# MDDM: Practical Message-Driven Generative Image Steganography Based on Diffusion Models

Zihao Xu $^1$  Dawei Xu $^{12}$  Zihan Li $^2$  Chuan Zhang $^2$ 

### Abstract

Generative image steganography (GIS) is an emerging technique that conceals secret messages in the generation of images. Compared to GAN-based or flow-based GIS schemes, diffusion model-based solutions can provide high-quality and more diverse images, thus receiving considerable attention recently. However, previous GIS schemes still face challenges in terms of extraction accuracy, controllability, and practicality. To address the above issues, this paper proposes a practical message-driven GIS framework based on diffusion models, called MDDM. Specifically, by utilizing Cardan Grille, we encode messages into Gaussian noise, which serves as the initial input for image generation, enabling users to generate diverse images via controllable prompts without additional training. During the information extraction process, receivers only need to use the pre-shared Cardan Grille to perform exact diffusion inversion and recover the messages without requiring the image generation seeds or prompts. Experimental results demonstrate that MDDM offers notable advantages in terms of accuracy, controllability, practicality, and security. With flexible strategies, MDDM can always achieve almost 100% accuracy. Additionally, MDDM demonstrates certain robustness and exhibits potential for application in watermarking tasks.

# 1. Introduction

With the proliferation of the Internet and digital communications, information security issues have become increasingly prominent. As a crucial technique, steganography offers a secure and covert method of communication. With the increasing volume of image data on digital platforms, images have become a primary medium for covert data embedding (Bachrach & Shih, 2011; Subramanian et al., 2021; Mandal et al., 2022). Image steganography conceals information within images in a manner that prevents unauthorized access or detection. It is widely used in copyright protection (Altaay et al., 2012) and covert communication (Juneja, 2014).

However, achieving provably secure image steganography remains a major challenge. Traditional methods (van Schyndel et al., 1994; Liao et al., 2020; Yang et al., 2020; Su et al., 2021; Chan & Cheng, 2004) usually embed secret data through cover image modification. These methods often leave traces, making them susceptible to detection by steganalysis tools, and are therefore deemed empirically secure rather than provably secure (Hopper et al., 2002).

Provable security in image steganography requires that the distributions of the normal and stego images be indistinguishable (Weiming, 2023). Recently, the rise of generative artificial intelligence (AIGC) has provided a more diverse and flexible steganographic environment. Meanwhile, AI-generated data follows a controllable distribution. Leveraging this advantage, generative image steganography (GIS) has emerged and demonstrated strong resistance to typical steganographic attacks.

Compared with traditional methods, GIS schemes have better anti-detection performance against existing statistical feature-based steganalysis methods. GIS is mainly based on generative adversarial networks (GAN) (Yang et al., 2024b; Zhou et al., 2023a; Li et al., 2020; Su et al., 2024; You et al., 2022), flow-based models (such as Glow) (Zhou et al., 2023b; Wei et al., 2022; Xu et al., 2022b) and the latest diffusion models (DM) (Hu et al., 2024; Peng et al., 2023; 2024; Yu et al., 2024; Jois et al., 2024). Although GIS has achieved notable progress, it still faces limitations in practical applications. For example, GAN-based methods require much training (Yang et al., 2024b), and are costly and difficult to control; flow-based methods require high computational resources during training and inference, especially in high-resolution image generation tasks (Zhou et al., 2023b). Diffusion model-based methods have made breakthroughs in generation quality and provide a better steganographic environment due to their widespread application in image generation. However, they still suffer

<sup>&</sup>lt;sup>1</sup>Changchun University <sup>2</sup>Beijing Institute of Technology. Correspondence to: Chuan Zhang <chuanz@bit.edu.cn>.

Proceedings of the  $42^{nd}$  International Conference on Machine Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025 by the author(s).



*Figure 1.* The overall framework of MDDM. The sender needs to share the Cardan Grille with the receiver in advance to determine the location of the information in the initial noise. The Cardan Grille is randomly generated. For the same message and the same Cardan Grille, the sender can use MDDM to generate multiple pictures for selection.

from certain limitations. Specifically, existing diffusion model steganography methods can be divided into seed-free and seed-dependent types. Seed-free methods essentially amount to image modification (Yu et al., 2024; Yang et al., 2024a), causing the generated stego images to maintain similarity to the originals and thus reducing the diversity of image generation. On the other hand, seed-dependent methods (Peng et al., 2023; Jois et al., 2024; Peng et al., 2024) require a complete reproduction of the generation process. Following a one-time seed exchange, the sender is constrained in generating arbitrary images, potentially necessitating additional communication.

To address these issues, we propose MDDM, a practical, message-driven generative image steganography framework (see Figure 1). The sender generates a Cardan Grille and maps the uniformly distributed binary message onto noise that follows a standard normal distribution through a carefully designed encoding strategy. Then, the sender uses DDIM and conditional text to generate a stego image starting from noise. In this process, the sender can generate different noises while keeping the message and Cardan Grille unchanged and use different prompts to generate a variety of images for selection. Since the generation process is identical to normal image generation, the stego image is indistinguishable from a regular image. The receiver performs exact diffusion inversion on the stego image to reconstruct the noise and uses the pre-shared Cardan Grille to obtain the hidden message without knowing the image generation seed or the prompt.

The main contributions of MDDM are as follows:

- Message-driven GIS Framework: We leverage the reversibility of DDIM to develop a message-driven image steganography framework based on diffusion models. By utilizing a Cardan Grille, we encode secret message into the initial noise for image generation. This enables MDDM to achieve imperceptible steganography.
- **Controllability and Practicality**: Without compromising the concealment of secret information, MDDM uses controllable conditional text and random seeds as inputs to generate high-quality, diverse stego images without requiring any additional training. By adjusting the length of the secret message, MDDM can be flexibly applied to both information hiding and anticounterfeiting watermarking.
- Security and Robustness: The images generated by MDDM are consistent with randomly generated images in distribution, which is considered to be provably secure. At the same time, Cardan Grille in MDDM has good resistance to exhaustive attacks. In addition, a large number of experiments have shown that MDDM has certain robustness and can ensure a high extraction accuracy even under VAE compression and common image distortion conditions so that it can also be applied to watermarking tasks.

# 2. Related Work

The practice of image steganography has evolved considerably over time. With the advancement of digital technology, it has undergone significant evolution. The techniques have evolved from simple pixel modifications to the application of complex algorithms and advanced encryption methods, becoming increasingly sophisticated. In recent years, research has focused on improving concealment, enhancing security, and increasing resistance to attacks. Additionally, the application of artificial intelligence has further promoted innovation in image steganography.

### 2.1. Cover Image Based Steganography

From early research to the present day, most studies on image steganography have focused on embedding information directly into images. The oldest and most classic embedded steganography scheme is LSB (van Schyndel et al., 1994; Chan & Cheng, 2004), which embeds information into the least significant bit of image pixels. Furthermore, researchers have proposed manually designed methods (Liao et al., 2020; Su et al., 2021) and neural network-based techniques (Yang et al., 2020) as adaptive approaches to reduce image distortion caused by embedding. There are also frequency domain-based steganography methods, including DctDwt (Al-Haj, 2007) and DctDwtSvd (Navas et al., 2008). In recent years, deep learning based steganography (Tancik et al., 2020; Guan et al., 2022; Lu et al., 2021; Xu et al., 2022a; Jing et al., 2021; Zhu et al., 2018; Fadhil et al., 2023) has emerged. The first end-to-end trainable framework for data hiding, HIDDEN (Zhu et al., 2018), achieves embedding and extraction through autoencoders. Recent research on embedded steganography has largely focused on invertible neural networks (INNs) (Guan et al., 2022; Lu et al., 2021; Xu et al., 2022a; Jing et al., 2021). For example, HiNet (Jing et al., 2021) is the first attempt to utilize invertible neural networks for image hiding tasks, where information hiding and extraction can be accomplished through a single network. However, the aforementioned embedded steganography schemes often leave modification traces on the image, which can still be detected (Fu et al., 2024).

#### 2.2. Generative Image Steganography

Generative image steganography has gained traction and is generally considered more secure than traditional embedding-based methods (Yang et al., 2019). Generative image steganography, leveraging generative models, is commonly divided into three main categories: flow-based (Zhou et al., 2023b; Wei et al., 2022; Xu et al., 2022b), GAN-based (Yang et al., 2024b; Zhou et al., 2023a; Li et al., 2020; Su et al., 2024; You et al., 2022), and diffusion-based methods. Flow-based methods primarily exploit the reversible properties of such models, with SR2IT (Zhou et al., 2023b) being a representative study. Compared to flow-based methods, GAN-based research has been explored more extensively. For example, PARIS (Yang et al., 2024b) developed a provably secure method to image steganography and attached a noise module to the generator to improve robustness. Compared to GAN-based methods, diffusion models (Hu et al., 2024; Peng et al., 2023; 2024; Yu et al., 2024; Jois et al., 2024) have recently developed generative steganography methods due to their advantages in generation quality and diversity. Currently, most of the research can be divided into two categories: image translation-based and generative process recovery-based. However, these approaches still face challenges related to computational efficiency, practical deployment, and steganographic effectiveness.

#### 2.3. Cardan Grille

The Cardan Grille is a classical steganographic technique, typically used as a shared key between sender and receiver for hiding and extracting information, and has considerable security (Utepbergenov et al., 2013). Recently, some studies have also tried to apply it to generative steganography. Most of these studies are based on the Cardan Grille and use GAN to perform image restoration to generate stego images (Liu et al., 2018; Zhang et al., 2019c; Wang et al., 2021). Although these methods have demonstrated promising results, particularly in terms of security, there are still problems such as uncontrollable generation effects, low stego capacity, and poor versatility.

#### 2.4. Diffusion Models

The rapid advancement of diffusion models has led to a surge in AI-generated images across the internet. Early diffusion models were primarily based on pixel-space diffusion like DDPM (Ho et al., 2020). However, in recent years, models like Stable Diffusion (Rombach et al., 2022), based on latent diffusion, have become the dominant approach. DDIM (Song et al., 2020a), a widely used sampling method, has been applied to both unconditional and conditional diffusion models. Notably, DDIM inversion techniques have recently been introduced, enabling fine-grained editing of AI-generated images (Wallace et al., 2023; Zhang et al., 2025; Wang et al., 2024; Mokady et al., 2023).

# 3. Method

We regard the message hiding process Gen as generating a natural image from secret noise and the extraction process Ext as recovering the noise from a generated image:

$$\begin{aligned} & \texttt{Gen}\left(sk,m\right) = I_{stego}, \\ & \texttt{Ext}\left(sk, I_{stego}\right) = m, \end{aligned} \tag{1}$$

where sk is the secret key, m is the secret message, and  $I_{stego}$  is the generated image.

In the following paragraphs, we provide a detailed description of our implementation of the message-driven diffusion model steganography (MDDM).

### 3.1. The Basics of MDDM

In our work, we exploit the properties of DDIM (Song et al., 2020a) to construct the information hiding and extraction process. DDIM is an improved diffusion model that, unlike DDPM (Ho et al., 2020), defines a non-Markov forward process. DDIM can improve the speed and quality of image generation and always be used as a sampling method. We exploit the following properties of DDIM to design our MDDM.

**Deterministic Sampling.** For steganography, we are particularly interested in the backward process of DDIM, which is a key distinction between DDIM and DDPM. We denote the number of diffusion steps as T, the initial noise as  $x_T$ , the indices of intermediate steps decrease progressively, and the final result of diffusion as  $x_0$ . The sampling formula is as follows:

$$x_{t-1} = \sqrt{\alpha_{t-1}} \left( \frac{\boldsymbol{x}_t - \sqrt{1 - \alpha_t} \boldsymbol{\epsilon}_{\theta}^{(t)}(\boldsymbol{x}_t)}{\sqrt{\alpha_t}} \right) + \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \boldsymbol{\epsilon}_{\theta}^{(t)}(\boldsymbol{x}_t) + \sigma_t \boldsymbol{\epsilon}_t.$$
(2)

When  $\sigma_t = 0$ , randomness disappears from the formula and this process becomes deterministic. The initial noise  $x_T$  completely determines the final result  $x_0$ , which is a potential one-to-one relationship. In our work, we adopt deterministic DDIM sampling across various diffusion models.

The Lossless Exact Diffusion Inversion. Due to the above characteristics, DDIM Inversion performs reverse mapping based on the diffusion process. Its goal is to reverse the sampling process. Specifically, given an image, DDIM Inversion maps it back to the initial noise state in the diffusion process. The original image can be reconstructed by initiating the conventional sampling process from the inverted noise. Meanwhile, this method serves as a fundamental component of image editing.

However, DDIM inversion is known to be unreliable, often leading to substantial discrepancies between the inverted and original images. Our goal is to achieve precise and reliable inversion. For example, we can employ EDICT (Wallace et al., 2023), which uses a combination of coupling and averaging/dilating steps to achieve exact inversion of the diffusion process. It is worth noting that while the sampling quality of EDICT may not be optimal, our task does not involve image editing. Instead, we focus solely on ensuring the precision of inversion.



*Figure 2.* If the initial noise inputs are similar, then the images generated through deterministic sampling are also similar.

**Correlation between Generated Image and Initial Noise.** DDIM sampling is a deterministic process, which means that under the same conditions, two inputs with slight differences will follow almost identical sampling trajectories, as shown in Figure 2. Let  $x_T^{(1)}$  and  $x_T^{(2)}$  denote two Gaussian noise inputs with a small difference:

$$\|x_T^{(1)} - x_T^{(2)}\|_2 \approx \delta, \tag{3}$$

where  $\delta$  represents a small difference.

Since the space of the input is continuous (Song et al., 2020b),  $x_T^{(1)}$  and  $x_T^{(2)}$  will undergo similar sampling processes and thus yield similar output images  $x_0^{(1)}$  and  $x_0^{(2)}$ . As illustrated in Figure 3, when the noise follows a standard normal distribution, slight perturbations applied to select regions of the initial noise do not significantly alter the content of the generated image. Meanwhile, by employing the exact diffusion inversion method mentioned previously, the bias in the inversion process is effectively eliminated. Consequently, this establishes a correlation between the similarity of initial noise inputs and the similarity of their generated images.



*Figure 3.* Images generated using DDIM sampling on DDPM (top row) and Stable Diffusion (bottom row), with different perturbation ratios applied to the initial noise inputs.

The Information Loss Is Acceptable. The encoding and decoding process of an image inevitably leads to information loss, but our research shows that this loss is within an acceptable range. We study the loss of conditional generation based on the Stable Diffusion (SD) and unconditional generation based on the DDPM under DDIM sampling using EDICT inversion. The original noise is called  $x_T$ , and the inverted noise is called  $x'_T$ . In order to evaluate the information loss, we flatten the high-dimensional space into a one-dimensional representation and calculate the absolute difference  $|x'_{T}[i] - x_{T}[i]|$  element by element, as shown in Figure 4. The results show that, whether based on the latent diffusion or the pixel-space diffusion, the absolute difference between the original noise and the inverted noise is smaller than that of random noise and most of them remain within 1.0. However, we also found that the unconditional generation based on the pixel-space diffusion model may be unstable, and the range of absolute differences is slightly larger than that of the latent diffusion. At the same time, considering that conditional generation based on latent diffusion models is more widely used in practical applications, the method in this paper will focus on latent diffusion models.



*Figure 4.* Distribution of the absolute value of the difference between the original noise and the noise after diffusion inversion for different cases.

#### 3.2. Framework of MDDM

MDDM requires no training but relies on multiple pretrained modules. Specifically, it leverages latent diffusion models (e.g., Stable Diffusion) or pixel-space diffusion models (e.g., DDPM), together with a DDIM scheduler and any exact diffusion-inversion technique (such as EDICT), all of which are readily obtainable.

In MDDM, the sender and receiver first share a randomly generated Cardan Grille, equivalent to sk in Equation (1), which specifies the location of the hidden message but does

not contain the message itself, as shown in Figure 5. Using our encoding strategy, the binary message is mapped to a noise following a standard normal distribution. Then, the pre-trained diffusion model can generate images using this noise. Since the entire process mirrors general image generation, the stego image remains indistinguishable from a normally generated image. The receiver uses the obtained image to perform exact diffusion inversion, reconstructing the noise and recovering the hidden binary data through the pre-shared Cardan Grille without requiring knowledge of the image generation seed or prompt. The benefit of MDDM is that the sender can generate multiple images without changing the secret information and Cardan Grille and select the images with high information extraction accuracy and good quality for transmission.



*Figure 5.* Visualization of Cardan Grille where the blue grid area represents the Cardan Grille. It is equivalent to a symmetric key that does not alter the distribution of the noise.

Message Hiding. In the hiding stage, considering that information loss may reduce the accuracy of steganographic content extraction, we first determine the area suitable for steganography. As shown in Section 3.1, regardless of latent diffusion or pixel-space diffusion, the absolute difference between the values before and after the accurate inversion of the diffusion is usually in the range of -1 to 1. Therefore, in order to ensure that potential errors do not affect the extraction accuracy, we define the highly robust area of steganography as data with a standard normal distribution in the range of  $(-\infty, -1)$  and  $(1, +\infty)$ . Here, we introduce the standard normal distribution cumulative distribution function :

$$\Phi(x) = \int_{-\infty}^{x} \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}} dt.$$
 (4)

For a noise of size  $1 \times c \times h \times w$  (Batch  $\times$  Channel  $\times$  Height  $\times$  Width), the amount of data in the highly robust area  $L_{max}$  can be calculated by the following formula:

$$L_{max} = 2 \times (1 - \Phi(1)) \times 1 \times c \times h \times w.$$
 (5)

For example, for noise of size  $1 \times 4 \times 64 \times 64$ , there are 5199 elements in the highly robust area. When we determine that the length l of the secret message m to be transmitted is less than  $L_{max}$ , we can calculate the truncation threshold k by the inverse function of  $\Phi(x)$  as follows:

$$k = \Phi^{-1} \left(1 - \frac{l}{2 \times 1 \times c \times h \times w}\right).$$
(6)

Next, the Cardan Grille CG is a list of length l with unique elements, generated by randomly sampling from the integer set  $\{1, 2, ..., 1 \times c \times h \times w\}$ . Let  $\{\zeta_r\}_{r=1}^{1 \times c \times h \times w}$  denote the elements of this integer set, then:

$$CG = \texttt{RandomSample}(\{\zeta_r\}_{r=1}^{1 \times c \times h \times w}, l), \qquad (7)$$

where RandomSample $(\cdot, l)$  denotes selecting l distinct elements at random.

It is worth noting that the secret message m is a uniformly distributed binary string formed by encryption or other means. Since this process can be implemented in many ways and is not the focus of this article, this article does not describe how to generate a uniformly distributed binary string in detail.

At this point, the preparatory steps are complete. We next generate the initial noise by performing truncated sampling from the standard normal distribution at the previously determined Cardan Grille positions CG. In each iteration, we draw a fresh batch of samples from  $\mathcal{N}(0, \mathbf{I})$  and partition them into three pools according to the truncation threshold k. For each bit of the secret message m, we then randomly sample from the  $(-\infty, -k)$  pool if the bit is 0, or from the  $(k, +\infty)$  pool if the bit is 1, and insert the sampled values into the one-dimensional noise vector at the indices specified by CG. After filling all message-bearing positions, the remaining entries are populated by sampling from the truncated interval (-k, k) for each new draw and placing them into the unused slots of the noise vector. Finally, we reshape this one-dimensional vector to the diffusion model's required noise size  $1 \times c \times h \times w$ . In Section 4, we demonstrate that the constructed noise is statistically indistinguishable from Gaussian noise following a standard normal distribution, thereby offering a provable security guarantee for our method. Algorithm 1 details this procedure, and Appendix A presents an optimized variant designed to accelerate noise generation for high-dimensional inputs.

After applying the above encoding rules to map the secret message m to the initial noise  $x_T$ , the sender directly uses  $x_T$  for DDIM sampling, where the prompt can be arbitrary, resulting in the stego image  $I_{stego}$ . Algorithm 2 clearly shows this process.

Algorithm 1 Generate initial noise

**Input:** The noise size (1, c, h, w), secret message m with length l, Cardan Grille CG, truncation threshold k**Output:** Initial noise  $x_T$ **Function**  $F_1(n \in \{0, 1\})$ *list* ~  $\mathcal{N}(0, \mathbf{I})$ if n = 0 then  $y = \text{RandomChoice}\{u \mid u < -k, u \in list\}$ else  $y = \texttt{RandomChoice}\{u \mid u > k, u \in list\}$ end if **Return** y **Function** *F*<sub>2</sub>  $list \sim \mathcal{N}(0, \mathbf{I})$  $y = \texttt{RandomChoice}\{u \mid -k \le u \le k, u \in list\}$ **Return** *y* Init noise  $x_T$ for i < l do  $x_T[CG[i]] = F_1(m[i])$ end for for  $j < 1 \times c \times h \times w$  do if  $j \notin CG$  then  $x_T[j] = F_2$ end if end for

Algorithm 2 Generate stego image based on LDM
<b>Input:</b> The pretrained LDM $pipe$ , VAE decoder $\mathcal{D}$ and
DDIM sampler DDIMSampling, the initial noise $x_T \sim$
$\mathcal{N}(0,\mathbf{I})$ , prompt $Text$
<b>Output:</b> Stego image $I_{stego}$
$x_0 = \text{DDIMSampling}(pipe, Text, x_T)$
$I_{stego} = \mathcal{D}\left(x_0\right)$

**Message Extraction.** The receiver first obtains the stego image from the public channel and then applies the preselected precise inversion method to revert it through the same number of steps (usually 50 by default), without requiring reconstruction, to quickly obtain the inverted noise. Using the pre-shared Cardan Grille, the receiver can sequentially extract the corresponding values and derive the binary message based on whether they are greater or less than zero:

$$m'[i] = \begin{cases} 0, & \text{if } x'_T[CG[i]] < 0, \\ & & \\ 1, & \text{if } x'_T[CG[i]] \ge 0. \end{cases}$$
(8)

Ideally, the m' obtained by the receiver's inversion should be consistent with the secret message m hidden by the sender.





Acc: 100%

jungle, muted the bed colors

Acc: 100%

Acc: 100% small island by hair, front view and vibrant colors the sea





Acc: 100%

Acc: 99.90% Acc: 100%

Astronaut in a A cat is lying on Lighthouse on a A girl with blond A Van Gogh inspired landscape with swirling skies

Acc: 100%

(a) Guidance Scale: 3.5



jungle, muted the bed colors the sea

Acc: 99.80%

small island by

Astronaut in a A cat is lying on Lighthouse on a A girl with blond A Van Gogh inspired landscape with swirling skies hair, front view and vibrant colors

(b) Guidance Scale: 7.5

Figure 6. The performance of MDDM under different guidance scales. The default guidance scale setting of 7.5 can generally achieve better visual effects.

# 4. Security Analysis

We first examine how MDDM steganography affects the distribution of generated images to demonstrate its provable resistance to steganalysis. To formally describe the sampling process, let  $x_T \in \mathbb{R}^{1 \times c \times h \times w}$  denote the initial noise, and let CG from Equation (7) denote the set of indices corresponding to the Cardan Grille positions, which are determined uniformly at random according to the length of the binary message string. For each index  $i \in CG$ , we embed one bit  $m[i] \in \{0, 1\}$  by sampling the corresponding value  $x_T^{(i)}$  independently from a truncated standard normal distribution. Specifically, we define the conditional distribution for  $x_T^{(i)}$  as follows:

$$x_T^{(i)} \sim \begin{cases} \mathcal{N}(0,1) \text{ truncated to } (k,\infty), & \text{if } m[i]=1, \\ \mathcal{N}(0,1) \text{ truncated to } (-\infty,-k), \text{if } m[i]=0. \end{cases}$$
(9)

The value of threshold k is chosen according to the ratio of the payload length to the total embedding capacity, such that larger payloads require smaller k to maintain coverage, while smaller payloads allow for more extreme truncation and thus potentially stronger security. For all remaining positions  $j \notin CG$ , we sample the values  $x_T^{(j)}$  independently from the central region of the same distribution:

$$x_T^{(j)} \sim \mathcal{N}(0,1)$$
 truncated to  $(-k,k)$ . (10)

Since all samples are drawn independently, and the truncation preserves the symmetry and local density properties of the original distribution, the resulting noise vector  $x_T$ maintains the overall statistical characteristics of a standard normal distribution. Moreover, because the Cardan Grille positions are chosen uniformly at random and the embedded bits are uniformly distributed, the sampling strategy does not introduce any detectable structural bias. As such, the steganographic modification remains statistically indistinguishable from standard noise under conventional detection methods. According to Appendix B, MDDM is provably secure.

Next, we need to prove that the Cardan Grille used by MDDM possesses privacy security. According to Appendix C, even when reusing the Cardan Grille for steganography across multiple images, a third party cannot obtain the ordering privacy information of the Cardan Grille, thus unable to threaten the privacy security of the hidden information.

# **5.** Experiments

We use publicly available pre-trained diffusion models as the pipeline based on PyTorch and diffusers <sup>1</sup>, DDIM (Song et al., 2020a) for sampling, and EDICT (Wallace et al., 2023) as the exact diffusion inversion method to perform MDDM.

<sup>&</sup>lt;sup>1</sup>https://huggingface.co/docs/diffusers/index

Table 1. The robustness of MDDM is compared with that of baseline methods in terms of accuracy at a size of  $512 \times 512$ . The evaluation includes results for the "PNG" image (i.e., without attack), along with various attacks such as "JPEG" (image compression with an image quality of 95%), "Resize" (image scaling by a factor of 0.5), "Gaussian Blur (GB)" (with a radius of 1.0), "Drop" (randomly adding a white square of 0.5% area), "Brightness" (with a factor of 0.15), and "Rotation" (an angle of 0.3 degrees).

Method	Version	Capacity (Bits)	PNG	JPEG (95)	Resize (0.5)	GB (1.0)	Drop (0.005)	Brightness (0.15)	Rotation (0.3)
DwtDct (Al-Haj, 2007)	-	32	0.83	0.65	0.63	0.68	0.83	0.83	0.72
DwtDctSvd (Navas et al., 2008)	-	32	1.00	1.00	1.00	1.00	1.00	0.78	1.00
RivaGAN (Zhang et al., 2019a)	-	32	0.99	0.99	0.99	0.99	0.99	0.99	0.99
LaWa (Rezaei et al., 2025)	SD-V1.4	48	1.00	1.00	1.00	0.99	1.00	1.00	1.00
MDDM (Ours)	SD-V1.4	48	1.00	1.00	0.99	1.00	1.00	1.00	0.99
GS (Yang et al., 2024c)	SD-V2.1	256	1.00	1.00	1.00	1.00	1.00	1.00	1.00
MDDM (Ours)	SD-V2.1	256	1.00	1.00	0.99	1.00	1.00	1.00	0.99

*Table 2.* Ablation experiment. Comparison of steganography ability of MDDM under different truncation thresholds.

Туре	Truncation Threshold	Bits	Face	Accuracy Bedroom	Cat
Uncondition -3 256 × 250	al $\begin{array}{c} 3\sigma\\ 2\sigma\\ 6\\ 1\sigma\end{array}$	531 8946 62386	99.92% 98.53% 92.56%	98.73% 95.04% 84.95%	99.97% 98.21% 95.19%
Conditional $512 \times 512$	$rac{3\sigma}{2\sigma}$ $1\sigma$	44 745 5199	100% 99.87% 99.46%	99.89% 99.73% 99.20%	99.95% 99.88% 99.13%

The forward diffusion and backward diffusion processes each contain 50 steps. All experiments are conducted on an NVIDIA RTX 3090. MDDM does not require fine-tuning or further training of the pre-trained diffusion model, but slightly improves the DDIM sampling and EDICT inversion processes. For the general latent diffusion, we utilize the Stable Diffusion v1.5<sup>2</sup>, and in the comparative experiment based on the latent diffusion model on watermarking, we use Stable Diffusion v1.4 and Stable Diffusion v2.1 for experiments due to the constraints of the comparative method, and in the comparative experiment on steganography, we use the pre-trained DDPM models publicly available on huggingface <sup>3</sup>.

### 5.1. Overall Performance

Our experiments cover various evaluations, including ablation studies, quantitative comparisons, and visualizations.

**Image Quality.** The visualization experiment in Figure 6 demonstrates the performance of conditional generation under different guidance scales. These images hide 1024

bits of information at a size of  $512 \times 512$ , and the upper and lower images in each column are generated using the same seed. The visualization results indicate that at the default guidance scale of 7.5, the image generation effect and message accuracy achieve a high level. Based on this, in subsequent experiments based on Stable Diffusion, we use the guidance scale of 7.5 by default to conditionally generate stego images. Figure 9 provides additional examples of generated images, demonstrating that our results are of high quality and diversity, closely resembling randomly generated images. In addition, we also evaluate the FID (Heusel et al., 2017) on the generated images, as detailed in Appendix D.1.

Steganography Capabilities. First, we conduct ablation experiments, as shown in Table 2, to assess the performance of MDDM with varying secret message lengths in both unconditional and conditional diffusion. Here, accuracy (Acc) represents the correctness of the extracted message, while capacity reflects the amount of information that can be embedded. The specific settings of the ablation study are detailed in Appendix D.2. Experimental results indicate that as the length of the secret message increases, extraction accuracy declines but remains at a reasonable level. In addition, compared with the unconditional diffusion based on pixel space, the conditional diffusion based on latent space achieves higher extraction accuracy. This further confirms our findings in Section 3.1. Therefore, our method is more suitable for latent diffusion. Moreover, we find that although variations in CPU and GPU configurations may yield differences in the initially generated noise, these variations have no significant impact on the overall performance of MDDM.

In Appendix D.3, we compare MDDM with the latest diffusion model-based generative steganography methods under identical conditions. All comparison methods involve embedding hidden bit messages. The above experiments demonstrate the effectiveness of our proposed method.

<sup>&</sup>lt;sup>2</sup>https://huggingface.co/stable-diffusion-v1-5/stablediffusion-v1-5

<sup>&</sup>lt;sup>3</sup>https://huggingface.co/

#### 5.2. Robustness

In Section 3.1, we demonstrate a correlation between the similarity of the initial noise and the similarity of the final generated image. Consequently, when the final generated image remains largely unchanged, MDDM is not significantly affected, particularly by common image distortions such as JPEG compression. As shown in Table 1, MDDM resists various image-based attacks and offers more flexible information capacity. Note that since DwtDct, DwtDctSvd, and RivaGAN are not generative image watermark methods, we use the Stable Diffusion to generate cover images. For the comparison of conditional watermark generation methods based on latent diffusion, we utilize prompts from Stable-Diffusion-Prompts <sup>4</sup>. Results indicate that MDDM can, to some extent, serve as a watermarking method.

#### 5.3. Controllability and Practicality

MDDM enables random image generation while keeping the secret message and Cardan Grille unchanged, allowing the sender to continuously generate diverse images for selection, as shown in Figure 9. Table 5 in Appendix D.4 presents a comparison of our method with other methods in terms of diversity, where diversity refers to the number of images that can be generated after a single communication and the sharing of a secret key or seed. Results show that existing regeneration-based methods (Peng et al., 2023; Jois et al., 2024; Peng et al., 2024) require changing the noise seed if the generated image is unsatisfactory, which may lead to increased communication overhead and potential risks. In contrast, our method synchronizes the encryption key for the message and the random seed for the Cardan Grille during initial setup. After that, subsequent exchanges no longer require direct point-to-point communication. For example, the sender can use the initial seed to generate the Cardan Grille and then post the coded images to their social media page. The receiver can then download these images and extract the information using the initial seed. For subsequent transmissions, the seed can be changed sequentially, and the receiver can do the same; thus, communication can proceed even asynchronously. In summary, with MDDM, each image sent by the sender is arbitrarily controllable, and multiple images can be generated for selection, which avoids arousing suspicion from potential attackers and ensures the accuracy of message extraction. This further demonstrates the safety and practicality of MDDM.

#### 5.4. Steganalysis

We test MDDM on ZhuNet (Zhang et al., 2019b), SiaSteg-Net (You et al., 2020), XuNet (Xu et al., 2016), YeNet



*Figure 7.* Results of detection using multiple advanced steganalysis methods.

(Ye et al., 2017), StegNet (Deng et al., 2019) and SRNet (Boroumand et al., 2018), as shown in Figure 7. The results are generally around 50%, indicating that the distribution cannot be distinguished, which is consistent with the theoretical analysis.

### 6. Conclusion

We propose MDDM, a practical message-driven generative image steganography method based on diffusion models. The process of generating stego images using MDDM is identical to the general image diffusion generation process, and MDDM performs well in terms of controllability, security, image quality, robustness, and extraction accuracy. In the future, we will introduce error-correcting codes and explore the performance of MDDM in other diffusion models.

### Acknowledgements

This work is partially supported by the 111 Project under Grant B21044, the Sichuan Science and Technology Program under Grants 2021ZDZX0011, the National Natural Science Foundation of China under Grant 62376175 and 62472032, and the Young Elite Scientists Sponsorship Program by CAST (Grant No. 2023QNRC001).

# **Impact Statement**

Generative image steganography (GIS) has recently emerged as a prominent research direction within the broader field of steganography. The research presented in this paper not only advances the development of GIS but also demonstrates the potential for real-world applications. For instance, the proposed MDDM is training-free, which may facilitate the rapid adoption of GIS technologies. Overall, our work is anticipated to make an impact in domains such as information security and digital copyright protection.

<sup>&</sup>lt;sup>4</sup>https://huggingface.co/datasets/Gustavosta/Stable-Diffusion-Prompts

### References

- Al-Haj, A. Combined dwt-dct digital image watermarking. *Journal of computer science*, 3(9):740–746, 2007.
- Altaay, A. A. J., Sahib, S. B., and Zamani, M. An introduction to image steganography techniques. In 2012 international conference on advanced computer science applications and technologies (ACSAT), pp. 122–126. IEEE, 2012.
- Bachrach, M. and Shih, F. Y. Image steganography and steganalysis. Wiley Interdisciplinary Reviews: Computational Statistics, 3(3):251–259, 2011.
- Boroumand, M., Chen, M., and Fridrich, J. Deep residual network for steganalysis of digital images. *IEEE Transactions on Information Forensics and Security*, 14(5): 1181–1193, 2018.
- Cachin, C. An information-theoretic model for steganography. In *International Workshop on Information Hiding*, pp. 306–318. Springer, 1998.
- Chan, C.-K. and Cheng, L.-M. Hiding data in images by simple lsb substitution. *Pattern recognition*, 37(3):469–474, 2004.
- Deng, X., Chen, B., Luo, W., and Luo, D. Fast and effective global covariance pooling network for image steganalysis. In *Proceedings of the ACM workshop on information hiding and multimedia security*, pp. 230–234, 2019.
- Fadhil, A. M., Jalo, H. N., and Mohammad, O. F. Improved security of a deep learning-based steganography system with imperceptibility preservation. *International journal of electrical and computer engineering systems*, 14(1): 73–81, 2023.
- Fu, T., Chen, L., Jiang, Y., Jia, J., and Fu, Z. Image steganalysis based on dual-path enhancement and fractal downsampling. *IEEE Transactions on Information Forensics and Security*, pp. 1–1, 2024. doi: 10.1109/TIFS.2024.3493615.
- Guan, Z., Jing, J., Deng, X., Xu, M., Jiang, L., Zhang, Z., and Li, Y. Deepmih: Deep invertible network for multiple image hiding. *IEEE Transactions on Pattern Analysis* and Machine Intelligence, 45(1):372–390, 2022.
- Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., and Hochreiter, S. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Proceedings* of the 31st International Conference on Neural Information Processing Systems, NIPS'17, pp. 6629–6640, Red Hook, NY, USA, 2017. Curran Associates Inc. ISBN 9781510860964.

- Ho, J., Jain, A., and Abbeel, P. Denoising diffusion probabilistic models. Advances in neural information processing systems, 33:6840–6851, 2020.
- Hopper, N. J., Langford, J., and Von Ahn, L. Provably secure steganography. In Advances in Cryptology—CRYPTO 2002: 22nd Annual International Cryptology Conference Santa Barbara, California, USA, August 18–22, 2002 Proceedings 22, pp. 77–92. Springer, 2002.
- Hu, X., Li, S., Ying, Q., Peng, W., Zhang, X., and Qian, Z. Establishing robust generative image steganography via popular stable diffusion. *IEEE Transactions on Information Forensics and Security*, 19:8094–8108, 2024. doi: 10.1109/TIFS.2024.3444311.
- Jing, J., Deng, X., Xu, M., Wang, J., and Guan, Z. Hinet: Deep image hiding by invertible network. In 2021 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 4713–4722, 2021. doi: 10.1109/ICCV48922. 2021.00469.
- Jois, T. M., Beck, G., and Kaptchuk, G. Pulsar: Secure steganography for diffusion models. In *Proceedings of* the 2024 on ACM SIGSAC Conference on Computer and Communications Security, pp. 4703–4717, 2024.
- Juneja, M. A covert communication model-based on image steganography. *International Journal of Information Security and Privacy (IJISP)*, 8(1):19–37, 2014.
- Karras, T., Laine, S., and Aila, T. A style-based generator architecture for generative adversarial networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4401–4410, 2019.
- Li, J., Niu, K., Liao, L., Wang, L., Liu, J., Lei, Y., and Zhang, M. A generative steganography method based on wgangp. In Sun, X., Wang, J., and Bertino, E. (eds.), *Artificial Intelligence and Security*, pp. 386–397, Singapore, 2020. Springer Singapore. ISBN 978-981-15-8083-3.
- Liao, X., Yu, Y., Li, B., Li, Z., and Qin, Z. A new payload partition strategy in color image steganography. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(3):685–696, 2020. doi: 10.1109/TCSVT.2019. 2896270.
- Liu, J., Zhou, T., Zhang, Z., Ke, Y., Lei, Y., and Zhang, M. Digital cardan grille: A modern approach for information hiding. In *Proceedings of the 2018 2nd International Conference on Computer Science and Artificial Intelligence*, pp. 441–446, 2018.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.

- Lu, S.-P., Wang, R., Zhong, T., and Rosin, P. L. Largecapacity image steganography based on invertible neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10816– 10825, 2021.
- Mandal, P. C., Mukherjee, I., Paul, G., and Chatterji, B. Digital image steganography: A literature survey. *Information sciences*, 609:1451–1488, 2022.
- Mokady, R., Hertz, A., Aberman, K., Pritch, Y., and Cohen-Or, D. Null-text inversion for editing real images using guided diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6038–6047, 2023.
- Navas, K., Ajay, M. C., Lekshmi, M., Archana, T. S., and Sasikumar, M. Dwt-dct-svd based watermarking. In 2008 3rd international conference on communication systems software and middleware and workshops (COM-SWARE'08), pp. 271–274. IEEE, 2008.
- Peng, Y., Hu, D., Wang, Y., Chen, K., Pei, G., and Zhang, W. Stegaddpm: Generative image steganography based on denoising diffusion probabilistic model. In *Proceedings* of the 31st ACM International Conference on Multimedia, pp. 7143–7151, 2023.
- Peng, Y., Wang, Y., Hu, D., Chen, K., Rong, X., and Zhang, W. Ldstega: Practical and robust generative image steganography based on latent diffusion models. In *Proceedings of the 32nd ACM International Conference* on Multimedia, pp. 3001–3009, 2024.
- Rezaei, A., Akbari, M., Alvar, S. R., Fatemi, A., and Zhang, Y. Lawa: Using latent space for in-generation image watermarking. In *European Conference on Computer Vision*, pp. 118–136. Springer, 2025.
- Rombach, R., Blattmann, A., Lorenz, D., Esser, P., and Ommer, B. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- Song, J., Meng, C., and Ermon, S. Denoising diffusion implicit models. arXiv preprint arXiv:2010.02502, 2020a.
- Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., and Poole, B. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020b.
- Su, W., Ni, J., Hu, X., and Fridrich, J. Image steganography with symmetric embedding using gaussian markov random field model. *IEEE Transactions on Circuits and Systems for Video Technology*, 31(3):1001–1015, 2021. doi: 10.1109/TCSVT.2020.3001122.

- Su, W., Ni, J., and Sun, Y. Stegastylegan: Towards generic and practical generative image steganography. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(1):240–248, Mar. 2024. doi: 10.1609/ aaai.v38i1.27776. URL https://ojs.aaai.org/ index.php/AAAI/article/view/27776.
- Subramanian, N., Elharrouss, O., Al-Maadeed, S., and Bouridane, A. Image steganography: A review of the recent advances. *IEEE access*, 9:23409–23423, 2021.
- Tancik, M., Mildenhall, B., and Ng, R. Stegastamp: Invisible hyperlinks in physical photographs. In *Proceedings of* the IEEE/CVF conference on computer vision and pattern recognition, pp. 2117–2126, 2020.
- Utepbergenov, I., Kuandykova, J., Mussin, T., and Sagyndykova, S. Creating a program and research a cryptosystem on the basis of cardan grille. In 2013 Second International Conference on Informatics & Applications (ICIA), pp. 24–29. IEEE, 2013.
- van Schyndel, R., Tirkel, A., and Osborne, C. A digital watermark. In *Proceedings of 1st International Conference* on *Image Processing*, volume 2, pp. 86–90 vol.2, 1994. doi: 10.1109/ICIP.1994.413536.
- Wallace, B., Gokul, A., and Naik, N. Edict: Exact diffusion inversion via coupled transformations. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 22532–22541, 2023.
- Wang, F., Yin, H., Dong, Y., Zhu, H., Zhang, C., Zhao, H., Qian, H., and Li, C. Belm: Bidirectional explicit linear multi-step sampler for exact inversion in diffusion models. *arXiv preprint arXiv:2410.07273*, 2024.
- Wang, Y., Yang, X., and Liu, W. Generative image steganography based on digital cardan grille. In Security and Privacy in New Computing Environments: Third EAI International Conference, SPNCE 2020, Lyngby, Denmark, August 6-7, 2020, Proceedings 3, pp. 343–355. Springer, 2021.
- Wei, P., Luo, G., Song, Q., Zhang, X., Qian, Z., and Li, S. Generative steganographic flow. In 2022 IEEE International Conference on Multimedia and Expo (ICME), pp. 1–6, 2022. doi: 10.1109/ICME52920.2022.9859628.
- Weiming, Z. Provable secure steganography: Theory, application and prospects. Journal of Cybersecurity, 1(1):38, 2023. doi: null. URL https://www.journalofcybersec.com/ EN/Y2023/V1/I1/38.
- Xu, G., Wu, H.-Z., and Shi, Y.-Q. Structural design of convolutional neural networks for steganalysis. *IEEE Signal Processing Letters*, 23(5):708–712, 2016.

- Xu, Y., Mou, C., Hu, Y., Xie, J., and Zhang, J. Robust invertible image steganography. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7865–7874, 2022a. doi: 10.1109/CVPR52688.2022. 00772.
- Xu, Y., Mou, C., Hu, Y., Xie, J., and Zhang, J. Robust invertible image steganography. In 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 7865–7874, 2022b. doi: 10.1109/CVPR52688.2022. 00772.
- Yang, J., Ruan, D., Huang, J., Kang, X., and Shi, Y.-Q. An embedding cost learning framework using gan. *IEEE Transactions on Information Forensics and Security*, 15: 839–851, 2020. doi: 10.1109/TIFS.2019.2922229.
- Yang, K., Chen, K., Zhang, W., and Yu, N. Provably secure generative steganography based on autoregressive model. In Yoo, C. D., Shi, Y.-Q., Kim, H. J., Piva, A., and Kim, G. (eds.), *Digital Forensics and Watermarking*, pp. 55–68, Cham, 2019. Springer International Publishing. ISBN 978-3-030-11389-6.
- Yang, Y., Liu, Z., Jia, J., Gao, Z., Li, Y., Sun, W., Liu, X., and Zhai, G. Diffstega: towards universal training-free coverless image steganography with diffusion models. *arXiv preprint arXiv:2407.10459*, 2024a.
- Yang, Z., Chen, K., Zeng, K., Zhang, W., and Yu, N. Provably secure robust image steganography. *IEEE Transactions on Multimedia*, 26:5040–5053, 2024b. doi: 10.1109/TMM.2023.3330098.
- Yang, Z., Zeng, K., Chen, K., Fang, H., Zhang, W., and Yu, N. Gaussian shading: Provable performance-lossless image watermarking for diffusion models. In *Proceedings* of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 12162–12171, 2024c.
- Ye, J., Ni, J., and Yi, Y. Deep learning hierarchical representations for image steganalysis. *IEEE Transactions on Information Forensics and Security*, 12(11):2545–2557, 2017.
- You, W., Zhang, H., and Zhao, X. A siamese cnn for image steganalysis. *IEEE Transactions on Information Foren*sics and Security, 16:291–306, 2020.
- You, Z., Ying, Q., Li, S., Qian, Z., and Zhang, X. Image generation network for covert transmission in online social network. In *Proceedings of the 30th ACM International Conference on Multimedia*, MM '22, pp. 2834–2842, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450392037. doi: 10. 1145/3503161.3548139. URL https://doi.org/ 10.1145/3503161.3548139.

- Yu, F., Seff, A., Zhang, Y., Song, S., Funkhouser, T., and Xiao, J. Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop. arXiv preprint arXiv:1506.03365, 2015.
- Yu, J., Zhang, X., Xu, Y., and Zhang, J. Cross: diffusion model makes controllable, robust and secure image steganography. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*, NIPS '23, Red Hook, NY, USA, 2024. Curran Associates Inc.
- Zhang, G., Lewis, J. P., and Kleijn, W. B. Exact diffusion inversion via bidirectional integration approximation. In *European Conference on Computer Vision*, pp. 19–36. Springer, 2025.
- Zhang, K. A., Xu, L., Cuesta-Infante, A., and Veeramachaneni, K. Robust invisible video watermarking with attention. arXiv preprint arXiv:1909.01285, 2019a.
- Zhang, R., Zhu, F., Liu, J., and Liu, G. Depth-wise separable convolutions and multi-level pooling for an efficient spatial cnn-based steganalysis. *IEEE Transactions on Information Forensics and Security*, 15:1138–1150, 2019b.
- Zhang, Z., Liu, J., Ke, Y., Lei, Y., Li, J., Zhang, M., and Yang, X. Generative steganography by sampling. *IEEE access*, 7:118586–118597, 2019c.
- Zhou, Z., Dong, X., Meng, R., Wang, M., Yan, H., Yu, K., and Choo, K.-K. R. Generative steganography via auto-generation of semantic object contours. *IEEE Transactions on Information Forensics and Security*, 18:2751– 2765, 2023a.
- Zhou, Z., Su, Y., Li, J., Yu, K., Wu, Q. M. J., Fu, Z., and Shi, Y. Secret-to-image reversible transformation for generative steganography. *IEEE Transactions on Dependable* and Secure Computing, 20(5):4118–4134, 2023b. doi: 10.1109/TDSC.2022.3217661.
- Zhu, J., Kaplan, R., Johnson, J., and Fei-Fei, L. Hidden: Hiding data with deep networks. In Computer Vision – ECCV 2018: 15th European Conference, Munich, Germany, September 8-14, 2018, Proceedings, Part XV, pp. 682–697, Berlin, Heidelberg, 2018. Springer-Verlag. ISBN 978-3-030-01266-3. doi: 10.1007/ 978-3-030-01267-0\_40. URL https://doi.org/ 10.1007/978-3-030-01267-0\_40.

# A. Improved Algorithm for Initial Noise Generation

Analogous to the original algorithm, the enhanced method still uses each uniformly drawn binary bit (0 or 1) to decide whether its associated sample should lie in the left tail (x < -k) or the right tail (x > k) of the standard normal distribution. The key innovation is the adoption of vectorized rejection sampling for the central region: a large batch of standard normal variates is generated in a single operation, values falling within the interval (-k, k) are identified via Boolean masking, and those accepted samples are filled sequentially into all non-tail positions. By avoiding element-wise loops for random-number generation, this method greatly improves overall execution efficiency. For the tail regions, the algorithm similarly generates multiple candidate variates in batches for each tail position, filters the subset satisfying the prescribed left or right tail condition, and randomly selects one. If no candidate meets the criterion in a given batch, sampling continues iteratively until valid tail samples have been obtained for every position.

# **B.** Provable Security of MDDM

Assume that the distribution of carrier data is  $P_{\rm C}$ , the distribution of encrypted data is  $P_{\rm s}$ , and the relative entropy between carrier data and encrypted data is  $D(P_c || P_s)$ , which is defined as:

$$D(P_C || P_S) = \sum_{q \in \Omega} P_C(q) \lg \frac{P_C(q)}{P_S(q)}$$
(11)

If  $D(P_C || P_S) = 0$ , the steganography system is said to be absolutely secure (Cachin, 1998; Weiming, 2023).

First, since the hidden message is not in the pixels of the image, it is impossible to find the hidden message through steganalysis in the pixels. Second, since the hidden message does not change the distribution of the noise, it will not cause any characteristic difference between the stego images and other naturally generated images. To sum up, MDDM is absolutely safe.

# C. Privacy Security of Cardan Grille in MDDM

In MDDM, the privacy security of the Cardan Grille are closely related to the steganographic communication protocol. Specifically, if a new Cardan Grille is used for each image, the distribution of the reversed noise for each image will be different, and third parties will not have sufficient information to infer whether an image contains hidden information. However, when multiple images share the same Cardan Grille for steganography, the distribution patterns of their reversed noise will exhibit a certain degree of similarity. Therefore, we evaluated the privacy security of the Cardan Grille steganographic method under the scenario where the grille is used for only one communication and then reused in subsequent transmissions, across different message lengths. The diffusion model employed is Stable Diffusion v2.1, with image resolution set to  $512 \times 512$  and an empty prompt. We assume that an attacker has access to the stego images and is aware of the diffusion model, the steganographic method, and even the number of Cardan Grille positions *l* (i.e., the message length). The attacker's objective is to infer the position distribution and the ordering of the Cardan Grille.

The attacker extracts the inverted noise from multiple images and filters out the elements with larger absolute values, as positions with larger absolute values are more likely to correspond to locations used by the Cardan Grille. Two attack strategies are considered:

- Union Strategy: For each image, the attacker selects the *l* position elements with the largest absolute values from the inverted noise vector. The positions chosen across all images are then combined via a union operation to obtain the final set of selected locations.
- Top Strategy: For each image, the attacker selects the *l* position elements with the largest absolute values from the inverted noise vector. Subsequently, among all the selected positions from all images, the *l* positions that occur most frequently are chosen as the final set of selected locations.

The experimental results are shown in Figure 8. "Union" refers to the position hit rate of the first strategy; "Union Redundancy" represents the redundancy of the union set from the first strategy (calculated as the ratio of the number of Cardan Grille positions to the number of elements in the union set — lower values indicate higher redundancy); and "Top"

refers to the position hit rate of the second strategy. A high position hit rate in the Union strategy combined with low Union Redundancy indicates that the attacker has inferred a significant portion of the Cardan Grille's position distribution.

The experimental results show that when the Cardan Grille is used only once or reused up to two times, neither strategy poses a significant threat to the position distribution or ordering of the grille. When the grille is reused 3 to 5 times, the Union strategy achieves a higher hit rate, but the redundancy of the position set rapidly decreases, limiting the threat it poses to the grille's position distribution. The Top strategy also sees a modest improvement in hit rate, but neither strategy is able to compromise the ordering of the grille positions. When the grille is reused more than 5 times, the hit rate of the Union strategy shows no substantial improvement, while the position set redundancy continues to decline, paradoxically reducing its threat to the grille's positional information. The Top strategy shows only limited further improvement in hit rate, and neither strategy is ever able to compromise the ordering information of the Cardan Grille.

Since the Cardan Grille is ordered, the brute-force computational complexity of guessing the order of an n-position Cardan Grille is O(l!), where

$$l! \approx \sqrt{2\pi l} \left(\frac{l}{e}\right)^l,\tag{12}$$

according to Stirling's approximation. Since the difficulty of brute-forcing the ordering of the Cardan Grille increases exponentially with the number of positions, a more secure strategy is to use longer messages. We therefore recommend that, in steganographic scenarios where the Cardan Grille is reused, each image should carry at least 50 bits of hidden information.



*Figure 8.* Simulated attacks against the MDDM that employs repeatedly used Cardan Grille under varying message lengths based on Stable Diffusion v2.1. The percentages in the chart represent the proportion of simulated attack hit positions rather than the proportion of hit orders.

# **D.** Experiments

### **D.1.** Quantitative Evaluation of Image Quality

We compare MDDM with two GAN-based generative steganography models, StegaStyleGAN (Su et al., 2024) and CIS-Net (You et al., 2022), on the CelebA dataset (Liu et al., 2015). The results are shown in Table 3. Due to image size limitations of CIS-Net, all comparisons are conducted on images of size  $32 \times 32$ . Moreover, considering the instability often associated with GAN training, and to avoid potential bias from our own reimplementation, we cite the results of StegaStyleGAN and CIS-Net directly from their original papers.

To quantitatively evaluate the image quality of our proposed method, we employ the Fréchet Inception Distance (FID), a widely used metric that assesses the similarity between the distributions of generated and real images. It is worth noting that FID scores are influenced by various factors, such as image compression and cropping.

Experimental results show that our method achieves competitive performance not only in terms of steganography capacity and accuracy but also in image quality, with FID scores comparable to those of existing methods. This indicates good generalization ability. Unlike GAN-based steganographic models, MDDM does not require additional training. Consequently, the image quality depends on the generative quality of the original diffusion model, and using more advanced diffusion models further improves the image quality produced by MDDM.

Method	Size	Capacity (Bpp)*	Accuracy	FID
CIS-Net (You et al., 2022)	$32 \times 32$	0.03	0.98	6.20
StegaStyleGAN (Su et al., 2024)	$32 \times 32$	0.06	1.00	3.74
MDDM (Ours)	$32 \times 32$	0.08	0.99	11.71
MDDM-U3 (Ours)	$32 \times 32$	0.08	1.00	-

Table 3. Comparison of our method and baseline GAN methods.

<sup>\*</sup> Bpp means bits per pixel.

# **D.2.** Ablation Experiment Setup

We evaluate unconditional diffusion on facial images from the Bedroom and Cat datasets (Yu et al., 2015) and FFHQ (Karras et al., 2019). As described in Pulsar (Jois et al., 2024), MDDM is also a symmetric key scheme. The sender can easily know the accuracy of the information extracted from the generated image and can terminate and regenerate. A key advantage of our method is that when the two parties communicate and determine the Cardan Grille, the same secret message can be generated multiple times using different seeds without affecting the reception (see Section D.4 for details). Accordingly, we apply the MDDM-U3 optimization strategy for unconditional pixel-space diffusion to mitigate potential instabilities (see Section 3.1) by selecting, from three consecutive samples, the image with the smallest information loss. This is completely feasible in practical applications, because if the improved algorithm in Appendix A is used, each unconditional generation and extraction for a  $256 \times 256$  image takes only about 10 seconds on an NVIDIA RTX 3090.

Due to the lack of standard datasets corresponding to the above categories for evaluating conditional generation, we adopt common prompts corresponding to the above categories. The prompts for Face, Bedroom, and Cat are "Portrait photo, best quality, masterpiece, ultra detailed, UHD 4K, photographic, 1girl, face, looking at viewer, color photo, natural colors", "A photo of the bedroom", and "A cat", respectively.

# D.3. Quantitative Results Compared to Other Image Steganography Methods

Table 4 shows the quantitative comparison results of the steganography abilities. Since neither StegaDDPM nor LDStega has made the experimental code or seed public, the data of StegaDDPM and LDStega are from the data in the paper LDStega (Peng et al., 2024) for reference. For Pulsar, we utilize its open-source code to conduct tests. The specific experimental settings of MDDM-U3 are described in Appendix D.2. Using MDDM-U3 can be regarded as comparing the accuracy after each seed sharing.

# **D.4.** Diversity of MDDM

Table 5 shows the comparison of our method with other methods in terms of diversity, where diversity refers to the number of images that can be generated after a single communication of a shared key. Figure 9 shows a visualization of the generation effect of MDDM in various cases.



*Figure 9.* Examples of images generated by Stable Diffusion v2.1 for various combinations of Cardan Grille, messages, and prompts (**CG** refers to the Cardan Grille and **Msg** refers to the message). All images are  $512 \times 512$  pixels with a hidden-message length of 1024 bits. Message-extraction accuracy for every case is 100%.

Table 4. Comparison of extraction accuracy and capacity between MDDM and baseline methods.

Method	Method Type	Sizo	Capacity	Accuracy		
		5120	(Bits)	FFHQ	Bedroom	Cat
Pulsar (Jois et al., 2024)		256  imes 256	$\approx 4500$	97.00%	93.00%	96.00%
StegaDDPM (Peng et al., 2023)	CIC	256  imes 256	4096	93.45%	90.19%	90.81%
LDStega (Peng et al., 2024)	(Diffusion)	256  imes 256	4096	98.65%	<b>98.50</b> %	98.48%
MDDM (Ours)	(Diffusion)	256  imes 256	4096	94.00%	92.46%	92.68%
MDDM-U3 (Ours)		$256\times256$	4096	<b>98.81</b> %	96.32%	<b>99.87</b> %

Table 5. Comparison of the diversity of MDDM and baseline methods.

Method	Hiding Type	Diffusion Type	Diversity (Number of Images)
StegaDDPM (Peng et al., 2023)	Binary	Pixel-Space	1
Pulsar (Jois et al., 2024)	Binary	Pixel-Space	1
LDStega (Peng et al., 2024)	Binary	Latent	1
CRoSS (Yu et al., 2024)	Image	Latent	Several*
MDDM (Ours)	Binary	Pixel-Space + Latent	$\infty$

\* CRoSS can alter individual elements in the secret image but has difficulty changing the overall contents of the image.