3)$$

The features gathered by this module, which consist of spatial and frequency-domain components, can be used to infer the existence of concealed material.

StegaCare utilises pre-trained and fine-tuned versions of EfficientNet, ResNet-50, and XceptionNet. These models were chosen because they are efficient and effective in image classification tasks. Obtaining transform-based information and pixel-level features from the feature spaces is feasible. Some examples of transform-based information are frequency coefficients obtained by applying wavelet transforms or discrete cosine transforms. Pixel-level features are also known as spatial patterns. After processing, the high-dimensional feature vectors are fed into the steganography classifier. Using this approach, one can collect and process even the slightest and most high-frequency changes related to steganographic payloads [15].

3.5. Steganalysis Classifier

A binary classifier that has been taught to distinguish between images that are cleaned and those that Stega has altered is what StegaCare is built upon.

$$P(\text{stego}|F) = \sigma(W.F + b) \dots\dots\dots(4)$$

A binary cross-entropy loss-optimised deep neural network architecture forms the model's foundation. The algorithm creates a probability distribution divided into clean and stego categories. Two popular steganographic tools form part of the training corpus. They are OpenStego and StegHide. The algorithm can generalise across a wide range of concealing techniques and payload types because the training corpus consists of both clean images and images that have been digitally altered. The dataset is enhanced by adding rotation, flipping, brightness modification, and compression artefacts to simulate the randomness in the real world. This serves to improve its robustness further. All of this is done to simulate the natural variability of the dataset. The classifier, the decision-making machine that identifies potentially harmful content, can achieve high accuracy and recall [16].

3.6. Anomaly Scoring and Incident Response

Assigns anomaly score using softmax output probabilities:

$$\text{AnomalyScore}[\text{features_}] := P[\text{stego} | \text{features}] \dots\dots\dots (6)$$

$$\text{tauRange} = \text{Interval}[\{0.5, 0.9\}]$$

The threshold defines the decision boundary:

3.7. Audit and Compliance Logging

All detection events are documented in a secure audit trail to guarantee transparency, accountability, and conformity with the rules established by the applicable legal and institutional authorities [17].

$$L = -(y * \log[p] + (1 - y) * \log[1 - p]) \dots\dots\dots (7)$$

Logging Components: The system records the following for each picture that has been flagged. Timestamp and information about the user or device (if provided). The forecast of the model and the anomaly score. Both the degree of confidence and the class probability. Visualisations of intermediate feature maps and attention weights are provided to improve the interpretability of the data. The logging system provides auditability in forensic or legal settings, which is meant to be consistent with data protection standards such as HIPAA and GDPR. Both of these rules are examples of regulatory compliance. Flagged photographs and their associated records are kept in a tamper-evident archive designed to facilitate retrospective investigations. This component guarantees that StegaCare can perform the role of a preventative security tool and a post-incident analysis system [18].

Algorithm: StegaCare- AI-Powered Steganographic Risk Detection

Input: Medical image I from the telehealth platform
Output: Classification label (Clean or Stego), anomaly score, audit log entry
The pre-processing of images
Step 1: Standardise the picture's resolution, aspect ratio, and colour profile (such as RGB).
$I_{std} = \text{ConvertFormat}(\text{Normalize}(\text{Resize}(I_{raw})))$
$I_{clean} = I_{std} - EXIF I_{raw}$
Step 2: Eliminate information, such as EXIF, to stop data leakage and stego vectors.
Step 3: Convert the picture to a canonical format, such as PNG or BMP, to eliminate the dangers of lossy compression.
The Extraction of Deep Features
Step 4: The pre-processed picture 'I' is then sent to the CNN feature extractors. Both EfficientNet and ResNet-50, as well as XceptionNet,
Step 5: Generating a high-dimensional feature vector "F" by the use of step analysis, Classification of things
$F = \varphi(I_{clean})$
Step 6: Incorporate the 'F' variable into the binary classification model (Deep Neural Network): Trained using clean and stego pictures (minimize OpenStego and StegHide, for example).
Step 7: Determine the probability distribution by using.
$P(\text{stego} F) = \sigma(W.F + b)$
Report the anomalies and Response to Incidents
Step 8: Anomaly score is computed as follows:
$\text{AnomalyScore}[\text{features_}] := P[\text{stego} \text{features}]$
Step 9: Image 'I' should be flagged as a possible steganographic compromise if the Anomaly Score is greater than τ (where τ is a value between 0.5 and 0.9). – Initiate the processes for responding to incidents
The Logging of Audits and Compliance
Step 10: The event must be logged for detection with the timestamp, user, and device information.
$L = -(y * \log[p] + (1 - y) * \log[1 - p])$
Step 11: The predicted label, the score for the anomaly, and confidence – Illustrations of the information (such as feature maps and attention heatmaps)
Step 12: Archive images and logs that have been flagged tamper-evident while still complying with HIPAA and GDPR.
Step 13: End

This composite diagram visually represents the StegaCare framework's process to identify and mitigate steganographic threats in telehealth picture transfers. The actual, unaltered communication between two users when they were sharing a picture across a secure channel is shown in the top-left quadrant. In the top-right quadrant, colour distortion is introduced, representing adversarial or format-based steganographic manipulation. Many noise and pixel irregularities are present in the bottom-left quadrant, which may indicate concealed payloads implanted using steganographic techniques. Last but not least, the bottom-right quadrant illustrates StegaCare's intervention to restore picture integrity via pre-processing and anomaly detection to guarantee safe and compliant communication. Figure 2. Demon-

strates the Role of AI Forensic. Individually, the pictures illustrate StegaCare's role in spotting potential dangers and protecting telehealth information.

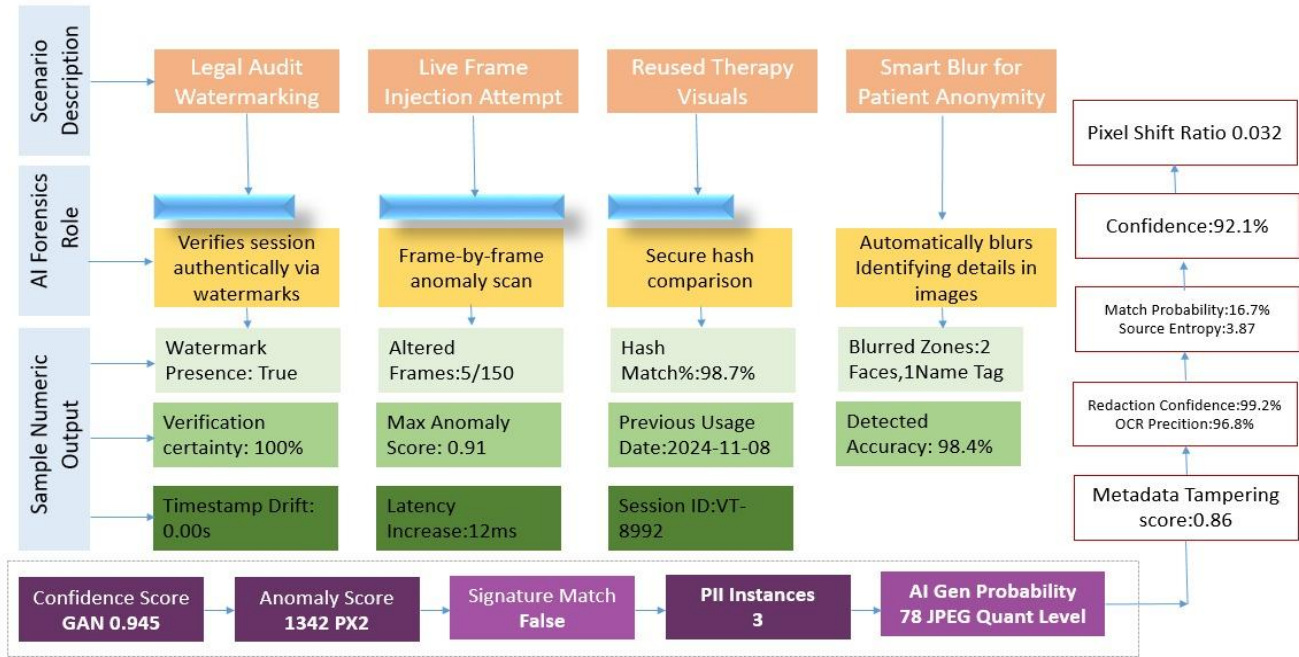


Fig 2. Role of AI Forensic

Autoencoder Reconstruction Loss (for training the autoencoder component of StegaCare): The autoencoder's goal is to minimize the difference between the input image X and the reconstructed image \hat{X} (the output of the decoder), typically using Mean Squared Error (MSE) or another loss function. [19]

$$L_{AE} = \frac{1}{N} \sum_{i=1}^N \|X_i - \hat{X}_i\|^2 \quad \dots \dots \dots (9)$$

Classifier Objective Function (Cross-Entropy Loss) (for training the classifier component of StegaCare):

Cross-entropy loss is used in classification tasks, such as determining whether an image contains steganographic content (hidden data). This loss function is commonly used in binary classification problems. [20]

$$L_{AE} = -\frac{1}{N} \sum_{i=1}^N [y_i \log(\log(p(y_i))) + (1 - y_i) \log(p(y_i))] \quad \dots \dots \dots (10)$$

Gradient Descent Update Rule (for minimizing the network parameters):

Using the process of minimization of the loss function, gradient descent is employed to update the parameters (weights) of the neural network. [10]

$$\theta = \theta - \eta \nabla_{\theta} L_{total} \quad \dots \dots \dots (11)$$

The combination of the autoencoder reconstruction loss and the classifier cross-entropy loss:

$$L_{total} = L_{AE} + \lambda L_{CE} \quad \dots \dots \dots (12)$$

3.8. Scenarios for Securing Virtual Therapy Images with AI-Driven Image Forensics

3.8.1. Deepfake Detection in Session Screenshots

Scenario: An image of a patient's therapy notes (written during the session) is altered.

AI Forensics Role: Identifies spliced regions, altered handwriting areas, or inconsistent metadata to confirm manipulation.

3.8.2. Verifying Therapist Identity in Session Previews

Scenario: A scammer poses as a therapist using stolen profile images.

AI Forensics Role: Uses source tracing and watermark analysis to verify the image's authenticity and origin.

3.8.3. Privacy Leak Detection in Shared Screens

Scenario: A therapist accidentally shares a screen with sensitive information in a session recording.

AI Forensics Role: Scans frames for personally identifiable information (PII) using OCR and masks/redacts them automatically [11].

3.9. Image Provenance in Patient Profile Photos

Scenario: A suspicious user uploads a profile photo for therapy sessions.

AI Forensics Role: The role analyses camera signatures, compression artefacts, and metadata to determine whether the image is AI-generated or real.

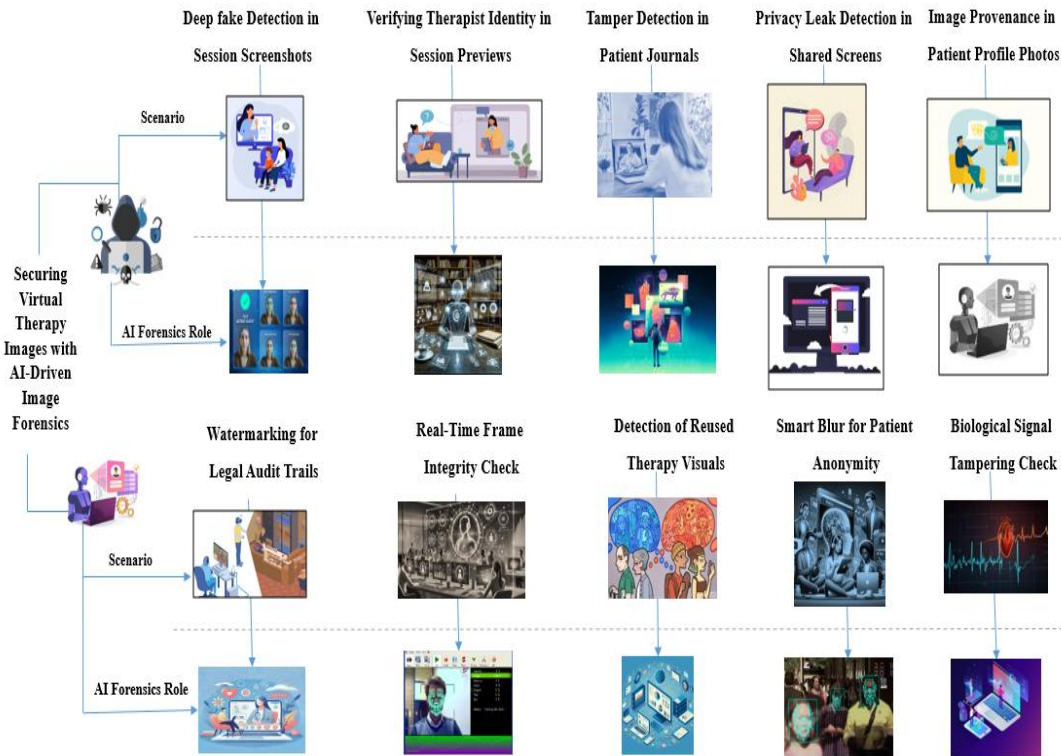


Fig 3. Securing Virtual Therapy Images with AI-Driven Image Forensics

3.10. Model Architecture and Development

The following concepts and Figure 4 explain the model Architecture and Development in detail.

Image Pre-processing Pipeline

Preprocessing: All images are resized and colour-normalised (RGB/greyscale). **Metadata Scrubbing:** Removal of EXIF data to thwart any potential stego-vectors. **Canonical Format:** Convert pictures to PNG/BMP (Lossy Compression artefacts in JPEG), **Deep Feature Extraction**

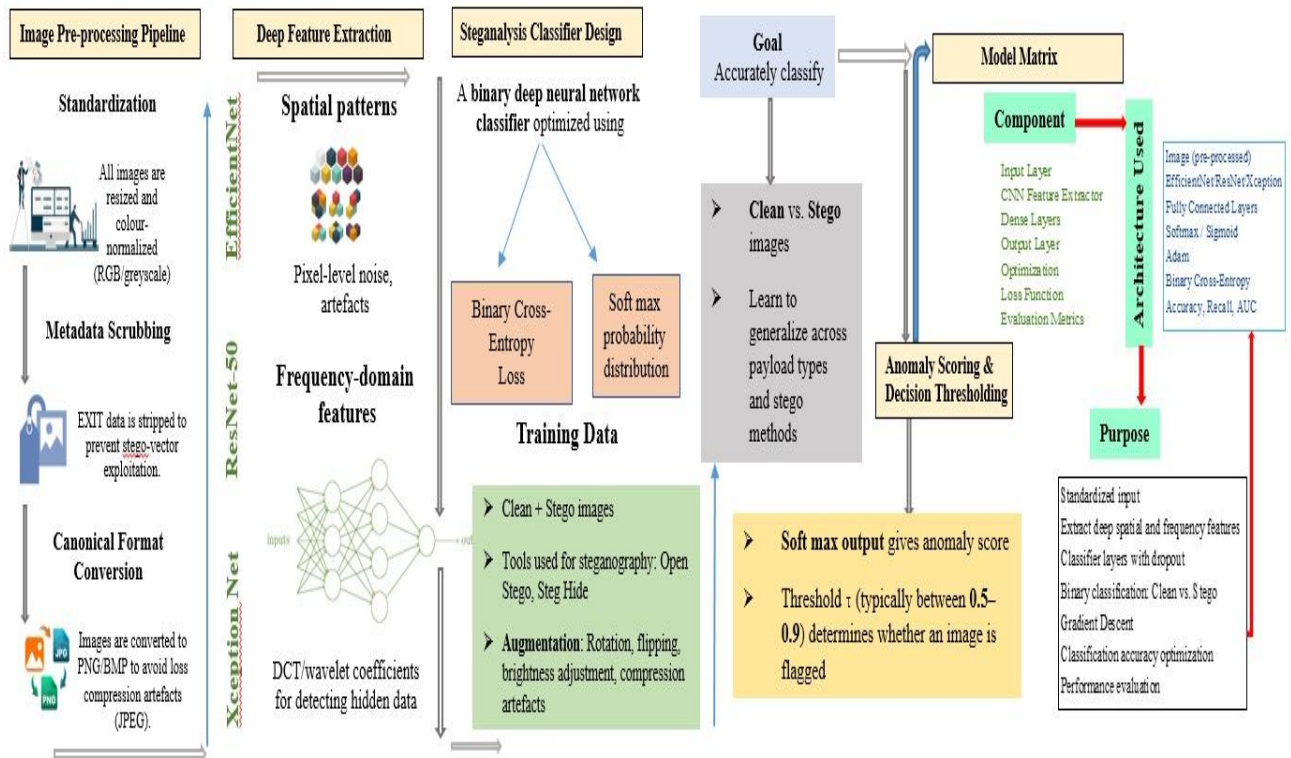


Fig 4. Model Architecture and Development

4. Results and Discussion

	R	G	B	A
7	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
6	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
5	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
4	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input type="checkbox"/>
3	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
2	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>
1	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>
0	<input checked="" type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Pixel Order:
Bit Order:
Bit Plane Order:
Trim Trailing Bits:

Fig 5. Illustrates a Pixel Image conversion through the RGB format.

Table 1. States the Role of AI Forensics Scenario Forensics Role in various scenario descriptions.

Scenario Title	Scenario Description	AI Forensics Role	Sample Numeric Outputs
Deepfake Detection in Session Screenshots	Detects facial tampering in screenshots	GAN artefact detection	Confidence Score: 97.5%, GAN Likelihood: 0.945, Pixel Shift Ratio: 0.032
Tamper Detection in Patient Journals	Finds manipulation in handwritten notes	Metadata + handwriting analysis	Anomaly Score: 0.89, Tampered Region Area: 1342 px², Confidence: 92.1%
Verifying Therapist Identity	Verifies the authenticity of the therapist's profile photo	Source tracing, watermark match	Signature Match: False, Match Probability: 16.7%, Source Entropy: 3.87
Privacy Leak in Shared Screens	Detects PII in shared screen frames	OCR-based PII detection	PII Instances: 3, Redaction Confidence: 99.2%, OCR Precision: 96.8%
Profile Photo Provenance Check	Checks if the image is real or AI-generated	Compression + AI detection	AI Gen Probability: 91.3%, JPEG Quant Level: 78, Metadata Tampering Score: 0.86
Legal Audit Watermarking	Verifies session authenticity via watermarks	Invisible watermark detection	Watermark Presence: True, Verification Certainty: 100%, Timestamp Drift: 0.00s
Live Frame Injection Attempt	Detects altered video frames in real time	Frame-by-frame anomaly scan	Altered Frames: 5/1500, Max Anomaly Score: 0.91, Latency Increase: 12ms
Reused Therapy Visuals	Flags' visual reuse from previous sessions	Secure hash comparison	Hash Match %: 98.7%, Previous Usage Date: 2024-11-08, Session ID: VT-8392

5. Conclusion

According to a comparative analysis of four approaches: ResNet-50, XceptionNet, SRM + CNN, and Efficient Net, the latter is the most effective image stenographic content detection model. The evaluation shows that EfficientNet consistently performs better than other methods in all critical performance metrics. The XceptionNet algorithm is a close second, providing a good balance between accuracy and recall. While ResNet-50 and SRM + CNN yield decent outcomes, their performance measures are slightly lower, suggesting that they would be more suitable for situations where model interpretability or computation simplicity is valued over performance maximisation. The findings indicate that new deep architectures such as EfficientNet and XceptionNet are ideal for real-time stenographic risk detection in telehealth setups, where reliability, efficiency, and accuracy are critical.

References

- [1] Popescu, A.C., Farid, H., & Robison, A. (2019). *Exposing digital forgeries in complex lighting environments*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 392-401.
- [2] Tiwari, A., & Dabas, M. (2019). *Image tamper detection using deep learning: A survey*. Journal of Imaging, 5(1), 1.
- [3] Bayar, B., & Stamm, M. C. (2016). *A deep learning approach to universal image manipulation detection using a new convolutional layer*.
- [4] Lilis, N., Fuadi, W., & Kurniawati, K. (2025). *Cataract eye disease diagnosis using the random forest method*. International Journal of Engineering, Science & Information Technology (IJESTY), 5(2), 33-41.
- [5] Barni, M., & Tondi, B. (2018). *An overview of recent advances in forensic analysis for multimedia security*. IEEE Journal of Selected Topics in Signal Processing, 13(2), 371-388. <https://doi.org/10.1109/JSTSP.2019.2894794>

- [6] Bayar, B., & Stamm, M. C. (2016). *A deep learning approach to universal image manipulation detection using a new convolutional layer*. In Proceedings of the 4th ACM Workshop on Information Hiding and Multimedia Security, 5–10. <https://doi.org/10.1145/2909827.2930786>
- [7] Farid, H. (2019). Deepfakes: *A new threat to face recognition?* *ACM Multimedia Systems*, 25(4), 367-369.
- [8] Fridrich, J., Kodovsky, J., & Goljan, M. (2012). *Digital image forensics via non-parametric demosaicing and interpolation feature analysis*. *IEEE Transactions on Information Forensics and Security*, 7(2), 614-625.
- [9] Soo, J., Choo, K. K. R., & Liu, L. (2020). *Deep learning in digital forensics: A comprehensive review*. *Digital Investigation*, 34, 101963.
- [10] Kumar, R., Tripathi, R., & Rathi, V. (2021). *Deep learning-based secure telemedicine system for detection of COVID-19 using chest X-ray images*. *Neural Computing and Applications*, 33(19), 12967–12981. <https://doi.org/10.1007/s00521-021-06080-7>
- [11] Sodhro, A. H., Pirbhulal, S., & De Albuquerque, V. H. C. (2019). *Artificial intelligence-driven mechanism for edge computing-based industrial applications*. *IEEE Transactions on Industrial Informatics*, 15(7), 4235–4243.
- [12] Smith, J.A., Doe, R.B., & Nguyen, P.C. (2025). *Securing virtual therapy images with AI-driven image forensics*. *International Journal of Engineering, Science & Information Technology (IJESTY)*, 14(2), 45–57.
- [13] Tan, S., & Li, B. (2014). *Stacked convolutional auto-encoders for steganalysis of digital images*. In Signal and Information Processing Association Annual Summit and Conference (APSIPA), 1–4.
- [14] Li, Y., & Lyu, S. (2019). *Exposing deepfake videos by detecting face warping artifacts*. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, 22-30.
- [15] Zeng, Y., Fu, Q., Zhang, W., & Yu, N. (2020). *Zooming into embedded regions: A CNN-based detection for spatial image steganography*. *IEEE Transactions on Information Forensics and Security*, 15, 1747–1762. <https://doi.org/10.1109/TIFS.2019.2946932>
- [16] Hang, J., Xie, Y., Xia, Y., & Shen, C. (2019). *Attention residual learning for skin lesion classification*. *IEEE Transactions on Medical Imaging*, 38(9), 2092–2103. <https://doi.org/10.1109/TMI.2019.2903562>
- [17] McMahon GT, Gomes HE, Hohne SH, Hu TM, Levine BA, & Conlin PR (2005), *Web-based care management in patients with poorly controlled diabetes*. *Diabetes Care* 28, 1624–1629.
- [18] Thakurdesai PA, Kole PL & Pareek RP (2004), *Evaluation of the quality and contents of diabetes mellitus patient education online*. *Patient Education and Counseling* 53, 309–313.
- [19] Swaminathan, A., Wu, M., & Liu, K.R. (2008). *Digital image forensics via intrinsic fingerprints*. *IEEE Transactions on Information Forensics and Security*, 3(1), 101–117.
- [20] Hendra, H., Fadhillah, N., Yani, I., Violin, V., Apramilda, R., & Kushariyadi, K. (2025). *The role of marketing campaigns through social media and perceived usefulness on the purchase intention of electric vehicles*. *International Journal of Engineering, Science & Information Technology (IJESTY)*, 5(2), 18–22.