# Fragile Watermarking for Image Certification Using Deep Steganographic Embedding

Davide Ghiani[1], Jefferson David Rodriguez Chivata[1], Stefano Lilliu[1], Simone Maurizio La Cava[1],
Marco Micheletto[1], Giulia Orrú[1], Federico Lama[2], Gian Luca Marcialis[1]

[1]University of Cagliari, Piazza d'Armi I - 09123 Cagliari (Italy), e-mail:
{davide.ghiani,jeffersond.rodriguez,simonem.lac,marco.micheletto,giulia.orru,marcialis}@unica.it
[2]Dedem S.p.A., Via Cancelleria 59 - 00072 Ariccia (Italy), e-mail: federico.lama@dedem.it

*Abstract*—**Modern identity verification systems increasingly rely on facial images embedded in biometric documents such as electronic passports. To ensure global interoperability and security, these images must comply with strict standards defined by the International Civil Aviation Organization (ICAO), which specify acquisition, quality, and format requirements. However, once issued, these images may undergo unintentional degradations (e.g., compression, resizing) or malicious manipulations (e.g., morphing) and deceive facial recognition systems. In this study, we explore fragile watermarking, based on deep steganographic embedding as a proactive mechanism to certify the authenticity of ICAO-compliant facial images. By embedding a hidden image within the official photo at the time of issuance, we establish an integrity marker that becomes sensitive to any post-issuance modification. We assess how a range of image manipulations affects the recovered hidden image and show that degradation artifacts can serve as robust forensic cues. Furthermore, we propose a classification framework that analyzes the revealed content to detect and categorize the type of manipulation applied. Our experiments demonstrate high detection accuracy, including cross-method scenarios with multiple deep steganography-based models. These findings support the viability of fragile watermarking via steganographic embedding as a valuable tool for biometric document integrity verification.**

*Index Terms*—**watermarking, image certification, morphing**

## I. INTRODUCTION

Nowadays, facial recognition (FR) technology plays a crucial role in identity verification, serving as a fundamental component in border control, secure identity management, and forensic applications [1]. To enhance security and streamline passenger processing, many countries have adopted electronic passports (ePass), which store biometric data to enable accurate and automated identity verification at border checkpoints. According to the International Civil Aviation Organization (ICAO) guidelines, facial images stored in machine readable travel documents, such as ePass, must comply with strict biometric standards to facilitate reliable authentication [2].

However, multiple factors can affect the reliability of these systems, including compression artifacts, noise, or other distortions introduced during acquisition, transmission, or storage [3]. In addition to these non-malicious alterations, the increasing sophistication of image manipulation techniques has introduced new security vulnerabilities that can compromise the integrity of identity verification [4].

One of the most pressing threats is the morphing attack, which leverages image synthesis techniques to blend facial features of multiple subjects, generating realistic composite images that can be falsely considered to belong to different individuals [5]. This vulnerability is particularly critical in border security, where identity verification is based on comparing a live subject with the ePass photograph. If a morphed image is successfully enrolled in an ePass, both contributing individuals can authenticate using the same document, bypassing security checks [6], [7]. To mitigate the risks associated with morphing attacks, Morphing Attack Detection (MAD) techniques have been developed to differentiate between genuine and manipulated facial images [8]. Despite significant progress, existing MAD methods face challenges in terms of generalization across novel morphing techniques, and adaptability to real-world conditions of identity recognition scenarios [9], [10].

To address these limitations, researchers have investigated proactive mechanisms that embed verification signals within the image itself at acquisition time, ensuring integrity throughout the document lifecycle. A promising direction involves active authentication mechanisms that introduce integrity markers directly within an image to facilitate the detection of manipulation. Among such techniques, digital watermarking has been widely adopted in multimedia security [11]; however, watermarking is designed primarily for copyright protection and may lack the adaptability required for ePass applications, where facial images must remain unchanged after issuance. In parallel, deep learning-based steganography has recently enabled the embedding of large amounts of data into cover images with minimal visual distortion. Although steganography is not designed for integrity verification, its capacity and perceptual quality open opportunities for alternative use cases. In this work, we hypothesize that steganographic models can be repurposed to implement active integrity verification mechanisms for facial images. Specifically, we propose a fragile watermarking framework based on deep steganographic embedding, in which any manipulation of the host image degrades the embedded content, allowing post-hoc integrity verification through reconstruction of a known marker.

The concealed image can only be retrieved through a dedicated decoding process, which generates a revealed image. Since the hidden data is embedded within the cover image, any modification applied to such an image will inevitably affect

the revealed image, therefore introducing artifacts that could potentially provide a forensic indicator of tampering [12].

Despite its long history in multimedia applications, the use of steganography-inspired embedding methods for forensic integrity verification in ICAO-compliant biometric documents has, to the best of our knowledge, yet to be explored. This represents a novel direction in the fight against digital attacks such as morphing: we hypothesize that tampering leaves subtle but detectable fingerprints in the revealed marker.

In this regard, this study proposes a twofold contribution: (i) assessing the feasibility of fragile watermarking via deep steganographic embedding as a mechanism to verify the integrity of ICAO-compliant facial images, ensuring that any post-issuance modifications, whether intentional or accidental, can be detected; (ii) developing a classification model capable of distinguishing between different types of alterations, such as morphing, compression artifacts, and noise addition, based on how these transformations affect the retrieved marker.

The aim is not to develop a novel steganographic method, but rather to investigate whether standard steganographic models, when repurposed as fragile watermarking tools, inherently offer resilience properties applicable to biometric image security. By examining how various manipulations impact the fidelity of the revealed content, we explore the potential of this approach for unauthorized modification detection.

The rest of this paper is organized as follows. Section II reviews the current literature on steganography and watermarking for digital images. Section III describes the proposed approach. Section IV reports the experimental protocol employed to conduct our evaluation, while Section V reports the obtained results. Finally, conclusions are drawn in Section VI.

## II. RELATED WORK

Steganography and watermarking are both data hiding frameworks that can be classified according to their tolerance to image modifications [13]. While fragile approaches are highly sensitive and signal any alteration in the host image, robust methods aim to resist various manipulations, including synthetic content generation such as deepfakes [14]. In the context of ICAO-compliant facial verification, fragile mechanisms are preferable, as their universal sensitivity enables the detection of any post-acquisition modification. This aligns with the proactive nature of ICAO security requirements, where even minor unauthorized alterations must be detected to preserve the integrity of biometric documents. Following this, it is important to distinguish how the information is embedded and interpreted in the frameworks. Digital watermarking typically relies on compact bitstrings embedded at predefined locations or patterns. These are effective in robust scenarios like copyright protection [15] or tamper detection [16], but often lack interpretability when the watermarked image is degraded or tampered with. Conversely, modern steganographic models enable the embedding of richer information, such as entire images, directly into the visual structure of the host. Originally developed for covert messaging, these techniques can be repurposed for fragile watermarking. In particular,

any modification to the host image corrupts the embedded payload, and this degradation can be visually observed in the recovered content, providing actionable forensic indicators of tampering. This repurposing is made possible by recent advances in deep learning, which have significantly improved the imperceptibility and recoverability of hidden information. Architectures based on convolutional neural networks (CNNs) [17], generative adversarial networks (GANs) [18], and autoencoders [19] have demonstrated high-capacity and visually stable embedding and extraction, making them suitable candidates for integrity verification in fragile settings.

## III. PROPOSED APPROACH

The proposed approach introduces a fragile watermarking mechanism, implemented through deep steganographic embedding, to certify the authenticity of ICAO-compliant facial images, ensuring long-term verifiability of photographs used in official identity systems by embedding a hidden integrity marker at the time of issuance. If the stego-image remains unaltered, the embedded content can be retrieved without distortion. Conversely, any modification to the stego-image, whether due to standard image processing or deliberate tampering, inevitably affects the extracted content. Our key hypothesis is that different manipulations introduce systematic and detectable artifacts in the recovered marker, which can signal both the presence and nature of the transformation. If these degradation patterns are consistent, they can serve as a forensic cue for tamper detection and manipulation classification. To evaluate this hypothesis, we define a three-stage methodology (Figure 1):

1) A steganographic process $E$ is applied to the ICAO-compliant facial image (cover image) $I_C$ to embed a secret marker image $I_S$ within it, producing a certified watermarked image $I_{stego}$:

$$I_{stego} = E(I_C, I_S) \qquad (1)$$

The hidden image $I_S$, imperceptibly hidden within the cover image, acts as a fragile integrity marker, ensuring that any future modification to the stego-image affects the embedded content.

2) A set of controlled transformations $T$ are applied to $I_{stego}$ to simulate real-world manipulations:

$$I_t = T(I_{stego}) \qquad (2)$$

The considered manipulations—resizing, compression, noise, blur, sharpening, and morphing—reflect both common post-processing operations and intentional biometric attacks. These transformations simulate real-world conditions where an image might be altered after issuance, allowing us to assess whether the hidden integrity marker can act as a forensic signal.

3) A decoder $D$ is used to extract the revealed image $I_r$ from the potentially modified image $I_t$:

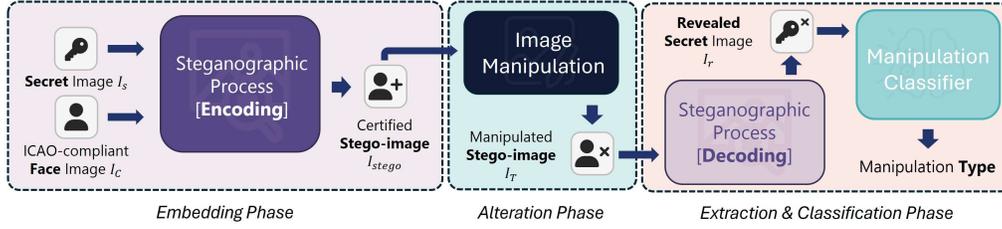$$I_r = D(I_t) \qquad (3)$$

Fig. 1: Overview of the proposed methodology, structured into three main phases: embedding, alteration and classification.



Fig. 2: Overview of the steganographic certification process and its impact on image quality: a) original input image; b) original secret image to be embedded; c) stego image generated using Steguz; d) secret image recovered from the Steguz stego image; e) stego image generated using Stegformer; f) secret image recovered from the Stegformer stego image.

If the stego-image remains unchanged (i.e., no transformation artifacts), the revealed image retains its expected structure. However, when modifications occur, artifacts emerge, reflecting the type and severity of the applied transformation. To systematically analyze these distortions, we introduce a classification model capable of distinguishing between different manipulation patterns.

The proposed system assumes that the integrity marker is present and matches the expected reference at verification time. In scenarios where the marker is missing, mismatched, or unrecoverable, the classification model may yield unreliable outputs. Handling such cases requires a separate detection stage, which is not addressed in the current work.

The next section presents the complete experimental pipeline, detailing the embedding architecture, transformation setup, classifier design, and evaluation criteria.

## IV. EXPERIMENTAL PROTOCOL

### A. Dataset

The goal is to assess the feasibility of our fragile watermarking approach on facial images compliant with ICAO guidelines. Accordingly, we selected the Chicago Face Database (CFD) [20]–[22], providing high-quality facial images of 827 men and women of varying ethnicity between the ages of 17-65. Specifically, it includes a single ICAO-compliant image per subject [23], with a $2444 \times 1718$ resolution.

To prepare the data for embedding, each image was cropped to a $1718 \times 1718$ square format, removing peripheral background while retaining the facial region. The cropped images were then resized to $224 \times 224$ pixels to comply with the input requirements of the steganography-based models used for analysis. As the integrity marker, we selected the ICAO logo, which was resized to $224 \times 224$ and embedded into each subject's facial image during the certification phase (Figure 2).

### B. Steganography Models

As discussed in Section II, recent deep learning-based steganography methods offer high-capacity and high-fidelity embedding mechanisms. Although originally developed for covert communication, these models can be repurposed for fragile watermarking, enabling integrity verification through the degradation of a hidden marker. To investigate the generality of our approach and support cross-model comparisons, we selected two representative state-of-the-art methods. The first, *Stegformer* [19], is a transformer-based autoencoder architecture designed for dense image-to-image embedding. The second, *SteGuz* [24], follows a more classical CNN-based design. Both are evaluated as embedding engines for our watermarking framework.

*Stegformer:* This model follows a U-Net-inspired architecture [25], where an autoencoder is used to encode a full image into a cover image and reconstruct it upon decoding. It includes a self-attention mechanism to enhance feature preservation and embedding quality (Figure 2e-f).

*SteGuz:* This model uses symmetry-aware CNNs to perform the embedding and recovery of a hidden image (Figure 2c-d). The architecture includes a preprocessing block, an encoder for information embedding, and a decoder for marker recovery. To ensure that the embedding process maintains the visual fidelity of the cover image, SteGuz introduces a custom loss function based on two image similarity metrics:

- **PSNR** (Peak Signal-to-Noise Ratio), which quantifies the ratio between the power of a signal and the power of corrupting noise, reflecting the fidelity between two compared images [26]. It is defined as:

$$\text{PSNR}(x,y) = 10 \cdot \log_{10}\left(\frac{MAX_x^2}{\text{MSE}(x,y)}\right) \quad (4)$$

where $x$ and $y$ are the original and reconstructed images respectively, both of size $m \times n$; $MAX_x$ is the maximum possible pixel value of the image (255 for 8-bit grayscale

images); and MSE is the Mean Squared Error between $x$ and $y$, defined as:

$$\text{MSE}(x,y) = \frac{1}{mn} \sum_{i=1}^{m} \sum_{j=1}^{n} [x(i,j) - y(i,j)]^2 \qquad (5)$$

In this formulation, $x(i,j)$ and $y(i,j)$ denote the pixel intensities at position $(i,j)$ in the images $x$ and $y$, respectively. The MSE measures the average squared difference between corresponding pixels, and PSNR expresses the result in decibel scale [27].

- **SSIM** (Structural Similarity Index Measure), which evaluates perceptual similarity between two images by comparing local patterns of pixel intensities normalized for luminance and contrast [28]. It is computed as:

$$\text{SSIM}(x,y) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \qquad (6)$$

where $\mu_x$, $\mu_y$ are the local means, $\sigma_x^2$, $\sigma_y^2$ are the variances, and $\sigma_{xy}$ is the covariance between the two images. Constants $C_1$ and $C_2$ are used to stabilize the division in case of weak denominators.

### C. Image Manipulations

To evaluate the resilience of the proposed integrity verification mechanism, we apply a series of controlled manipulations to the certified stego-image, simulating real-world modifications that may occur after issuance. These perturbations include both unintentional degradations (e.g., compression, resizing) and deliberate alterations (e.g., morphing, noise injection), allowing us to assess their impact on the embedded integrity marker. The exact parameters employed for each manipulation are reported in Table I.

*Compression:* Images may undergo re-encoding during digital storage or transmission, leading to quality degradation. Additionally, compression artifacts often emerge when images are processed through automated verification systems, where scanned photos are stored or analyzed in varying formats. To examine this effect, we apply JPEG and WebP compression, controlled by the quality factor $(Q_F)$, where $Q_F \in [80, 100]$. Lower values introduce moderate compression artifacts, while higher values result in minimal to no compression loss.

*Resizing:* ICAO-compliant images may be rescaled for different document formats, online submissions, or storage. To evaluate the resizing impact, each image is downscaled according to a resizing factor $(R_F)$, where $R_F \in [50\%, 99.9\%]$. Higher values result in minor resizing effects; lower values introduce severe downscaling, causing loss of detail. The image is then restored to its original dimension to observe potential degradation in the embedded marker.

*Noise Addition:* Low-bitrate encoding, repeated compression cycles, and scanning artifacts can introduce unwanted noise, affecting the overall integrity of an image. We simulate these effects using:

- *Gaussian noise*, where the standard deviation $(\sigma_G)$ is varied in the range $\sigma_G \in [2, 32]$. Low values introduce

minor pixel intensity variations, while high values lead to strong noise artifacts affecting fine details.

- *Salt-and-pepper noise*, parameterized by a corruption probability pair $(P_{SP})$, where $P_{SP} = (P_{\text{Salt}}, P_{\text{Pepper}})$. The first value represents the probability of a pixel turning white (Salt), and the second represents the probability of turning black (Pepper). Higher probabilities create more visible pixel corruption.

*Blurring:* Some post-processing techniques apply smoothing to remove noise or artifacts, which may also interfere with hidden data. We consider:

- *Gaussian blur*, which applies a weighted average of neighboring pixels, progressively diffusing fine details and potentially spreading steganographic patterns across a wider area. The kernel size $(K_G)$ is selected from $K_G \in \{3, 5, 7, 9\}$, with larger values causing stronger blurring.
- *Median blur*, which replaces each pixel with the median of its surrounding values, preserving edges better than Gaussian blur but disrupting embedded information by altering local pixel distributions. The kernel size $(K_M)$ is chosen from $K_M \in \{3, 5, 7, 9\}$, where higher values increase the filtering effect.

*Sharpening:* Certain document processing tools enhance image clarity by artificially increasing edge contrast. The sharpening intensity factor $(S_F)$ is adjusted within the range $S_F \in [0, 1]$. Low values produce no visible sharpening, while high values apply strong edge enhancement, which may introduce artificial artifacts.

*Morphing:* Unlike previous transformations, which may occur unintentionally, morphing is a deliberate biometric attack designed to deceive identity verification systems. In our study, we utilized *FaceMorpher*[1], an open-source tool based on facial landmarks to blend faces and create realistic morphed images. The blending factor $(\alpha_M)$ controls the degree of fusion between two source images; we set $\alpha_M = 0.9$ in our experiments, favoring the stego-identity while subtly incorporating features of the second. Therefore, while the outcome can be easily attributed to the most contributing individual, the second still has the chance to pass the identity verification [29]. This choice maximizes the attack success rate for the first individual while maintaining plausible deniability in border control scenarios thanks to an increased realism compared, for instance, to $\alpha_M = 0.5$, which may produce morphs too distant from either original biometric template.

### D. Image Quality And Manipulation Assessment

To assess the impact of image manipulations on the embedded integrity marker, we used three full-reference image quality metrics: Peak Signal-to-Noise Ratio (PSNR), Structural Similarity Index Measure (SSIM), and Mean Squared Error (MSE). The mathematical definitions of these metrics are provided in Section IV-B. Here, the focus is on their interpretation. PSNR and SSIM increase with image similarity,

---

[1]https://github.com/alyssaq/face_morpher

TABLE I: Summary of applied image manipulations and exact parameter values.

| Manipulation Class | Manipulation Type | Parameter | Values | # Samples |
|---|---|---|---|---|
| Compression | JPEG | Quality Factor ($Q_F$) | 100, 99, 90, 80 | 3308 |
| | WebP | Quality Factor ($Q_F$) | 100, 99, 90, 80 | 3308 |
| Resizing | Resizing | Scaling Factor ($R_F$) | 99.9%, 97.5%, 95%, 90%, 85%, 75%, 65%, 50% | 6616 |
| Gaussian Noise Addition | Gaussian Noise | Standard Deviation ($\sigma_G$) | 2, 4, 6, 8, 10, 16, 25, 32 | 6616 |
| Salt & Pepper Noise Addition | Salt & Pepper Noise | Corruption Probability ($P_{SP}$) | (0.01, 0.3), (0.03, 0.1), (0.1, 0.03), (0.3, 0.01), (0.01, 0.01), (0.03, 0.03), (0.1, 0.1), (0.3, 0.3) | 6616 |
| Blurring | Gaussian Blur | Kernel Size ($K_G$) | 3, 5, 7, 9 | 3308 |
| | Median Blur | Kernel Size ($K_M$) | 3, 5, 7, 9 | 3308 |
| Sharpening | Sharpening | Intensity Factor ($S_F$) | 0, 0.001, 0.01, 0.05, 0.1, 0.5, 0.75, 1 | 6616 |
| Morphing | FaceMorpher | Blending Factor ($\alpha_M$) | 0.9 | 6616 |

while MSE increases with distortion. An SSIM value close to 1 indicates high structural similarity; PSNR values above 50 dB typically correspond to minimal degradation. In contrast, higher MSE values reflect stronger pixel-wise differences. These metrics allow us to quantify how much the hidden marker is degraded after manipulation, and to evaluate whether an image has been altered (Figures 3 and 4).

*E. Classification protocol*

To train a model capable of identifying the type of manipulation applied to the image from the revealed secret, we employed ResNet-50 [30], pre-trained on ImageNet [31], as a backbone for feature extraction. We set the classification problem to seven classes: compression, resize, blur, gaussian noise, salt and pepper noise, sharpening and morph generation.

To perform the classification from the extracted embeddings, we concatenated a sequence of fully connected layers. First, a linear layer reduces the dimensionality from 2048 to 512 units. Next, a ReLU activation function is introduced to add non-linearity. Then, a dropout layer with an activation probability of 0.5 is inserted. Finally, a second linear layer reduces the output vector into the space of the target manipulation classes.

To assess the reliability of the classification model, we employed 70% of the user identities in the dataset (i.e., *578*) as the training set for fine-tuning and 30% as the test set (i.e., *249*). This allowed us to keep the manipulation classes balanced in the training and test sets and to simulate a real-world application context where the specific user face information is not known in training.

We evaluated the classification performance through metrics typically employed in pattern recognition: accuracy, precision, recall, and F1 Score. To assess the ability to generalize across different embedding models and under previously unseen manipulation conditions, we designed four evaluation protocols:

- *Intra-stega and intra-manipulation scenario:* both training and test images are embedded using the same model and subjected to the same types and strengths of manipulations.
- *Cross-stega scenario:* training and test images are embedded using different steganography-based models, while the manipulation types remain consistent.

It is important to highlight that while in a practical certification scenario, the embedding method would typically remain fixed throughout the system lifecycle, the cross-stega analysis can be useful to assess the generalizability and robustness of the proposed manipulation detection approach. In fact, it simulates potential real-world inconsistencies, such as re-certification with a different method or interoperability between systems using distinct embedding techniques. In addition, it provides information on the transferability of learned features between embedding strategies, which is essential for scalable or future-proof implementations.

Additionally, we applied two sub-protocols: 1) the first one, called P8-8, involves training and testing on the same eight variations of each manipulation type (e.g., levels of noise or compression); 2) the second, called P6-8, , involves training on six variations per manipulation type, while testing includes all eight—introducing two never-seen-before variations for each manipulation during testing.

## V. RESULTS

In this section, we present the results obtained from the previously described experiments. In Section V-A, we report the results obtained by analyzing the impact of the modifications in the images employed on the integrity verification process. In Section V-B, we discuss the capabilities of the classification model in detecting unauthorized modifications and distinguishing between different types of alterations.

*A. Image quality and manipulation assessment results*

To evaluate the feasibility of using secret image hiding and recovery for manipulation assessment, we analyzed the quality of the stego and recovered images before/after manipulations.

As shown in Table II, the average quality of the certified images and the revealed images, quantified by SSIM, MSE and PSNR, demonstrates that the distortion introduced by the certifying watermark is minimal and that the recovery processes produce high-fidelity outputs. This confirms that the embedding and extraction mechanisms, even when relying on different steganographic approaches (Steguz and Stegformer), do not introduce significant artifacts in both phases. For instance, the SSIM values remain above 0.92 for both models and stages, indicating that the certification and recovery
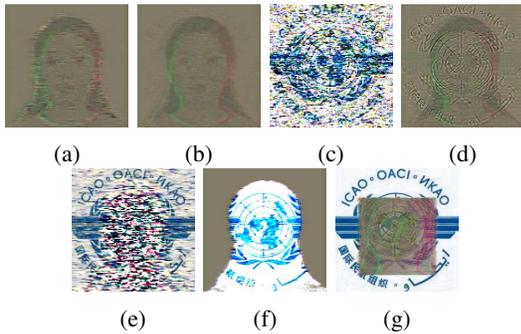
Fig. 3: Recovered secret images examples using Steguz after applying manipulations: a) JPEG compression ($Q_F = 80$), b) Gaussian blur ($K_G = 7$), c) Gaussian noise ($\sigma = 8$), d) resize ($R_F = 85\%$), e) salt & paper noise ($P_{SP} = (0.3, 0.01)$), f) sharpening ($S_F = 0.5$), g) morphing ($\alpha_M = 0.9$).
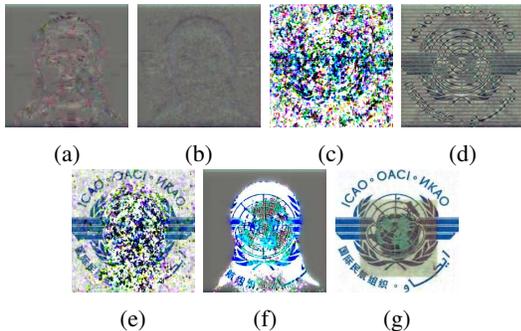


Fig. 4: Recovered secret images examples using Stegformer after applying manipulations: a) JPEG compression ($Q_F = 80$), b) Gaussian blur ($K_G = 7$), c) Gaussian noise ($\sigma = 8$), d) resize ($R_F = 85\%$), e) salt & paper noise ($P_{SP} = (0.3, 0.01)$), f) sharpening ($S_F = 0.5$), g) morphing ($\alpha_M = 0.9$).

processes do not noticeably corrupt the final images and are acceptable from a practical perspective. Comparing the two models, Stegformer provided better performance considering all the analyzed quality assessment metrics.

The impact of the manipulations on the secret image is, instead, appreciable and is shown in Figure 5. In all three metrics (SSIM, MSE, PSNR), manipulations such as Gaussian noise, morphing, and salt-and-pepper noise result in a substantial drop in quality. For example, the SSIM for morphing and salt-and-pepper noise drops below 0.5 in several cases, with a corresponding spike in MSE. These quality degradations are not only measurable, but also visually perceptible (see Figures 3 and 4), confirming that the manipulations leave strong and consistent fingerprints on the recovered secret image.

This behavior is the foundation for the proposed manipulation classifier. Despite being built on different design principles (Stegformer aiming for generalization, and Steguz for robustness) both models react similarly to post-embedding manipulations, suggesting that the introduced artifacts are sufficiently distinctive to be exploited for classification.

To support interpretation in practical deployments, we de-

TABLE II: Difference between original image and certified image (task certifying) and between original and recovered secret image (task recovery) in terms of SSIM, MSE and PSNR. Values are mean ± standard deviation.

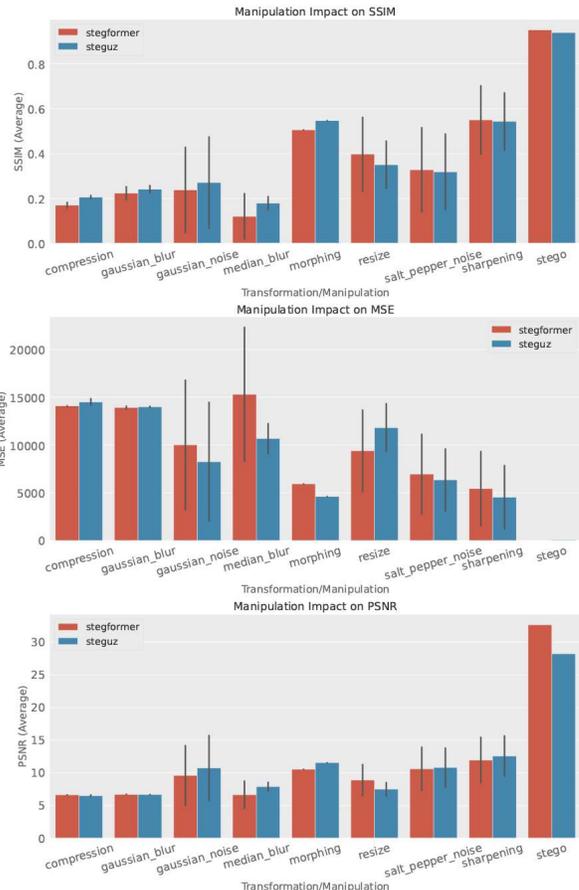| Task | Model | SSIM | MSE | PSNR |
|---|---|---|---|---|
| Certifyng | Steguz | 0.9266 ± 0.0097 | 123.27 ± 33.49 | 27.37 ± 1.14 |
| | Stegformer | 0.9696 ± 0.0035 | 7.23 ± 0.65 | 39.56 ± 0.40 |
| Recovery | Steguz | 0.9389 ± 0.0015 | 97.91 ± 3.51 | 28.23 ± 0.15 |
| | Stegformer | 0.9508 ± 0.0013 | 35.35 ± 1.81 | 32.65 ± 0.23 |



Fig. 5: Effect of transformations on image recovery. Error bars reflect variation within each manipulation type.

fine operational thresholds derived from the observed values in the unaltered (certified) case. Specifically, images with SSIM below 0.75 or PSNR below 22 dB are considered potentially manipulated.

### B. Classification Results

The classifier results reported in Table III show that manipulations can be detected by analyzing the recovered secret image. The intra-stega results for protocol P8-8 demonstrate that most manipulated samples ($\geq$99.95%) are correctly classified. However, the effects of manipulations vary depending on the embedding method used. In fact, the cross-stega results show an average performance drop of about 25%.

TABLE III: Classification performance in intra-stega and cross-stega scenarios (P8-8 sub-protocol).

| Training set | Test Set | Accuracy | Precision | Recall | F1 Score |
|---|---|---|---|---|---|
| Stegformer | Stegformer | 99.96% | 99.96% | 99.96% | 99.96% |
| | Steguz | 71.92% | 77.97% | 71.92% | 68.99% |
| Steguz | Steguz | 99.95% | 99.95% | 99.95% | 99.95% |
| | Stegformer | 80.51% | 86.16% | 80.51% | 77.35% |

A more detailed analysis through both the P8-8 and P6-8 sub-protocols highlights that the degradation in the generalization capability is strongly influenced by the type of manipulation applied, as shown by confusion matrices in Figure 6. In particular, the classification is still reliable on manipulation classes that generate well-defined and consistent structural and chromatic visual patterns, such as morphing, sharpening, and the addition of Gaussian noise. Most classification errors in cross-stega scenario involve the salt & pepper noise class being confused with Gaussian noise. Since these two types of manipulation are quite similar, such errors are potentially negligible in this application context.

Other manipulations often misclassified are resize and compression. Specifically, samples altered through resizing are, in some cases, misclassified as blurring. Regarding compression, the recovered images tend to be misclassified as either blurring or sharpening, depending on the method used.

In general, when comparing cross-stega P6-8 and P8-8 protocols, we observe that the performance drop primarily affects the model trained on images generated using Stegformer. In contrast, the model trained on Steguz maintains similar performance across both protocols. The latter is, therefore, able to generalize better on unknown variations.

In summary, the experimental results support the feasibility of using fragile watermarking, implemented via deep steganographic embedding, for image certification purposes. The proposed classification model was able to reliably identify the type of manipulation applied to the host image, even in challenging scenarios such as cross-stega settings and in the presence of unseen transformation variations. Despite the architectural and functional differences between the embedding models, the visual degradations produced on the revealed integrity marker were distinctive enough to enable generalization. This highlights the approach's potential for integration into real-world integrity verification pipelines, where robustness, interpretability, and scalability are critical requirements.

## VI. CONCLUSIONS

This work proposed a fragile watermarking framework for integrity verification of ICAO-compliant biometric images, based on deep steganographic embedding. A known visual marker is embedded into the facial image and later recovered to detect possible post-issuance manipulations through visible degradation. We evaluated this approach in the context of ICAO-compliant identity images, using two state-of-the-art steganography-based embedding models. A range of transformations, including compression, resizing, noise, and morphing, were applied to test the sensitivity and diagnostic
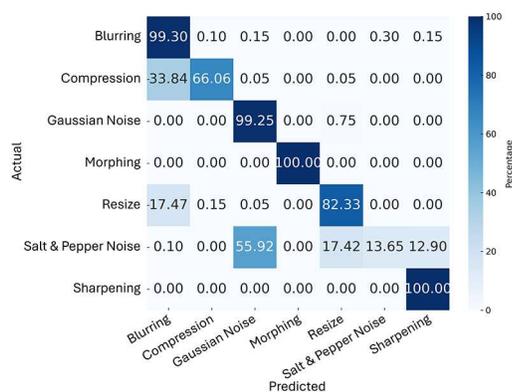
value of the revealed marker. Beyond tamper detection, we assessed the feasibility of classifying the type of manipulation by analyzing patterns of degradation in the recovered image. The findings demonstrate that steganography-based fragile watermarking can provide not only binary integrity verification but also actionable forensic information. To our knowledge, this is the first study to assess the use of standard deep steganographic models for this purpose in the context of biometric documents. To our knowledge, this is the first study to assess the use of standard deep steganographic models for this purpose in the context of document integrity. Future work will explore additional embedding architectures, extend the method to other biometric modalities, and evaluate robustness under adversarial conditions.

## REFERENCES

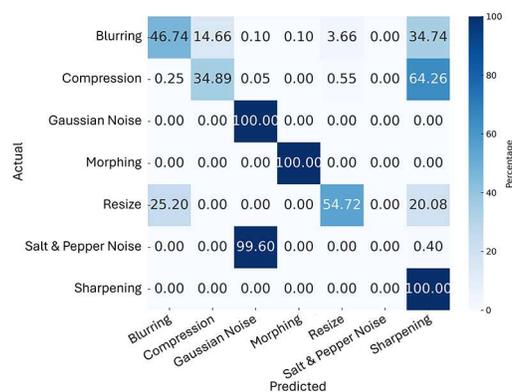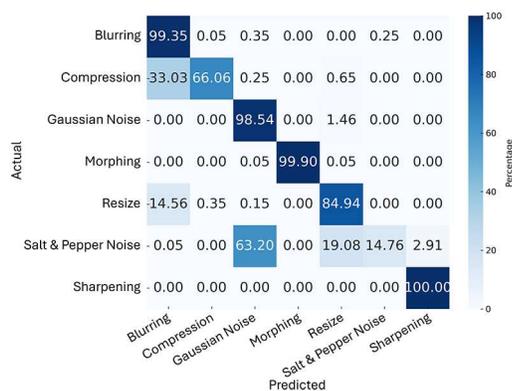[1] S. M. La Cava, G. Orrù, M. Drahansky, G. L. Marcialis, and F. Roli, "3d face reconstruction: the road to forensics," *ACM Computing Surveys*, vol. 56, no. 3, pp. 1–38, 2023.

[2] International Civil Aviation Organization (ICAO), *9303-machine readable travel documents-part 9: Deployment of biometric identification and electronic storage of data in emrtds*, ser. Doc 9303, Part 9, 2015, vol. 123.

[3] T. Schlett, C. Rathgeb, O. Henniger, J. Galbally, J. Fierrez, and C. Busch, "Face image quality assessment: A literature survey," *ACM Comput. Surv.*, vol. 54, no. 10s, Sep. 2022. [Online]. Available: https://doi.org/10.1145/3507901

[4] R. Tolosana, R. Vera-Rodriguez, J. Fierrez, A. Morales, and J. Ortega-Garcia, *An Introduction to Digital Face Manipulation*. Cham: Springer International Publishing, 2022, pp. 3–26. [Online]. Available: https://doi.org/10.1007/978-3-030-87664-7_1

[5] U. Scherhag, C. Rathgeb, J. Merkle, R. Breithaupt, and C. Busch, "Face recognition systems under morphing attacks: A survey," *IEEE Access*, vol. 7, pp. 23 012–23 026, 2019.

[6] D. J. Robertson, R. S. Kramer, and A. M. Burton, "Fraudulent id using face morphs: Experiments on human and automatic recognition," *PLoS One*, vol. 12, no. 3, p. e0173319, 2017.

[7] M. Ferrara, A. Franco, and D. Maltoni, "The magic passport," in *IEEE international joint conference on biometrics*. IEEE, 2014, pp. 1–7.

[8] M. Ngan, M. Ngan, P. Grother, K. Hanaoka, and J. Kuo, "Face recognition vendor test (frvt) part 4: Morph-performance of automated face morph detection," 2020.

[9] S. Venkatesh, R. Ramachandra, K. Raja, and C. Busch, "Face morphing attack generation and detection: A comprehensive survey," *IEEE transactions on technology and society*, vol. 2, no. 3, pp. 128–145, 2021.

[10] A. Panzino, S. M. La Cava, G. Orrù, and G. L. Marcialis, "Evaluating the integration of morph attack detection in automated face recognition systems," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 3827–3836.

[11] B. Song, P. Wei, S. Wu, Y. Lin, and W. Zhou, "A survey on deep-learning-based image steganography," *Expert Systems with Applications*, p. 124390, 2024.

[12] Y. Huang, Z. Liu, Q. Wu, and X. Liu, "Robust image steganography against jpeg compression based on dct residual modulation," *Signal Processing*, vol. 219, p. 109431, 2024.

[13] Z. Wang, O. Byrnes, H. Wang, R. Sun, C. Ma, H. Chen, Q. Wu, and M. Xue, "Data hiding with deep learning: A survey unifying digital watermarking and steganography," *IEEE Transactions on Computational Social Systems*, vol. 10, no. 6, pp. 2985–2999, 2023.
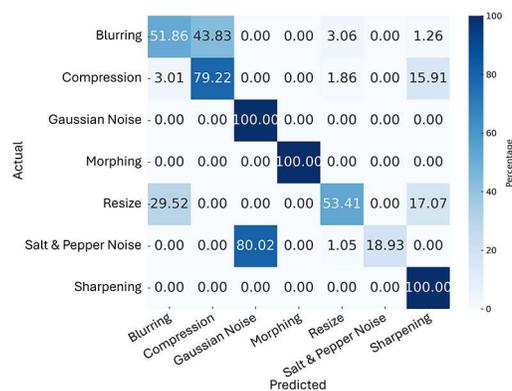
(a) P6-8 - Trained on Steguz, tested on Stegformer



(b) P6-8 - Trained on Stegformer, tested on Steguz



(c) P8-8 - Trained on Steguz, tested on Stegformer



(d) P8-8 - Trained on Stegformer, tested on Steguz

Fig. 6: Confusion matrices for Stegformer and Steguz models under cross-stega evaluation for P6-8 and P8-8.

[14] Y. Zhao, B. Liu, M. Ding, B. Liu, T. Zhu, and X. Yu, "Proactive deepfake defence via identity watermarking," in *Proceedings of the IEEE/CVF winter conference on applications of computer vision*, 2023, pp. 4602–4611.

[15] J. Wang, H. Wang, J. Zhang, H. Wu, X. Luo, and B. Ma, "Invisible adversarial watermarking: A novel security mechanism for enhancing copyright protection," *ACM Transactions on Multimedia Computing, Communications, and Applications*, vol. 21, no. 2, pp. 1–22, 2024.

[16] X. Zhang, R. Li, J. Yu, Y. Xu, W. Li, and J. Zhang, "Editguard: Versatile image watermarking for tamper localization and copyright protection," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 11 964–11 974.

[17] J. Zhu, R. Kaplan, J. Johnson, and L. Fei-Fei, "Hidden: Hiding data with deep networks," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 657–672.

[18] K. A. Zhang, A. Cuesta-Infante, L. Xu, and K. Veeramachaneni, "Steganogan: High capacity image steganography with gans," *arXiv preprint arXiv:1901.03892*, 2019.

[19] X. Ke, H. Wu, and W. Guo, "Stegformer: rebuilding the glory of autoencoder-based steganography," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 38, no. 3, 2024, pp. 2723–2731.

[20] D. S. Ma, J. Correll, and B. Wittenbrink, "The chicago face database: A free stimulus set of faces and norming data," *Behavior research methods*, vol. 47, pp. 1122–1135, 2015.

[21] D. S. Ma, J. Kantner, and B. Wittenbrink, "Chicago face database: Multiracial expansion," *Behavior Research Methods*, vol. 53, pp. 1289–1300, 2021.

[22] A. Lakshmi, B. Wittenbrink, J. Correll, and D. S. Ma, "The india face set: International and cultural boundaries impact face impressions and perceptions of category membership," *Frontiers in psychology*, vol. 12, p. 627678, 2021.

[23] C. Seibold, W. Samek, A. Hilsmann, and P. Eisert, "Accurate and robust neural networks for face morphing attack detection," *Journal of Information Security and applications*, vol. 53, p. 102526, 2020.

[24] A. Khalifa and A. Guzman, "Imperceptible image steganography using symmetry-adapted deep learning techniques," *Symmetry*, vol. 14, no. 7, p. 1325, 2022.

[25] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical image computing and computer-assisted intervention–MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III 18*. Springer, 2015, pp. 234–241.

[26] Q. Huynh-Thu and M. Ghanbari, "Scope of validity of psnr in image/video quality assessment," *Electronics letters*, vol. 44, no. 13, pp. 800–801, 2008.

[27] I. s. Avcıbaş, B. l. Sankur, and K. Sayood, "Statistical evaluation of image quality measures," *Journal of Electronic imaging*, vol. 11, no. 2, pp. 206–223, 2002.

[28] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE transactions on image processing*, vol. 13, no. 4, pp. 600–612, 2004.

[29] P. Korshunov and T. Ebrahimi, "Using face morphing to protect privacy," in *2013 10th IEEE international conference on advanced video and signal based surveillance*. IEEE, 2013, pp. 208–213.

[30] B. Koonce and B. Koonce, "Resnet 50," *Convolutional neural networks with swift for tensorflow: image recognition and dataset categorization*, pp. 63–72, 2021.

[31] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.

This figure "fig1.png" is available in "png" format from:

http://arxiv.org/ps/2504.13759v1