



ARTICLE

A Generative Image Steganography Based on Disentangled Attribute Feature Transformation and Invertible Mapping Rule

Xiang Zhang^{1,2,*}, Shenyan Han^{1,2}, Wenbin Huang^{1,2} and Daoyong Fu^{1,2}

¹School of Computer Science, Nanjing University of Information Science and Technology, Nanjing, 210044, China

²Engineering Research Center of Digital Forensics, Nanjing University of Information Science and Technology, Ministry of Education, Nanjing, 210044, China

*Corresponding Author: Xiang Zhang. Email: zhangxiang@nuist.edu.cn

Received: 12 November 2024; Accepted: 24 January 2025; Published: 26 March 2025

ABSTRACT: Generative image steganography is a technique that directly generates stego images from secret information. Unlike traditional methods, it theoretically resists steganalysis because there is no cover image. Currently, the existing generative image steganography methods generally have good steganography performance, but there is still potential room for enhancing both the quality of stego images and the accuracy of secret information extraction. Therefore, this paper proposes a generative image steganography algorithm based on attribute feature transformation and invertible mapping rule. Firstly, the reference image is disentangled by a content and an attribute encoder to obtain content features and attribute features, respectively. Then, a mean mapping rule is introduced to map the binary secret information into a noise vector, conforming to the distribution of attribute features. This noise vector is input into the generator to produce the attribute transformed stego image with the content feature of the reference image. Additionally, we design an adversarial loss, a reconstruction loss, and an image diversity loss to train the proposed model. Experimental results demonstrate that the stego images generated by the proposed method are of high quality, with an average extraction accuracy of 99.4% for the hidden information. Furthermore, since the stego image has a uniform distribution similar to the attribute-transformed image without secret information, it effectively resists both subjective and objective steganalysis.

KEYWORDS: Image information hiding; generative information hiding; disentangled attribute feature transformation; invertible mapping rule; steganalysis resistance

1 Introduction

Image information hiding is a technique for embedding secret information into an image, playing a crucial role in covert communication [1]. Existing methods can be categorized into cover modification-based methods and cover construction-based methods, depending on whether the cover image is modified [2]. Cover modification-based image information hiding typically alters the pixels or coefficients of the cover image to conceal secret information [3]. These modifications are imperceptible to the human eye but can be susceptible to detection by steganalysis [4]. To enhance resistance against steganalysis, researchers have developed cover construction-based image information hiding methods [5].

Within cover construction-based image information hiding, there are two primary types: mapping-based image information hiding and generative image steganography. Mapping-based image information hiding establishes a mapping relationship between the secret information and the image [6]. However, the



limited capacity of this method often requires multiple images to fully represent the secret information, significantly impacting its practicality. To address the low-capacity issue in mapping-based approaches, generative image steganography has been proposed, offering a solution for embedding larger amounts of secret information [7].

Generative image steganography embeds secret information during the image generation process. Early techniques used traditional image generation algorithms such as texture fusion [8–12], Turkish marbling [13,14], and fractals [15] to hide information. However, these methods generated secret images with simple content and texture. As deep learning continues to progress, new deep neural network architectures have been introduced, capable of generating a variety of high-quality images, which is compatible with the theory of generative image steganography. Consequently, researchers have utilized generative adversarial networks (GANs) [16–19], glow models [20,21], diffusion models [22,23], and autoencoders [24,25] to generate images and hide information. Compared to methods like texture fusion, these approaches generate more complex and diverse stego images, but they still require improvements in the quality of stego images and the accuracy of secret information extraction.

To improve imperceptibility and extraction accuracy of secret information, we propose a generative image steganography algorithm based on attribute feature transformation and invertible mapping rule. Firstly, we propose the image disentanglement to disentangle a reference image into content features and attribute features. Secondly, we propose a new invertible mapping rule to map the secret information into the new attribute feature. Finally, we substitute the original attribute feature with the new attribute feature to generate the stego image. The primary contributions of this paper are as follows:

- (1) Attribute feature transformation is implemented for generative image steganography. By separating content features and attribute features through image disentanglement and hiding secret information within attribute features, the algorithm realizes information hiding in the process of attribute feature transformation, generating high-quality stego images.
- (2) An invertible mapping rule for secret information and attribute features is proposed. The secret information is encoded into a noise vector consistent with the distribution of attribute features using the mean mapping rule, which is then input to the generator to obtain the stego image. Due to the fact that the noise vector is consistent with the attribute feature distribution, it effectively resists steganalysis detection.
- (3) Three types of loss functions are proposed to train the model: reconstruction loss, adversarial loss, and image diversity loss. These three losses ensure the extraction accuracy of secret information, the quality of the generated stego image, and its diversity, respectively. They also can help the model to converge quickly.
- (4) A generative image steganography framework based on attribute feature transformation and invertible mapping rule is proposed. With the image disentanglement and mapping rule, the proposed framework can achieve high extraction accuracy of secret information and steganalysis resistance. Experimental results demonstrate that the proposed algorithm improves the extraction accuracy and the quality of generated images compared to existing generative image steganography algorithms.

2 Related Work

Generative image steganography can be classified into two classes: one based on artificial design and the other based on deep learning.

(1) Generative Image Steganography Based on Artificial Design

Initially, researchers developed artificially designed generative algorithms to generate stegos. Inspired by texture fusion, Otori et al. [8] proposed a method to hide secret information during texture synthesis. However, it has a limited hiding capacity. Wu et al. [9] introduced a block-based image information hiding method. This algorithm sorts candidate blocks by mean square error; it improves capacity but reduces imperceptibility. Xu et al. [10] proposed an image steganography that deforms an image containing secret information to generate a stego image. However, the secret information must be an image, which limits the hiding ability of this method. Zhou et al. [11] proposed an image steganography algorithm based on seed region growing and least significant bit (LSB). This algorithm uses a seed region growing algorithm to determine the order of texture synthesis, improving imperceptibility but lacking robustness. Wei et al. [12] addressed the robustness issue by proposing a texture synthesis steganography algorithm based on super-pixel structure and SVM. Lee et al. [13] developed a generative image steganography algorithm based on pattern synthesis. However, the resulting image texture is simple, and its imperceptibility is limited. Li et al. [14] proposed a method based on fingerprint image construction. This method maps secret information into the detailed information of the fingerprint during the synthesis process, which can improve robustness and imperceptibility. However, its capacity remains limited. Zhang et al. [15] introduced a generative image steganography algorithm based on fractal theory. This algorithm hides secret information in the generation process of recursion and escape time, but it has deficiencies in the extraction accuracy of secret information.

In summary, early generative image information hiding methods based on artificial design generally suffer from poor quality of stego images, insufficient imperceptibility, and low capacity due to the limitations of image generation algorithms and information hiding methods.

(2) Generative Image Steganography Based on Deep Learning

With the swift advancement of deep learning, image generation techniques have made significant breakthroughs. Deep neural networks can generate highly realistic images. Thus, generative image steganography based on deep learning has emerged as a prominent research topic. Early on, scholars primarily used GANs for image construction and secret information hiding. For instance, Duan et al. [16] were the first to propose a generative image steganography method using GAN. It improves the quality, imperceptibility, and capacity compared with artificially designed methods. However, the information extraction relies on pairing with a generator, resulting in poor generalization. Hu et al. [17] introduced a generative image steganography method based on Deep Convolutional GAN (DCGAN). While this algorithm has better generalization compared to [16]. However, it suffers from poor image quality and insufficient extraction accuracy. Wei et al. [18] introduced a generative steganography network that consists of a generator, discriminator, steganalyzer, and extractor. Since the generation processes for the cover image and the stego image are identical, it can effectively resist steganalysis. However, as the hiding capacity increases, the extraction accuracy of secret information significantly decreases. Peng et al. [19] proposed a method combining GAN with gradient descent approximation. Although this algorithm has high capacity and extraction accuracy, it suffers from high complexity due to the updating of the noise vector. In [20,21], flow model based generative image steganography methods have been proposed. Similarly, some diffusion model based methods have been developed [22,23]. They typically map secret information into noise vectors, inputting them into the diffusion model to obtain the stego image. These approaches generally achieve high image generation quality and extraction accuracy, but at the cost of high complexity. Zhou et al. [24] proposed a generative image steganography method that relies on semantic object contours. This method inputs secret information into a long short-term memory (LSTM) network to generate contour lines, which are then synthesized with a GAN to produce the stego image. However, the visual quality of the generated images and the extraction accuracy of the secret information need further improvement.

Some scholars have explored hiding information through intermediate features of the image generation process. For example, Liu et al. [25] proposed a generative image steganography based on image disentanglement. Sun et al. [26] proposed a method based on guidance feature. These methods have high robustness, but the generated images are relatively homogeneous, and the extraction accuracy of secret information can still be further improved. We summarize special characteristics of some typical generative image steganography methods in Table 1.

Table 1: Summary of some typical generative image steganography methods

Category	Refs.	Years	Implement details	Advantages	Disadvantages
Artificial design based methods	[8]	2007	Hiding information by texture synthesis	Low complexity	Low capacity
	[9]	2015	Hiding information in overlapping region	Higher capacity	Low imperceptibility
	[14]	2019	Fingerprint image construction and image synthesis	Higher robustness and imperceptibility	Low capacity
	[15]	2020	Embedding by fractal generation process	High capacity and robustness	Low imperceptibility
Deep learning based methods	[18]	2022	Mapping secret information to noise and generate stego by GAN	High steganalysis resistance	Insufficient extraction accuracy
	[21]	2022	Generate stego image by flow model	High imperceptibility	High complexity
	[23]	2023	Using diffusion model to embed secret information	High capacity and robustness	High complexity
	[25]	2022	Embed secret in structure feature and generate style transferred stego	High capacity	Insufficient imperceptibility and accuracy
	[26]	2023	Embed secret in style feature and generate style transferred stego	High robustness	Insufficient extraction accuracy

In summary, due to the constraints of image generation mechanisms and the limitations of methods for embedding secret information, most existing generative image information hiding algorithms struggle with issues related to the quality of stego images and the accuracy of secret information extraction. To address

these issues, this paper proposes a generative image information hiding framework based on attribute feature transformation and invertible mapping rule.

3 Preliminaries

3.1 Generative Image Steganography Theory

Image steganography and hiding detection improve performance through a dynamic interplay. On one hand, scholars develop novel image steganography methods, which optimize the quality of the stego images to resist detection. On the other hand, hiding detection methods extract more precise features to distinguish between the stego image and the cover image. Let C represents the cover image, C' represents the stego image, and S represents the secret information image. Conventional image information hiding can be expressed by Eq. (1):

$$C' = EMD(C, S) \quad (1)$$

where $EMD(\cdot)$ denotes the information hiding method, which modifies the cover image C to get the stego image C' based on the secret information S . While the hiding detection algorithm is used to determine whether an image contains secret information. Assuming that the hiding detection method is $De(\cdot)$, X is the input image to be detected, and the classification probability of the image $De(X)$ can be obtained by the information hiding detection method:

$$De(X) = [p_0, p_1]^T \quad (2)$$

where p_0 and p_1 represent the probabilities that the image to be detected is a cover and a stego, respectively. The final discrimination result $F(De(X))$ is represented by the following equation:

$$F(De(X)) = \begin{cases} cover & \text{if } p_0 > p_1 \\ stego & \text{if } p_1 \leq p_0 \end{cases} \quad (3)$$

Usually, we measure the performance of information hiding detection through P_E , and its formula is as follows:

$$P_E = (P_{FA} + P_{MD})/2 \quad (4)$$

where P_{FA} denotes the false alarm rate and P_{MD} denotes the miss detection rate. From the above analysis, traditional image information hiding methods, which modify the cover image to obtain the stego image, lead to minor variations in the distributions of the cover and stego images. These differences can be easily detected by state-of-the-art hiding detection algorithms. Therefore, effective resistance to detection can theoretically be achieved only when no cover image exists or when the distributions of the cover image and the stego image are identical. The following equation expresses generative image steganography:

$$C' = HID(S) \quad (5)$$

where $HID(\cdot)$ represents a generative information hiding algorithm. From Eq. (5), it can be noticed that there is no cover image in generative image steganography. Instead, it generates the stego image C' directly from the secret information. Consequently, the information hiding detection method $DE(\cdot)$ cannot recognize the stego image. However, directly generating a stego image from the secret information is a big challenge. A compromise approach is using the intermediate features that are generated according to the secret information and then synthesizing the stego image by these intermediate features. Inspired by attribute

transformation, we propose to map secret information to attribute features and then generate stego by the attribute feature.

3.2 Information Hiding Based on Disentangled Feature

It is mentioned in [24] that assuming there exists a data space ε that can be encoded to a group of different feature spaces $\{f^1, f^2, \dots, f^n\}$. For a given feature space f^i , if it always satisfies:

$$P(f^i) = P(f^i|f^j) \quad (6)$$

where $1 \leq j \leq n$ and $j \neq i$, the features in this space can be considered independent of those in other spaces and are referred to as disentangled features. Conversely, features that are dependent are called entangled features. However, extracting entangled features from the stego image is challenging. In existing generative information hiding methods based on attribute transformation, the stego attribute feature and the reference image are typically input together into the generator to produce the attribute-transformed stego image, as illustrated in Fig. 1.

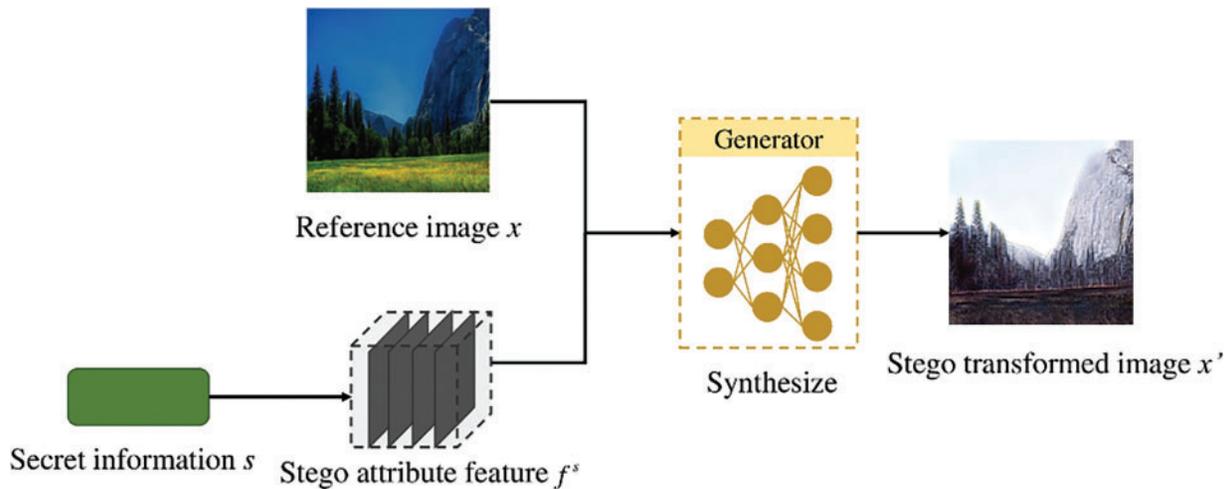


Figure 1: Existing generative image steganography framework

However, since the reference image itself contains attribute features, it means that the feature space in this information hiding system is entangled. Therefore, accurately recovering the attribute feature from the stego image is challenging. To address this, we disentangle the reference image x into the attribute feature and the content feature. The stego attribute feature and the content feature are then input into the generator to get the stego image, as shown in Fig. 2. As illustrated in the figure, unlike existing generative image steganography methods based on attribute transformation, this paper disentangles the image into two independent features and encodes the secret information into one of these features. It significantly enhances the extraction accuracy of secret information.

4 The Proposed Generative Image Steganography Method

In this paper, we construct disentanglement features during attribute transformation and use these features to hide secret information and generate stego images. Attribute feature transformation is a technology that uses deep neural networks to convert image attributes. It is related to image-to-image translation. The specific framework is shown in Fig. 3.

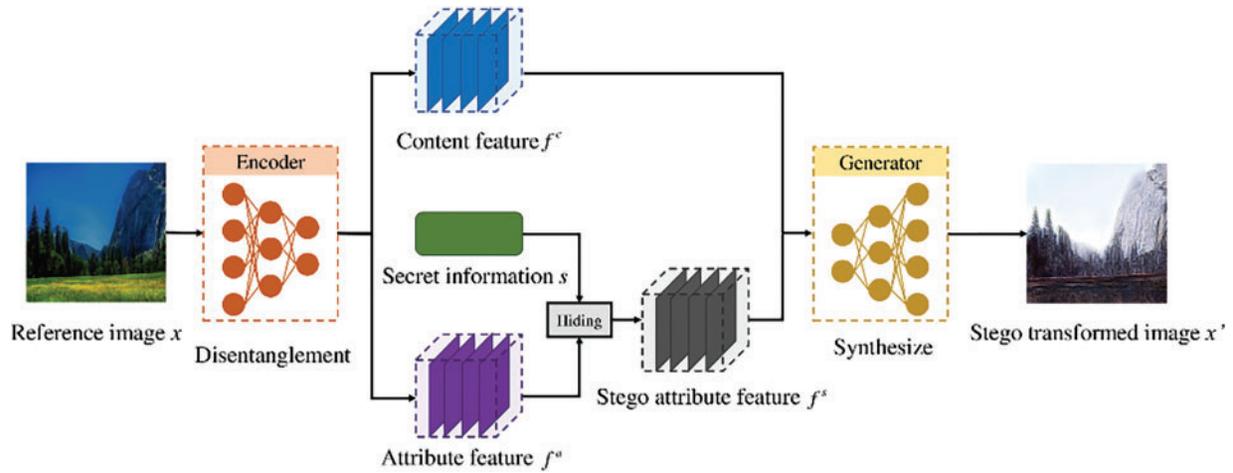


Figure 2: The proposed generative image steganography framework

On the sender side, the reference image x is first disentangled into independent content features f^c by the content encoder $E^c(\cdot)$. Meanwhile, we map the secret information m into the stego attribute feature f^s by the proposed mean mapping rule $map(\cdot)$. The stego attribute feature f^s and the content feature f^c are input to the generator $G(\cdot)$ to synthesize the transformed stego image x' . Simultaneously, we randomly generate a set of attribute feature f^{a2} , which are also input into the generator $G(\cdot)$ together with the content features f^c to generate another transformed image x'_2 . This process is used to construct the diversity loss to heighten the diversity of the generated images. On the receiver side, the transformed stego image x' is disentangled into the stego attribute features \bar{f}^s and content features f^c by the attribute encoder $E^a(\cdot)$ and content encoder $E^c(\cdot)$, respectively. The secret information m' is recovered by the inverse mapping rule $demap(\cdot)$ from the stego attribute feature \bar{f}^s . Finally, a discriminator $D(\cdot)$ is used to construct adversarial loss. As depicted in the figure, the main components of the proposed method include the invertible mean mapping rule, the loss function, and the structure of the encoder and generator.

4.1 Invertible Mean Mapping Rule

Invertible Mean Mapping Rule is a mechanism that maps secret information to a noise vector that satisfies a specific distribution. Meanwhile, the secret information can be extracted from the noise vector. It consists of a mean mapping rule and an inverse mapping rule.

(1) Mean Mapping Rule $map(\cdot)$

To guarantee the accuracy of secret information extraction and the imperceptibility of the stego image, we propose a mean mapping rule. This rule maps the secret information into noise vectors that are consistent with the distribution of attribute features. The detailed process is as follows:

Step 1. Divide the secret message into n segments; the length of each segment is l :

$$m = \{m_1, m_2, m_3, \dots, m_n\} \tag{7}$$

Step 2. Convert all the segments from binary to decimal and map it between $-d$ and d with the following formula:

$$f_k^s = \left(\frac{2 \times ten(m_k) + 1}{2^l} - 1 \right) \times d \tag{8}$$

where m_k is the k th segment, f_k^s is the k th converted and mapped segment, and $ten(\cdot)$ represents a function for converting binary to decimal.

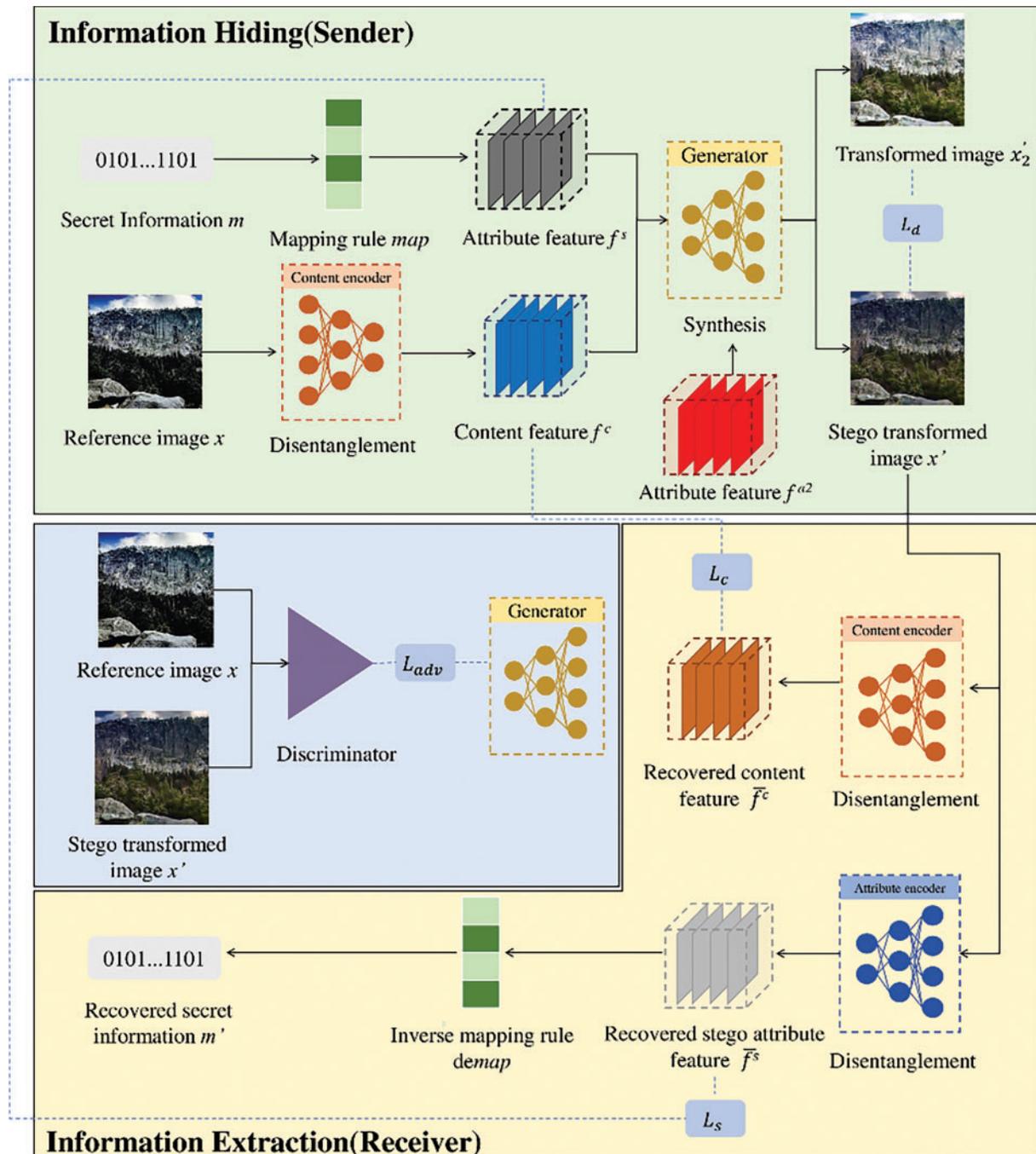


Figure 3: The overall structure of the proposed generative image steganography method

Step 3. By combining all the converted and mapped segments together in order, the stego attribute feature f^s can be obtained:

$$f^s = \{f_1^s, f_2^s, f_3^s, \dots, f_n^s\} \tag{9}$$

When $l = 3, d = 2$, an example of the proposed mapping rule is presented in Table 2.

Table 2: The mapping rule when $l = 3, d = 2$

m	000	001	010	011	100	101	110	111
f	-1.75	-1.25	-0.75	-0.25	0.25	0.75	1.25	1.75

(2) Inverse Mapping Rules $demap(\cdot)$

To accurately extract the secret information hidden in the stego attribute feature f^s , an inverse mapping rule based on the mean mapping rule is designed. The specific procedure is as follows:

Step 1. Extract each subvector from the stego attribute feature $\overline{f^s}$:

$$\overline{f^s} = \{\overline{f_1^s}, \overline{f_2^s}, \overline{f_3^s}, \dots, \overline{f_n^s}\} \tag{10}$$

Step 2. Each subvector is converted into a binary number of length l using the following formula:

$$m'_k = bin\left(\frac{(\overline{f_k^s} + d) \times 2^l - d}{2d}\right) \tag{11}$$

where f_k^s is the k th subvector, m'_k is the k th recovered secret information segment, and $bin(\cdot)$ represents a function for converting decimal to binary.

Step 3. By combining m'_k together in order, the hidden secret information m' can be recovered as:

$$m' = \{m'_1, m'_2, m'_3, \dots, m'_n\} \tag{12}$$

When $l = 3, d = 2$, an example of the inverse mapping rule is presented in Table 3.

Table 3: The inverse mapping rule when $l = 3, d = 2$

f	-2~-1.5	-1.5~-1	-1~-0.5	-0.5~0	0~0.5	0.5~1	1~1.5	1.5~2
m	000	001	010	011	100	101	110	111

The invertible mapping rule is closely related to the extraction accuracy of secret information. As l in Eq. (8) increases, the capacity grows, but the extraction accuracy of the secret information gradually decreases.

4.2 Loss Function

In order to guarantee the preciseness of secret information extraction and the quality of the stego images, we propose three types of loss functions to train the model: adversarial loss L_{adv} , reconstruction losses L_c and L_s , image diversity loss L_d .

(1) Adversarial Loss

To improve the quality of the generated transformed stego images, we first define the adversarial loss L_{adv} between the generator and discriminator according to generative adversarial theory [16], which is shown as follows:

$$L_{adv} = \mathbb{E} [\log D(x)] + \mathbb{E} [\log (1 - D(G(f^c, f^s)))] \quad (13)$$

where $D(x)$ represents the discrimination results of the reference image x . $D(G(f^c, f^s))$ represents the discrimination results of the stego image generated by the generator $G(\cdot)$.

(2) Reconstruction Loss

To guarantee the content of the transformed stego image remains unchanged, we use the content feature reconstruction loss L_c :

$$L_c = \|E^c(x) - E^c(x')\|_1 \quad (14)$$

where $\|\cdot\|_1$ represents L_1 Loss. $E^c(x)$ and $E^c(x')$ represent the content features extracted from the reference image x and the stego image x' . Meanwhile, since attribute features are mapped from secret information, the recovery rate of attribute features will directly affect the extraction accuracy of secret information. In order to enhance the extraction accuracy of secret information, we further design the attribute feature reconstruction loss L_s :

$$L_s = \mathbb{E} [\|f^s - E^a(x')\|_1] \quad (15)$$

where $E^a(x')$ represents the attribute features extracted from the stego image x' .

(3) Image Diversity Loss

We propose an image diversity loss through f^s and f^{a2} to enrich the visual effects of the generated transformed stego images:

$$L_d = \mathbb{E} \left[1 / \frac{\|G(f^c, f^s) - G(f^c, f^{a2})\|_1}{\|f^s - f^{a2}\|_1} \right] \quad (16)$$

where $G(f^c, f^s)$ and $G(f^c, f^{a2})$ represent the two transformed images with different attribute features generated by the generator $G(\cdot)$.

Finally, the overall loss L of the proposed model can be obtained through three types of losses:

$$L = \lambda_1 L_{adv} + \lambda_2 L_c + \lambda_3 L_s + \lambda_4 L_d \quad (17)$$

where $\lambda_1, \lambda_2, \lambda_3$ and λ_4 are the weights of each loss function.

4.3 The Structure of Encoder and Generator

(1) Content Encoder

Fig. 4 illustrates the content encoder structure of the proposed model, which primarily consists of ReflectionPad2d layers, Conv2d layers, InstanceNorm2d layers, and ReLU layers.

(2) Attribute Encoder

Fig. 5 shows the attribute encoder structure of the proposed model, which primarily consists of ReflectionPad2d layers, Conv2d layers, ReLU layers, and AdaptiveAvgPool2d layers.

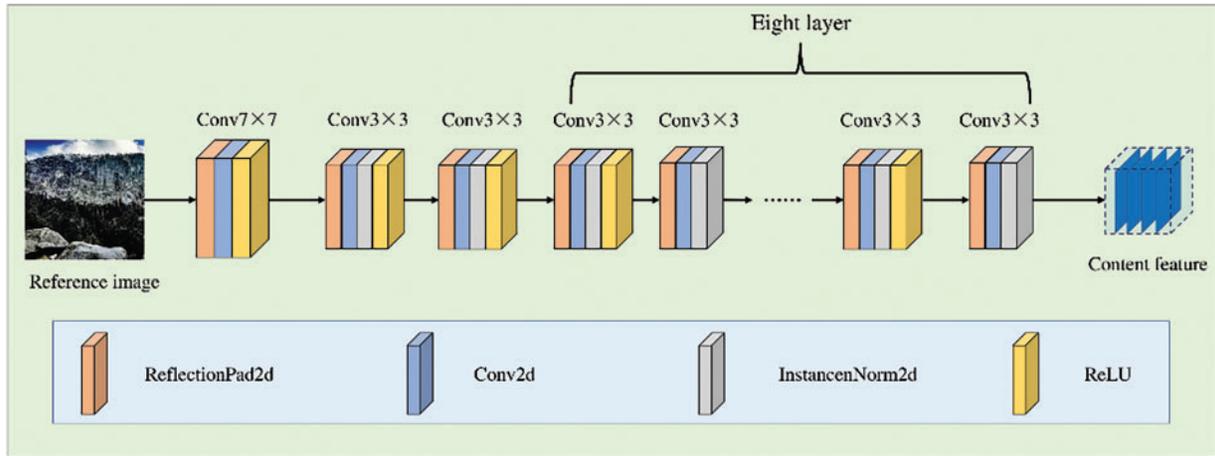


Figure 4: Content encoder structure

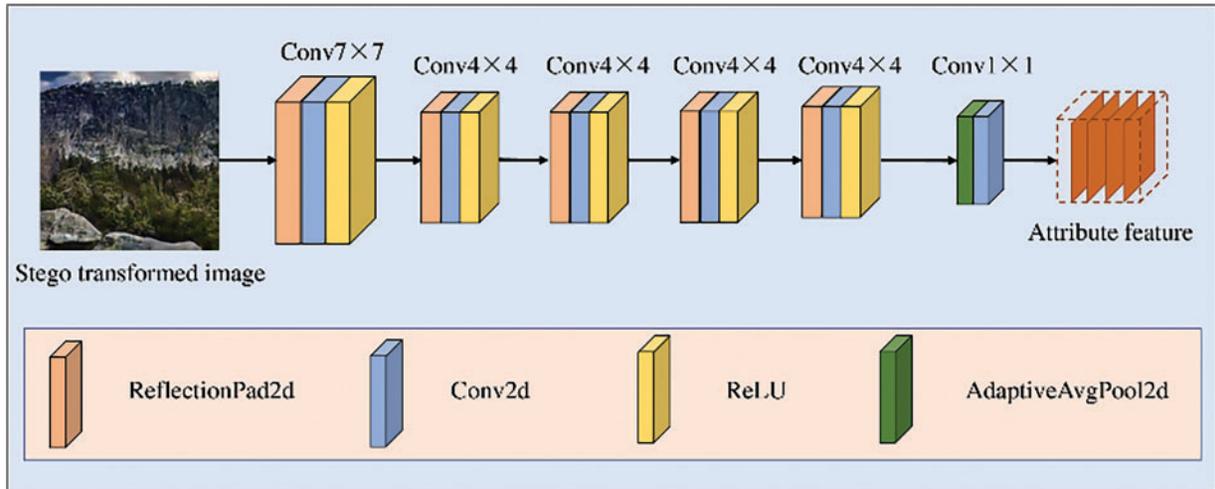


Figure 5: Attribute encoder structure

(3) Generator

Fig. 6 shows the generator encoder structure of the proposed model, which is mainly composed of ReflectionPad2d layers, Conv2d layers, InstanceNorm2d layers, ReLU layers and ConvTranspose2d layers. Where \oplus represents the connection operation between two tensors. The attribute features are connected with the content features every two layers, and finally a transformed stego image is generated.

4.4 Information Hiding and Extraction Process

(1) Information Hiding Stage

The detailed process of information hiding is as follows:

Step 1. The sender inputs the reference image x into the content encoder $E^c(\cdot)$ and disentangles it into the content feature f^c :

$$\{f^c\} = \{E^c(x)\} \tag{18}$$

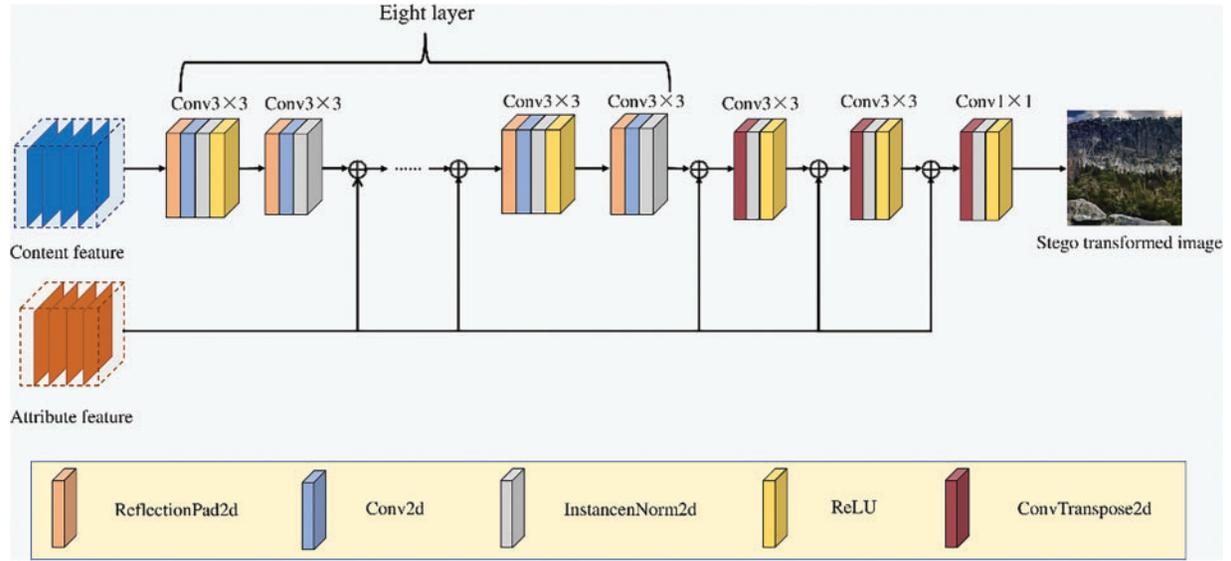


Figure 6: Generator structure

Step 2. The secret information m is mapped to the stego attribute feature f^s using the mean mapping rule $map(\cdot)$:

$$f^s = map(m) \quad (19)$$

Step 3. The stego attribute feature f^s is input into the generator $G(\cdot)$ along with the content feature f^c to generate the transformed stego image x' through attribute feature transformation.

$$x' = G(f^c, f^s) = G(E^c(x), map(m)) \quad (20)$$

Through the above steps, secret information can be hidden during the attribute feature transformation process, which can get a high-quality transformed stego image.

(2) Information Extraction Stage

The specific process of information extraction is as follows:

Step 1. After receiving the transformed stego image x' , the receiver inputs it into the attribute encoder $E^a(\cdot)$ and disentangles it to obtain the stego attribute feature \bar{f}^s :

$$\bar{f}^s = E^a(x') = E^a(G(f^c, f^s)) \quad (21)$$

Step 2. The secret information m' is reconstructed from \bar{f}^s by the inverse mapping rule $demap(\cdot)$:

$$m' = demap(\bar{f}^s) = demap(E^a(x')) \quad (22)$$

Through the above steps, secret information can be extracted from the transformed stego image.

5 Experiment Results and Analysis

5.1 Experiment Environment and Parameter Settings

The experiments were conducted on a server equipped with an Intel(R) Core(TM) i9-12900KF CPU at 3.19 GHz, an NVidia GTX3090Ti GPU, 32 GB of RAM, and the Windows 10 operating system. The testing environment includes Python 3.6.13 and Pytorch 1.10.2. The learning rate lr is initially set to 0.0001, the training weights λ_1 , λ_2 , λ_3 and λ_4 of the loss function are 0.1, 0.1, 10 and 0.1, respectively, and the L_1 regularization loss function is used. The parameter d in Eq. (8) takes the value of 2, and l is set as 1, 2, 3, and 4 when capacity is 16, 32, 48, and 64, respectively. The network is trained for 14,000 epochs, and the network parameters are trained by the Adam optimization algorithm. The training ratio of the generator to the discriminator is 1:1. The datasets are Yosemite and Landscape Pictures [27] and the batch size is set to 16. We analyze the proposed model from hidden capacity, extraction accuracy, image quality, and steganalysis.

5.2 Analysis of Hiding Capacity and Extraction Accuracy

Since the hiding capacity of generative image steganography differs from that of traditional image information hiding. We use bits per image (bpi) to measure the hiding capacity. The formula for bpi is shown in Eq. (23):

$$bpi = \frac{bits}{NS} \quad (23)$$

where NS denotes the number of tested stego images and bits represents the total number of bits hiding in NS stego images. Additionally, we use acc to measure the extraction accuracy of secret information with different capacities, which is shown in Eq. (24):

$$acc = \frac{\sum_i^{NS} \frac{nr_i}{bpi}}{NS} \times 100\% \quad (24)$$

where nr_i represents the number of bits of secret information correctly extracted from the i th stego image.

In the experiments, the extraction accuracy with capacities of 16, 32, 48, and 64 bpi on Yosemite and Landscape Pictures datasets is tested, respectively. The results are presented in Table 4.

Table 4: Hiding capacity and extraction accuracy

Dataset	Hiding capacity (bpi)	Extraction accuracy (acc)
Yosemite	16	100%
	32	99.9%
	48	99.7%
	64	99.4%
Landscape pictures	16	100%
	32	99.9%
	48	99.8%
	64	99.3%

From the table, it can be found that due to the independence of the feature disentanglement and the robustness of the proposed mean mapping rule, a high extraction accuracy of secret information can be obtained under different hiding capacities. When bpi is 16, the extraction accuracy of our model can reach

100% on both datasets. When bpi gradually increases, the extraction accuracy is slightly reduced. The reason is that the proposed model hides the secret information in float tensor data and then uses the generator to convert the data into an integer image. The entire conversion process has truncation loss. However, it's worth noting that when bpi increases to 64, the extraction accuracy of secret information still remains above 99%.

5.3 Analysis of Objective Visual Quality

Objective visual quality of the generated stego image is an important metric for evaluating the security of generative information hiding. However, since the stego image generated by the proposed algorithm is an attribute-transformed image, traditional image quality indicators such as SSIM and PSNR are not suitable to measure its performance. Therefore, we use three no-reference image quality assessment metrics: NIQE [28], PI (Perceptual Index) [29], and BRISQUE [30]. Among them, NIQE relies on the statistical characteristics of natural images to evaluate the quality of images and is particularly suitable for general distortion types. BRISQUE uses the spatial statistical characteristics of images for reference-free quality evaluation. PI is based on a perceptual model and comprehensively considers a variety of low-level visual features to provide a quality evaluation that is more in line with human perception. Lower values in these metrics indicate higher quality of the generated images. In the experiment, we analyze the objective visual quality of the generated stego images under different hiding capacities on the Yosemite and Landscape Pictures datasets. The reference image (Reference) represents the original image in the dataset. Table 5 displays the specific results.

Table 5: Objective visual quality

Dataset	Hiding capacity (bpi)	NIQE (\downarrow)	PI (\downarrow)	BRISQUE (\downarrow)
Yosemite	Reference	5.5	3.4	14.1
	16	4.8	3.2	14.6
	32	4.7	3.2	14.5
	48	4.8	3.3	14.5
	64	4.7	3.2	14.4
	Landscape pictures	Reference	7.4	4.6
16		5.0	3.9	9.7
32		5.1	3.9	9.4
48		5.1	3.8	9.6
64		5.0	3.8	9.6

Based on the results in Table 5, it is clear that for the Yosemite dataset, the NIQE and PI of the transformed stego images generated by our method are lower than those of the reference images, while the BRISQUE metrics are slightly higher. For the Landscape Pictures dataset, all three metrics of the transformed stego images are lower than those of the reference images. It indicates that the objective visual quality of images generated by our algorithm is high. This is because our method hides the secret information only in the attribute feature without interfering with the content feature. Meanwhile, the design of the image diversity and adversarial loss also contributes to the high quality of the generated stego images.

5.4 Analysis of Subjective Visual Quality

To further assess the visual quality of the stego images generated by our algorithms, the analysis of subjective visual quality is given. Fig. 7 shows some examples of transformed stego images generated on the Yosemite and Landscape Pictures datasets when the bpi is 64. As presented in Fig. 7, the transformed stego images generated by the proposed algorithm have excellent quality. Meanwhile, the visual effect is obviously different from that of the reference image. Despite these differences, the image's content information remains unchanged. It indicates that the proposed algorithm successfully achieves the disentanglement of attribute and content features. This ensures that the attribute feature can be modified without altering the content, resulting in a high-quality transformed stego image.

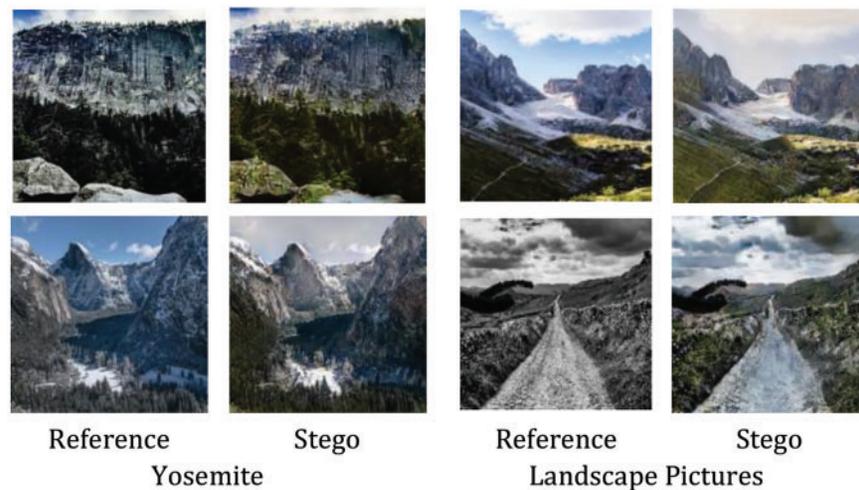


Figure 7: Example of transformed stego images

Additionally, we conducted two sets of subjective visual experiments: (1) the effect of hiding different secret information with the same reference images. (2) the effect of hiding the same secret information with different reference images. Figs. 8 and 9 present the obtained results. In Fig. 8, the first column displays two reference images, while the second to fourth columns show the transformed stego images generated by different secret information. It can be realized that the visual effects of the transformed stego images produced from the same reference image vary with different secret information, primarily in terms of color and shades. In Fig. 9, the first column presents different reference images, and the second column shows the transformed stego images generated by the same secret information. The results indicate that the visual effects of the transformed stego images generated by the same secret information have consistent style attributes.

In summary, the transformed stego images generated by our algorithm have excellent performance in terms of both subjective and objective visual quality. This further indicates that the algorithm can effectively resist subjective steganography analysis, and it is difficult to detect whether the generated images contain secret information through human vision.

5.5 Analysis of Steganalysis Resistance

Steganalysis resistance is also performed on the Yosemite and Landscape Pictures datasets for transformed stego images generated with the capacities of 16, 32, 48, and 64, respectively. Meanwhile, three state-of-the-art steganalysis networks XuNet [31], SRNet [32], and SiaStegNet [33] are regarded as steganalyzers to measure the effectiveness of the proposed algorithm. Since the method proposed in this paper

belongs to generative image steganography without the cover image, we define the positive samples of the steganalyzer as the attribute-images without hiding secret information and the negative samples as the attribute transformed stego images. Finally, we use P_E , whose formula appears in Eq. (4), to measure the steganalysis resistance performance. When the value of P_E approaches 0.5, it indicates excellent steganalysis resistance performance. Table 6 presents the experimental results.

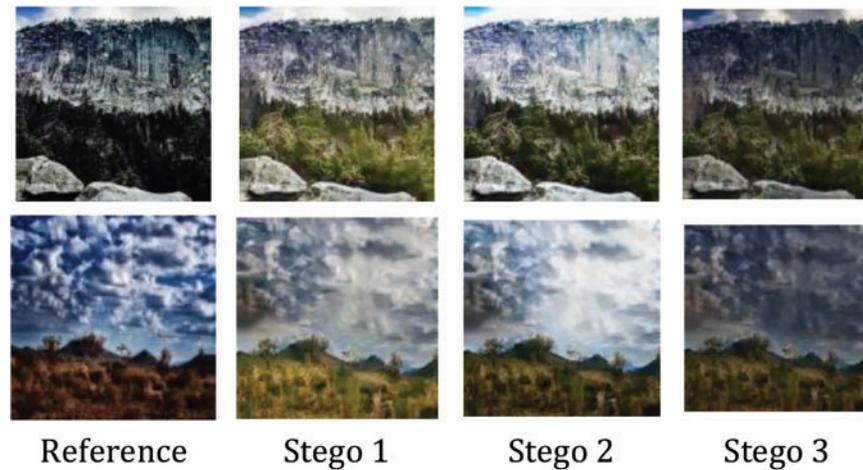


Figure 8: The visual effect of hiding different secret information with the same reference images

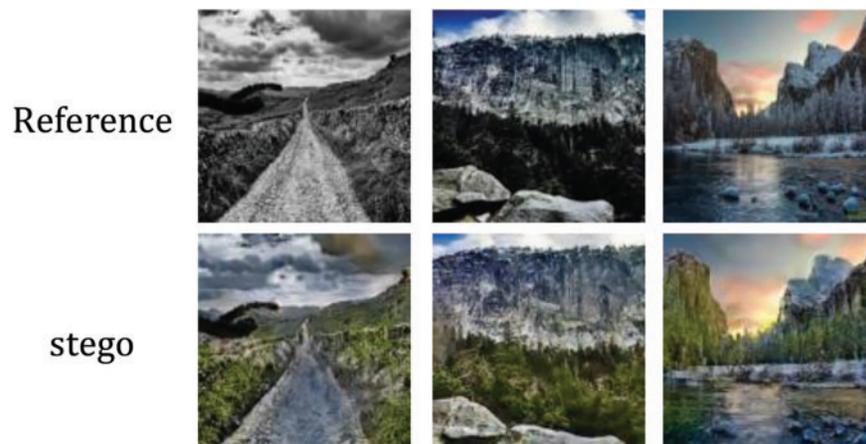


Figure 9: The visual effect of hiding the same secret information with different reference images

Table 6: The results of steganalysis resistance

Dataset	Hiding capacity (<i>bpi</i>)	Steganalysis	P_E ($\rightarrow 0.5$)
	16	XuNet	0.49
		SRNet	0.50
	32	SiaStegNet	0.53
		XuNet	0.51

(Continued)

Table 6 (continued)

Dataset	Hiding capacity (<i>bpi</i>)	Steganalysis	P_E ($\rightarrow 0.5$)
Yosemite	48	SRNet	0.51
		SiaStegNet	0.49
		XuNet	0.50
		SRNet	0.49
	64	SiaStegNet	0.48
		XuNet	0.51
		SRNet	0.49
		SiaStegNet	0.51
Landscape pictures	16	XuNet	0.52
		SRNet	0.50
		SiaStegNet	0.49
	32	XuNet	0.53
		SRNet	0.49
		SiaStegNet	0.51
	48	XuNet	0.52
		SRNet	0.49
		SiaStegNet	0.51
	64	XuNet	0.48
		SRNet	0.51
		SiaStegNet	0.52

As shown in Table 6, it is clear that the proposed algorithm outperforms steganalysis across various datasets and hiding capacities. The values of P_E are maintained at about 0.5. The reason is that we utilize the uniformly distributed noise vector as a stego attribute feature, ensuring that the distribution of the transformed stego images aligns with that of ordinary transformed images. Consequently, the proposed algorithm could effectively resist steganalysis.

5.6 Analysis of Computational Complexity

In deep neural network models, computational complexity is a key factor that shows the application effectiveness of the model. Therefore, we enumerate our model's computational and training complexity from various perspectives, and Table 7 displays the results.

Table 7: Computational and training complexity

Index	Result
Total parameters	27.06 M
FLOPs	778.29 G
Input size	2.30 MB
Forward/backward pass size	7297.22 MB
Parameters size	103.24 MB

(Continued)

Table 7 (continued)

Index	Result
Estimated total size	7402.76 MB
Training time	83.33 epoch/h

As shown in Table 7, the FLOPs and total parameters of the proposed model fall within an acceptable range. The time complexity and space complexity are not particularly high.

5.7 Analysis of Diversity Loss

To demonstrate the influence of image diversity loss L_d , we retrained our model without image diversity loss L_d , and the visual quality of the generated stego image, both with and without image diversity loss, is displayed in Table 8.

Table 8: The visual quality of the generated stego image with and without image diversity loss

	NIQE (↓)	PI (↓)	BRISQUE (↓)
With L_d	4.7	3.2	14.4
Without L_d	4.8	3.2	14.6

It is found from the table that image diversity loss L_d does indeed affect the visual quality of the generated stego image. With L_d , the three metrics NIQE, PI, and BRISQUE all obtain the best results compared with no L_d .

5.8 Analysis of Comparison

To further demonstrate the superiority of the proposed algorithms in terms of extraction accuracy of secret information, visual quality of stego images, and steganalysis resistance, we conducted comparative experiments against RoSteALS [3], IDEAS [25], and RoSteGFS [26] on the Yosemite and Landscape Pictures datasets. Among them, RoSteALS is a traditional image steganography algorithm designed to achieve high steganography security, and it contains a simple network structure with few parameters. While IDEAS maps the secret information into structure feature, therefore it can obtain larger capacity. RoSteGFS maps the secret information into guidance feature without consideration of image disentanglement. In the experiments, the hiding capacity for each algorithm was set to 64 *bpi*. The best results are highlighted in bold, and the less optimal ones are underlined. The detailed comparison outcomes are provided below.

(1) Comparison of Extraction Accuracy

Table 9 presents a comparison of the extraction accuracy of each algorithm on the two datasets with a capacity of 64 *bpi*. It is evident that the proposed algorithm surpasses all others in terms of secret information extraction accuracy. Specifically, the proposed algorithm achieves extraction accuracies of 99.4% and 99.3% on the Yosemite and Landscape Pictures datasets, respectively. This superior performance is attributed to the proposed feature disentanglement and invertible mapping rule, which enable the proposed algorithm to achieve high extraction accuracy.

Table 9: The comparison results of the extraction accuracy of secret information

Dataset	Model	Extraction accuracy (<i>acc</i>)
Yosemite	RoSteALS [3]	99.0%
	IDEAS [25]	<u>99.1%</u>
	RoSteGFS [26]	98.5%
	Ours	99.4%
Landscape pictures	RoSteALS [3]	99.1%
	IDEAS [25]	99.0%
	RoSteGFS [26]	98.5%
	Ours	99.3%

Compared with the proposed method, the extraction accuracy of secret information in the other three methods is slightly inferior. Firstly, RoSteGFS is an entanglement method; the stego image is generated by inputting a reference image and a stego guidance feature. Therefore, the guidance feature of the reference image will affect the extraction of the stego guidance feature, which results in a decrease in extraction accuracy of secret information. Secondly, the hiding information method in IDEAS maps the secret information to structure features. However, the structure feature seems more difficult to recover, and the mapping rule in IDEAS are also not robust. Finally, RoSteALS is a traditional image information hiding algorithm. Its information hiding network structure is simple and easy to train, but its special iterative extraction method results in insufficient extraction accuracy of secret information. Overall, the secret information extraction accuracy of the proposed method outperforms the other three methods.

(2) Comparison of Objective Visual Quality

In the experiments, we assess the objective visual quality of the four algorithms on the Yosemite and Landscape Pictures datasets with a hidden capacity of 64 *bpi*. We utilize the three no-reference image quality assessment metrics described in Section 5.3 for objective visual comparisons. The results are displayed in Table 10. From this table, it is evident that the proposed algorithm outperforms IDEAS and RoSteALS in terms of objective visual quality on both datasets. On the Landscape Pictures dataset, the proposed algorithm achieves the best NIQE and BRISQUE values, and the PI is slightly lower than RoSteGFS. On the Yosemite dataset, the proposed algorithm performs comparably to RoSteGFS, which results from the well-designed network structure, mapping rule, and disentangled model. While in IDEAS, secret information is mapped to structure features, changing structure features leads to insufficient quality of stego images. In RoSteALS, secret information is directly hidden in the cover image. The modification of information hiding results in lower quality compared to our model.

Table 10: The comparison results of objective visual quality

Dataset	Model	NIQE (↓)	PI (↓)	BRISQUE (↓)
Yosemite	Reference	5.5	3.4	14.1
	RoSteALS [3]	5.3	3.4	13.2
	IDEAS [25]	5.0	3.4	13.7
	RoSteGFS [26]	4.5	2.8	14.6
	Ours	<u>4.7</u>	<u>3.2</u>	<u>14.4</u>

(Continued)

Table 10 (continued)

Dataset	Model	NIQE (\downarrow)	PI (\downarrow)	BRISQUE (\downarrow)
Landscape pictures	Reference	7.4	4.6	22.8
	RoSteALS [3]	<u>6.3</u>	4.2	22.1
	IDEAS [25]	6.6	4.1	27.2
	RoSteGFS [26]	5.0	3.2	<u>17.9</u>
	Ours	5.0	<u>3.8</u>	9.6

(3) Comparison of Subjective Visual Quality

Since IDEAS alters the structure feature of the image, the content of the stego is fundamentally different from that produced by other models. Meanwhile, RoSteALS is a traditional image information hiding algorithm. Therefore, it is also not suitable for subjective visual comparison with the proposed model. Therefore, we only compare the subjective visual quality of the proposed model with RoSteGFS. The results are shown in Fig. 10.



Figure 10: Comparison results of subjective visual quality. RoSteGFS is from [26]

From Fig. 10, it can be observed that compared with RoSteGFS, the proposed model exhibits greater visual differences in the generated stego images, particularly due to the proposed image diversity loss. Additionally, the proposed model introduces less noise in smoothing texture regions and better visual quality since the proposed model uses disentanglement features to hide information. In summary, the proposed model demonstrates superior performance in both subjective and objective visual quality compared to existing generative image steganography algorithms.

(4) Comparison of Steganalysis Resistance

In this experiment, we compare the four algorithms on the Yosemite and Landscape Pictures datasets for steganalysis resistance with a hidden capacity of 64 bpi . P_E is used to measure the steganalysis resistance,

and the results are presented in Table 11. By comparing the results on both datasets, it can be seen that the four algorithms exhibit consistent performance. Among them, RoSteGFS, IDEAS, and the proposed algorithm are generative image steganography methods; their P_E values are around 0.5 due to the absence of cover images. This indicates that the proposed method can completely resist steganalysis. In contrast, RoSteALS is a traditional image information hiding method; the P_E is less than 0.1, demonstrating its poor steganalysis resistance.

Table 11: The comparison results of steganalysis resistance

Dataset	Method	Steganalysis	P_E ($\rightarrow 0.5$)
Yosemite	RoSteALS [3]	XuNet	0.07
		SRNet	0.02
		SiaStegNet	0.01
	IDEAS [25]	XuNet	0.52
		SRNet	0.52
		SiaStegNet	0.49
	RoSteGFS [26]	XuNet	0.51
		SRNet	0.50
		SiaStegNet	0.48
	Ours	XuNet	0.51
		SRNet	0.49
		SiaStegNet	0.51
Landscape pictures	RoSteALS [3]	XuNet	0.03
		SRNet	0.01
		SiaStegNet	0.01
	IDEAS [25]	XuNet	0.51
		SRNet	0.48
		SiaStegNet	0.49
	RoSteGFS [26]	XuNet	0.51
		SRNet	0.49
		SiaStegNet	0.49
	Ours	XuNet	0.48
		SRNet	0.51
		SiaStegNet	0.52

6 Conclusion

In this paper, we propose a generative image steganography model based on attribute feature transformation and invertible mapping rule. After analyzing existing generative image steganography algorithms, we reveal significant issues such as insufficient imperceptibility and robustness when embedding information into non-disentangled features. To address these issues, we encode the secret information into a noise vector that conforms to the attribute feature distribution by feature disentanglement and mean mapping rule. This noise vector is then fused with the content features of a reference image to produce the stego image. It can be further demonstrated through experiment results that the proposed feature disentanglement and invertible mapping rule achieve higher extraction accuracy and superior image quality compared to existing generative image steganography methods. However, the proposed method also has limitations, such as the balance

between capacity and extraction accuracy, and the generalization of other types of datasets. Therefore, our future work will focus on its generalization and larger capacity.

Acknowledgement: The authors would like to thank the anonymous reviewers for their kind comments and suggestions for improving the paper and thank their colleague Dr. Main Uddin for his contribution on the improvement of English writing.

Funding Statement: This work was supported in part by the National Natural Science Foundation of China (Nos. 62202234, 62401270), the China Postdoctoral Science Foundation (No. 2023M741778), the Natural Science Foundation of Jiangsu Province (Nos. BK20240706, BK20240694).

Author Contributions: The authors confirm contribution to the paper as follows: Conceptualization, Xiang Zhang; methodology, Xiang Zhang and Shenyan Han; software, Xiang Zhang and Shenyan Han; validation, Xiang Zhang, Shenyan Han and Wenbin Huang; formal analysis, Daoyong Fu; investigation, Xiang Zhang and Shenyan Han; data curation, Wenbin Huang and Daoyong Fu; writing—original draft preparation, Xiang Zhang and Shenyan Han; writing—review and editing, Wenbin Huang and Daoyong Fu; visualization, Wenbin Huang; supervision, Wenbin Huang and Daoyong Fu; project administration, Daoyong Fu; funding acquisition, Xiang Zhang. All authors reviewed the results and approved the final version of the manuscript.

Availability of Data and Materials: Not applicable.

Ethics Approval: Not applicable.

Conflicts of Interest: The authors declare no conflicts of interest to report regarding the present study.

References

1. Wang F, Zhang X, Fu Z. An iterative two-stage probability adjustment strategy with progressive incremental searching for image steganography. *IEEE Trans Circ Syst Video Technol.* 2024;34(10):9428–44. doi:10.1109/TCSVT.2024.3398222.
2. Zhang X, Peng F, Long M. Robust coverless image steganography based on DCT and LDA topic classification. *IEEE Trans Multimed.* 2018;20(12):3223–38. doi:10.1109/TMM.2018.2838334.
3. Bui T, Agarwal S, Yu N, Collomosse J, RoSteALS: robust steganography using autoencoder latent space. In: 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW); 2023 Jun 17–24; Vancouver, BC, Canada.
4. Meng L, Jiang X, Sun T. A review of coverless steganography. *Neurocomputing.* 2024;566:126945. doi:10.1016/j.neucom.2023.126945.
5. Liu Q, Xiang X, Qin J, Tan Y, Zhang Q. A robust coverless steganography scheme using camouflage image. *IEEE Trans Circ Syst Video Technol.* 2022;32(6):4038–51. doi:10.1109/TCSVT.2021.3108772.
6. Liu Q, Xiang X, Qin J, Tan Y, Tan J, Luo Y. Coverless steganography based on image retrieval of DenseNet features and DWT sequence mapping. *Knowl Based Syst.* 2020;192:105375. doi:10.1016/j.knsys.2019.105375.
7. Chen X, Zhang Z, Qiu A, Xia Z, Xiong NN. Novel coverless steganography method based on image selection and StarGAN. *IEEE Trans Netw Sci Eng.* 2022;9(1):219–30. doi:10.1109/TNSE.2020.3041529.
8. Otori H, Kuriyama S. Data-embeddable texture synthesis. In: The 8th International Symposium on Smart Graphics (SG); 2007 Jun 25–27; Kyoto, Japan.
9. Wu KC, Wang CM. Steganography using reversible texture synthesis. *IEEE Trans Image Process.* 2014;24(1):130–9. doi:10.1109/TIP.2014.2371246.
10. Xu J, Mao X, Jin X, Jaffer A, Lu S, Li L, et al. Hidden message in a deformation-based texture. *Vis Comput.* 2015;31(12):1653–69. doi:10.1007/s00371-014-1045-z.
11. Zhou Q, Qiu Y, Li L, Lu J, Yuan W, Feng X, et al. Steganography using reversible texture synthesis based on seeded region growing and LSB. *Comput Mater Contin.* 2018;55(1):151–63. doi:10.3970/cmc.2018.055.151.

12. Wei WY, Wang LZ, Ma HF. A texture synthesis steganography scheme based on super-pixel structure and SVM. In: The 10th IFIP TC 12 International Conference on Intelligent Information Processing (IIP); 2018 Oct 19–22; Nanning, China.
13. Lee WK, Ong S, Wong K, Tanaka K. A novel coverless information hiding technique using pattern image synthesis. In: 2018 Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC); 2018 Nov 12–15; Honolulu, HI, USA.
14. Li S, Zhang X. Toward construction-based data hiding: from secrets to fingerprint images. *IEEE Trans Image Process.* 2019;28(3):1482–97. doi:10.1109/TIP.2018.2878290.
15. Zhang X, Peng F, Lin Z, Long M. A coverless image information hiding algorithm based on fractal theory. *Int J Bifurcation Chaos.* 2020;30(4):2050062. doi:10.1142/s0218127420500625.
16. Duan X, Song H. Coverless information hiding based on generative model. arXiv:1802.035282018. 2018. doi: 10.48550/arXiv.1802.03528.
17. Hu D, Wang L, Jiang W, Zheng S, Li B. A novel image steganography method via deep convolutional generative adversarial networks. *IEEE Access.* 2018;6:38303–14. doi:10.1109/ACCESS.2018.2852771.
18. Wei P, Li S, Zhang X, Luo G, Qian Z, Zhou Q. Generative steganography network. In: Proceedings of the 30th ACM International Conference on Multimedia; 2022 Oct 10–14; Lisbon, Portugal.
19. Peng F, Chen G, Long M. A robust coverless steganography based on generative adversarial networks and gradient descent approximation. *IEEE Trans Circ Syst Video Technol.* 2022;32(9):5817–29. doi:10.1109/TCSVT.2022.3161419.
20. Wei P, Luo G, Song Q, Zhang X, Qian Z, Li S. Generative steganographic flow. In: 2022 IEEE International Conference on Multimedia and Expo (ICME); 2022 Jul 18–22; Taipei, Taiwan.
21. Zhou Z, Su Y, Li J, Yu K, Wu QMJ, Fu Z, et al. Secret-to-image reversible transformation for generative steganography. *IEEE Trans Dependable Secure Comput.* 2022;20(5):4118–34. doi:10.1109/TDSC.2022.3217661.
22. Wei P, Zhou Q, Wang Z, Qian Z, Zhang X, Li S. Generative steganography diffusion. arXiv:2305.03472. 2023. doi: 10.48550/arXiv.2305.03472.
23. Peng Y, Hu D, Wang Y, Chen K, Pei G, Zhang W. StegaDDPM: generative image steganography based on denoising diffusion probabilistic model. In: Proceedings of the 31st ACM International Conference on Multimedia (ACM MM); 2023 Oct 29–Nov 3; Ottawa, ON, Canada.
24. Zhou Z, Dong X, Meng R, Wang M, Yan H, Yu K, et al. Generative steganography via auto-generation of semantic object contours. *IEEE Trans Inf Forensics Secur.* 2023;18:2751–65. doi:10.1109/TIFS.2023.3268843.
25. Liu X, Ma Z, Ma J, Zhang J, Schaefer G, Fang H. Image disentanglement autoencoder for steganography without embedding. In: 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR); 2022 Jun 18–24; New Orleans, LA, USA.
26. Sun Y, Liu J, Zhang R. A robust generative image steganography method based on guidance features in image synthesis. In: 2023 IEEE International Conference on Multimedia and Expo (ICME); 2023 Jul 10–14; Brisbane, Australia.
27. Zhu JY, Park T, Isola P, Efros AA. Unpaired image-to-image translation using cycle-consistent adversarial networks. In: 2017 IEEE International Conference on Computer Vision (ICCV); 2017 Oct 22–29; Venice, Italy.
28. Mittal A, Soundararajan R, Bovik AC. Making a “completely blind” image quality analyzer. *IEEE Signal Process Lett.* 2013;20(3):209–12. doi:10.1109/LSP.2012.2227726.
29. Blau Y, Mechrez R, Timofte R, Michaeli T, Zelnik-Manor L. The 2018 PIRM challenge on perceptual image super-resolution. In: The European Conference on Computer Vision (ECCV) Workshops; 2018 Sep 8–14; Munich, Germany.
30. Mittal A, Moorthy AK, Bovik AC. No-reference image quality assessment in the spatial domain. *IEEE Trans Image Process.* 2012;21(12):4695–708. doi:10.1109/TIP.2012.2214050.
31. Xu G, Wu HZ, Shi YQ. Structural design of convolutional neural networks for steganalysis. *IEEE Signal Process Lett.* 2016;23(5):708–12. doi:10.1109/LSP.2016.2548421.
32. Boroumand M, Chen M, Fridrich J. Deep residual network for steganalysis of digital images. *IEEE Trans Inf Forensics Secur.* 2019;14(5):1181–93. doi:10.1109/TIFS.2018.2871749.
33. You W, Zhang H, Zhao X. A Siamese CNN for image steganalysis. *IEEE Trans Inf Forensics Secur.* 2020;16:291–306. doi:10.1109/TIFS.2020.3013204.