

Towards Blind Watermarking: Combining Invertible and Non-invertible Mechanisms

Rui Ma
rui_m@stu.pku.edu.cn
Peking University
Haidian, Beijing, China

Fan Yang
fyang.eecs@pku.edu.cn
Peking University
Haidian, Beijing, China

Mengxi Guo
guomengxi.qoelab@bytedance.com
Bytedance Inc.
Shanghai, China

Yuan Li*
yuanli@pku.edu.cn
Peking University
Haidian, Beijing, China

Xiaodong Xie
donxie@pku.edu.cn
Peking University
Haidian, Beijing, China

Yi Hou
yihou@pku.edu.cn
Peking University
Haidian, Beijing, China

Huizhu Jia
hzjia@pku.edu.cn
Peking University
Haidian, Beijing, China

ABSTRACT

Blind watermarking provides powerful evidence for copyright protection, image authentication, and tampering identification. However, it remains a challenge to design a watermarking model with high imperceptibility and robustness against strong noise attacks. To resolve this issue, we present a framework Combining the Invertible and Non-invertible (CIN) mechanisms. The CIN is composed of the invertible part to achieve high imperceptibility and the non-invertible part to strengthen the robustness against strong noise attacks. For the invertible part, we develop a diffusion and extraction module (DEM) and a fusion and split module (FSM) to embed and extract watermarks symmetrically in an invertible way. For the non-invertible part, we introduce a non-invertible attention-based module (NIAM) and the noise-specific selection module (NSM) to solve the asymmetric extraction under a strong noise attack. Extensive experiments demonstrate that our framework outperforms the current state-of-the-art methods of imperceptibility and robustness significantly. Our framework can achieve an average of 99.99% accuracy and 67.66 dB PSNR under noise-free conditions, while 96.64% and 39.28 dB combined strong noise attacks. The code will be available in <https://github.com/rmpku/CIN>.

CCS CONCEPTS

• Security and privacy → Digital rights management.

KEYWORDS

Robust blind watermarking; Invertible network

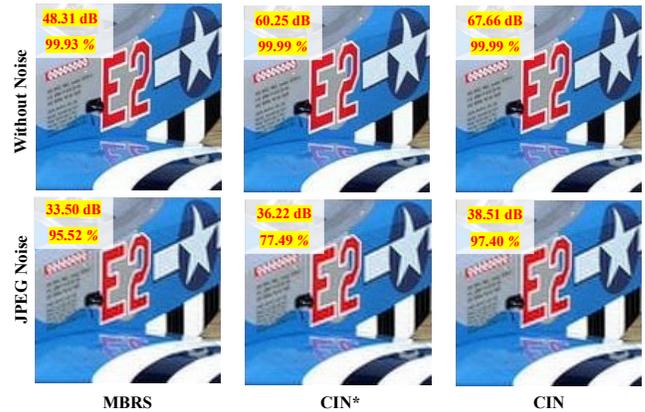


Figure 1: The top and bottom are the watermarked images with noise-free and *Jpeg* noise (training with combined noise). From left to right are the MBRS [21], the invertible-only baseline CIN*, and the proposed CIN, respectively. The red bars in the picture are PSNR with the input image and the accuracy *Acc* of the extracted watermark, respectively.

1 INTRODUCTION

Digital watermarking utilizes data concealment techniques to embed some form of identification into a digital medium that can be transmitted together and authenticated by the property owner. Watermarking has the characteristics that the information embedding should be robust, tamper-resistant, and for authentication [7, 47]. HiDDeN [50] is the first watermarking framework that enabled end-to-end training, and numbers of works are subsequently derived, which can be simply classified as CNN-based [21, 26, 32, 42, 50] and GAN-based [46, 48, 49]. The end-to-end joint training of the models enabled the incorporation of the embedding and extraction efficiently and ensured the effectiveness of the pipeline. As shown in the top part of Fig.2, the key to guaranteeing robustness is the adversarial training with the differential noise layer. There are some limitations in the end-to-end framework. The decoder

*Corresponding authors.

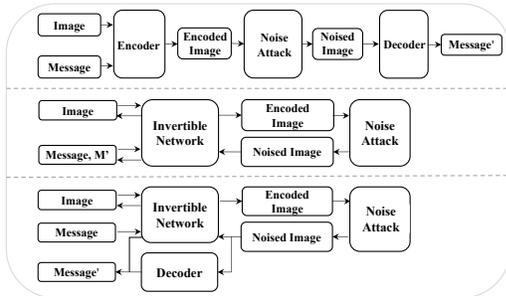


Figure 2: Framework of digital watermarking model. Top is the End-to-end method. Middle is the baseline CIN^{*}. Bottom is the proposed CIN.

and its latent variables are approximately likelihood evaluation inferred by data, which means the entire training objective is not an exact form. And if the model contains a Bottleneck structure, such as in auto-encoder based watermarking [22], the manipulation of the features will result in non-invertible losses of information that is detrimental to the watermark restoration. In addition, more information is uncontrollably removed when watermarked images are subject to noise attacks, which leads to the inevitable sacrifice of imperceptibility to improve robustness.

We propose a framework combining the invertible and non-invertible mechanisms, as shown the bottom of Fig.2. For the invertible part $f_{\theta_1}(\cdot)$, we introduce an invertible neural network (INN) based module that significantly improves the imperceptibility of the watermark and is robust to common additive noise. Denoting the input image and watermarking as I_m and W_m , and $f_{\theta_1}(I_m, W_m)$ as z . The inverse function $f_{\theta_1}^{-1}(\cdot)$ can be trivially obtained, such that $P(W_m)$ can be easily sampled with $W_m = f_{\theta_1}^{-1}(z)$, where $z \sim P_Z(z)$ and $P(W_m)$ refers to the probability distribution of the watermark. In a generic INN, the probability distribution of $P_Z(z)$ can be explicitly defined as a Gaussian prior since z poses no additional limitation [19], and yet we define it as the input image-based prior, which helps to reduce the introduction of errors destroying reversibility and helps to stabilize the overall training process [9]. Benefiting from the invertibility, the probability distribution of the latent variable in the inverse process is a full posterior probability, ensuring the accuracy of the restored watermark. Since the property that the INN shares a set of parameters for both the embedding and the extraction processes, we can train and learn the forward embedding process and obtain the inverse extraction process "for free", as opposed to the end-to-end, in which the decoder has a separate train and learn process.

For the high imperceptibility of the watermark overlaid with the input image, the distribution z should be numerically small. As shown in Fig. 3, since $P_Z(z)$ is more variable when subjected to lossy compression noise leading to a more fragile watermark. And the property of sharing parameters between the embedding and extraction process of the invertible module leads to the extraction process being only sampled according to $P_Z(z)$ of the embedding, which also limits the decoder's ability to adapt to the strong and

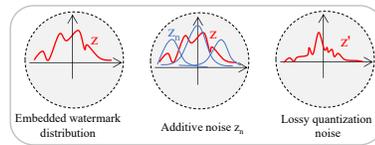


Figure 3: Noise description. Although the INN can fit the distribution z_n of additive noise, the lossy quantization noise will change the embedded watermark distribution.

non-differentiable noise. For the non-invertible part $f_{\theta_2}(\cdot)$, we introduce a non-invertible attention-based module (NIAM) and the noise-specific selection module (NSM) to solve the asymmetric extraction of watermarks under a lossy compression noise. The distribution of the noised image is denoted as $z' \sim P_{Z'}(z')$, and we expect using the NIAM to approximate $W_m \approx f_{\theta_2}^{-1}(z')$ from $P_{Z'}(z')$. We introduce approximately differentiable and non-differentiable compression noise in the training step to enable NIAM to guide the encoder as well. The gradients are backward to the encoder when the differentiable noise is selected in the noise pool. In contrast, only the NIAM is updated when the non-differentiable noise is chosen. Therefore, we can effectively combine invertible and non-invertible modules in digital watermarking.

The contributions of this paper can be summarized as follows:

1. To the best of our knowledge, we are the first to incorporate an INN with blind watermarking, while most of the existing deep learning-based watermarking approaches focus on encoder-decoder pipeline or adversarial training.
2. To compensate for the deficiency of the INN in combating quantization loss noise, we introduce NIAM as a parallel decoder to improve the robustness of the model against compression.
3. We propose the diffusion and extraction module (DEM) and the fusion and split module (FSM) for more efficient and robust embedding and extraction of watermarks.
4. We conduct extensive experiments on various image datasets and compare our approach against the state-of-the-art watermarking methods. Our method achieves excellent performance in terms of imperceptibility and robustness.

2 RELATED WORK

2.1 Watermarking

The research on digital watermark is first proposed in [40] in 1994, and it can be generally classified into two categories: traditional algorithms based on transform domain and a deep learning-based approach driven by data. Where the traditional watermarking methods include algorithms based on singular value decomposition [31, 37, 38], moment-based watermarking algorithms [17, 18] and transform domain watermarking algorithms [3, 15, 33]. Accordingly, the deep learning-based watermarking model is first introduced by Hamidi et al. [22] in 2017, whose method brings superior imperceptibility and robustness over traditional methods by employing an auto-encoder convolutional neural network (CNN). HiDDeN [50] is the first to introduce the adversarial network to blind watermarking and also the first end-to-end method using neural networks. Subsequently, Ahmadi et al. [2] propose a digital watermarking

framework based on residual networks and achieved excellent robustness and imperceptibility. Liu et al. [26] propose the TSDL framework, which is composed of two-stage: noise-free end-to-end adversary training and noise-aware decoder-only training. This method is effective against black-box noise and can introduce non-differentiable noise attacks in the end-to-end network. Soon, Jia et al. [21] propose a novel Mini-Batch of Simulated and Real *Jpeg* compression method to enhance robustness against *Jpeg* compression, which performs excellent performance in various noises. In addition, there are works studying deep learning-based steganography and encryption [12, 28, 35, 43, 49] and a review of research on deep learning-based watermarking and steganography can be found [7, 47].

2.2 Invertible Neural Network

Invertible neural network is the first learning-based normalizing flow framework for modeling complex high-dimensional densities, proposed by Dinh et al. [10] in 2014. To improve the efficiency and performance of image processing, Dinh et al. [11] introduce convolutional layers in coupling models by modifying the additive coupling layers to both multiplication and addition, called real NVP. To further improve the coupling layer for density estimation and achieve better generation results, Kingma et al. [24] propose a novel generative flow model based on the ActNorm layer and generalize channel-shuffle operations with invertible 1×1 convolutions. Real NVP and 1×1 convolutions are two frequently used structures in image tasks employing INN. Normalizing flow-based INN has become a popular choice in image generation tasks, and evolved with various similar deformations [6, 8, 13, 20, 24]. And there are also several approaches that incorporate INN with other methods, such as INN combined with self-attention [16] and INN constructed with masked convolutions [36]. Due to the flexibility and effectiveness of INN, it is also used for image super-resolution [29, 44] and video super-resolution[51]. In addition, Ouderaa et al.[39] applied INN to image-to-image translation, Wang et al.[41] applied INN in digital image compression, Ardiszone et al. [5] introduce conditional INN for colorization, Liu et al. [27] propose an invertible denoising network, Xing et al. [45] propose an invertible image signal processing and Pumarola et al. [34] apply INN for image and 3D point cloud generation.

3 METHOD

3.1 Overall Architecture

Fig. 5 shows the architecture of our proposed *CIN*, which is divided into the following parts: a Diffusion and Extraction Module (DEM), an Invertible Module (IM), a Fusion and split Module (FSM), a Non-invertible Attention-based Module (NIAM), and Noise-specific Selection Module (NSM). The embedding and extraction of *CIN* are defined as $f_{CIN}(\theta)$ and $f_{CIN}^{-1}(\theta)$.

3.2 Diffusion and Extraction Module

The watermark W_m is a binary sequence of length L randomly sampled from $W_m \sim \{0, 1\}^L$. Embedding the watermark into an RGB image I_m with length and width of H and W , respectively. Then the watermark and the image of the input model are $W_m \in \mathbb{R}^{B \times L}$

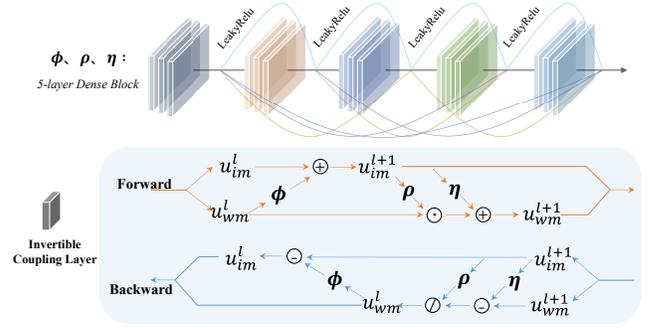


Figure 4: The structure diagram of the invertible coupling layer. Top is the functions, ϕ , ρ and η constructed with 5-layer dense block, of the invertible coupling layer. Bottom is the exact form of the affine coupling layer.

and $I_m \in \mathbb{R}^{B \times C \times H \times W}$ respectively, where B and C refer to batchsize and channel number.

As shown in fig. 6, the top and bottom parts show the diffusion and extraction processes, respectively. In the diffusion processing, to align the watermark with the number of channels of the image, we first replicate the watermark W_m in three copies. The different fully connected (FC) branches produce redundant watermarks of longer length, respectively. Subsequently, reshape and upsample to the same scale size as the cover image by two-dimensional transpose convolution (ConvT). After passing through the FC layer, the watermark length \hat{L} is 256, the kernel size and stride of ConvT are both 2, and the block number is 3. Finally, the output of the three branches is concatenated and fed into the invertible module after the Haar transform. For forward embedding operations:

$$\Psi_{DEM} = \Gamma_{haar}(\text{Ocat}(\Gamma_{convT}(\Gamma_{fc}(\text{Ocopy}(W_m)))))) \quad (1)$$

where $\text{Ocopy} \in 3 \times \mathbb{R}^{B \times L}$, $\Gamma_{fc} \in \mathbb{R}^{B \times \hat{L}}$, $\Gamma_{convT} \in \mathbb{R}^{B \times 1 \times H \times W}$, $\text{Ocat} \in \mathbb{R}^{B \times 3 \times H \times W}$ and $\Gamma_{haar} \in \mathbb{R}^{B \times 12 \times H/2 \times W/2}$ refer to operations Copy, FC, ConvT, Concatenate and Haar Transform, respectively. And Ψ_{DEM} is the output tensor of Diffusion and Extract Module.

In the extraction process, the operation Ψ_{DEM}^{-1} , which is the opposite of the embedding process, is taken for extraction. In contrast to Copy in the watermark embedding step, the final result is output by Average Pooling in the extraction process. The formula is as follows:

$$W_{m_1} = \Psi_{DEM}^{-1}(\cdot) \quad (2)$$

3.3 Invertible Module

The coupling layer in the IM is an additive affine transformation, which was first proposed in NICE [10]. Recently, invertible architecture has been applied to information hiding with excellent representational capacity in works [14, 28, 44], from which we were inspired. We use the watermark W_m and the image I_m as the two inputs of the invertible module, respectively. Our goal is to embed the W_m into the I_m with excellent imperceptibility and robustness.

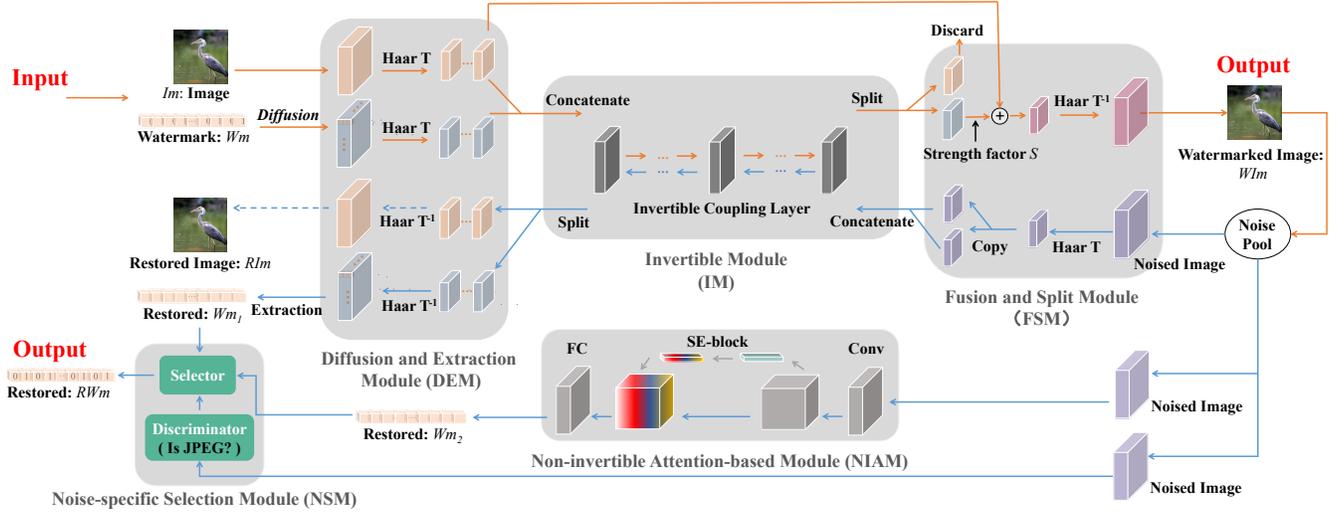


Figure 5: Overall model architecture. The DEM diffuses the watermark to the same dimension as the image using FC, Convolution, and Haar Transform. IM maps the diffused watermark to an embeddable distribution. FSM scales the watermark to be embedded and stacks it with the input image in the frequency domain. The noise pool introduces a variety of traditional noises. NIAM is used to enhance the robustness against lossy compression noise. NSM is used to output the best result of IM and NIAM.

The invertible module is shown in Fig. 4. The embedding and extraction correspond to the forward and backward of the bijection structure [44], respectively. In the coupling layer of the l^{th} , u_{wm} and u_{im} denote the input watermark and image, respectively. The corresponding u_{wm}^{l+1} and u_{im}^{l+1} denote the output watermark and image after passing through the current coupling layer. The invertible module is formulated as:

$$\mathbf{u}_{im}^{l+1} = \phi(\mathbf{u}_{wm}^l) + \mathbf{u}_{im}^l \quad (3)$$

$$\mathbf{u}_{wm}^{l+1} = \mathbf{u}_{wm}^l \odot \exp\left(\rho\left(\mathbf{u}_{im}^{l+1}\right)\right) + \eta\left(\mathbf{u}_{im}^{l+1}\right) \quad (4)$$

where $\exp(\cdot)$ is exponential operator, $\phi(\cdot)$, $\rho(\cdot)$ and $\eta(\cdot)$ are arbitrary functions, \odot is the Hadamard product. The corresponding backward propagation of the extraction process is formulated as:

$$\mathbf{u}_{im}^l = \mathbf{u}_{im}^{l+1} - \phi\left(\mathbf{u}_{wm}^l\right) \quad (5)$$

$$\mathbf{u}_{wm}^l = \left(\mathbf{u}_{wm}^{l+1} - \eta\left(\mathbf{u}_{im}^{l+1}\right)\right) \odot \exp\left(-\rho\left(\mathbf{u}_{im}^{l+1}\right)\right) \quad (6)$$

3.4 Fusion and Split Module

The output $\Psi_{inv} \in \mathbb{R}^{B \times 24 \times H/2 \times W/2}$ of the invertible module (IM) can be split into two parts, $\Psi_{inv}^I(x; LR, HR)$ and $\Psi_{inv}^{II}(x; LR_-, HR_-) \in \mathbb{R}^{B \times 12 \times H/2 \times W/2}$. After Haar transform, the LR and HR represent the image's low and high-frequency components. The left part of Fig. 7 is similar to the Channel Squeeze module of work [9]. The corresponding channels of the two outputs of the IM are averaged, which can be fused and squeezed to the size of the image. It is, however, difficult to trade off the watermark robustness with the imperceptibility. Therefore, we propose the fusion method as shown on the right in Fig. 7.

In the embedding process, we discard the image part of the IM output, keep only the mapped watermark part, and then add it to the image after scaling by the strength factor S to obtain the final watermarked image. The formula is as follows:

$$WIm = \Gamma_{haar}^{-1}\left(\Psi_{inv}^{II}(x; LR_-, HR_-) \times S + \Gamma_{haar}(I_m(x; LR, HR))\right) \quad (7)$$

where S is the strength of the watermark. To restore the embedded watermark by invertible branch, the inputs are

$$\Psi_{DEM} = \text{Ocat}\left(\text{Ocopy}\left(\Gamma_{haar}^{-1}(WIm)\right)\right) \quad (8)$$

where $\Gamma_{haar}^{-1}(\cdot)$ is inverse Haar transform.

3.5 Non-invertible Module

The embedding and extraction of the watermark in invertible networks has a deterministic mapping relationship, which makes excellent results for watermark extraction accuracy in scenes without or with additive noise. However, when subjected to lossy compression or complex non-additive noise, since the forward and backward of the invertible network share the same set of parameters, the parameters of the decoder are updated along with the encoder, which limits the ability of the decoder to cope with complex noise. Therefore, an additional decoder is introduced in our framework to enhance the robustness of the invertible module against non-differentiable noise attacks, such as lossy compression noise. The non-invertible module uses SENet as the backbone to extract the watermark information Wm_2 :

$$Wm_2 = \Gamma_{fc}\left(\Phi_{SE}\left(\Gamma_{conv}(\cdot)\right)\right) \quad (9)$$

where $\Gamma_{fc}(\cdot)$, $\Phi_{SE}(\cdot)$, and $\Gamma_{conv}(\cdot)$ are FC, SENet and convolution layer, respectively.

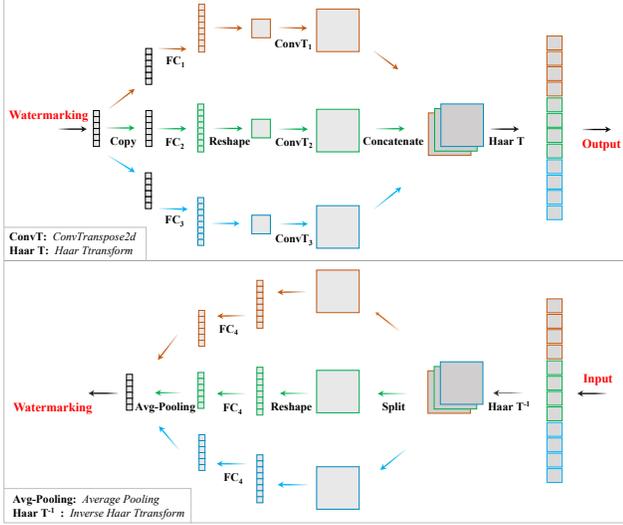


Figure 6: Diffusion and Extraction Module. Top: watermark embedding. Bottom: watermark extraction.

Inspired by the article [21], we introduce differentiable and non-differentiable compression noise into the noise pool to push the NIAM robustness by encountering lossy compression. The network can update the parameters of both the IM and NIAM when introducing differentiable compression attack and only the NIAM for non-differentiable noise. For NSM, we employ a CNN-based noise discriminator to distinguish whether the attack is *Jpeg* noise or not. If it is *Jpeg*, the selector exports Wm_2 extracted by NIAM; otherwise, return Wm_1 decoded by IM.

3.6 Noise Pool

The robustness of the watermark is improved by introducing a noise layer in the architecture in [26, 30, 50] et al.. To optimize the network parameters against the noise attack, it is generally necessary to use a differentiable noise layer trained jointly with the other basic module. In this work, the following 14 types of common noises:

$$N_{pool} = \{Identity, JpegMask, RealJpeg, Crop, Cropout, Resize, GaussianBlur, SaltSalt\&Pepper, GaussianNoise, Dropout, Brightness, Contrast, Saturation, Hue\}$$

where *JpegMask* is the simulated differentiable *Jpeg* noise [21]. For resisting specific noise *RealJpeg*, we employ the following noise pool to train the model:

$$N_{pool}^{c_j} = \{JpegMask, RealJpeg\}$$

To test the robustness of the model against simultaneous superimposed attacks of multiple noises, we use the following noise pool:

$$N_{pool}^{s_i} = \{Identity, Cropout, Resize, Saturation, Hue, Dropout, GaussianBlur, SaltSalt\&Pepper, GaussianNoise\}$$

For comparison with other works:

$$N_{pool}^{c_{p1}} = \{Identity, RealJpeg, Dropout, Cropout, Resize\}$$

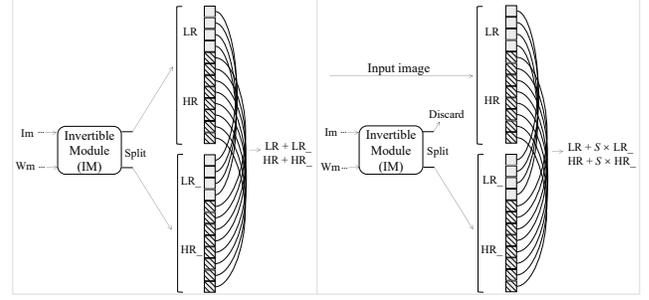


Figure 7: Left: fusion to the average value of the corresponding image channels. Right: fusion with the input image and discards the corresponded image part.

$$N_{pool}^{c_{p2}} = \{Identity, RealJpeg, Crop, Cropout, GaussianBlur, Dropout\}$$

3.7 Loss Functions

The loss functions constrain two parts: the watermarked image and the extracted watermark. Since the INN shares parameters for the embedding and extraction and has the same input and output dimensions, the loss constraint on the restored image in the noise-free version can also accelerate the convergence [4].

Watermarked Image We employ L_2 loss to guide the watermarked image WI_m to be visually alike to the reference image I_m :

$$\mathcal{L}_{WI_m} = \|I_m - WI_m\|_2^2 = \|I_m - f_{CIN}(\theta, I_m, W_m)\|_2^2 \quad (10)$$

Restored Watermark Calculate the L_2 distance for each pair of input watermark W_m and the extracted watermark RW_m :

$$\mathcal{L}_{RW_m} = \|W_m - RW_m\|_2^2 = \|W_m - f_{CIN}^{-1}(\theta, N_{pool}(WI_m))\|_2^2 \quad (11)$$

Restored Image When training the *CIN* with *Identity* noise layer, we employ L_2 distance to constrain the difference between the restored image RI_m and the reference image I_m :

$$\mathcal{L}_{RI_m} = \|I_m - RI_m\|_2^2 = \|I_m - f_{CIN}^{-1}(\theta, N_{pool}(WI_m))\|_2^2 \quad (12)$$

Total Loss To sum up, our *CIN* is optimized by minimizing the compact loss \mathcal{L}_{total} , with the corresponding weight coefficients λ_{WI_m} , λ_{RW_m} and λ_{RI_m} :

$$\mathcal{L}_{total} = \lambda_{WI_m} \mathcal{L}_{WI_m} + \lambda_{RW_m} \mathcal{L}_{RW_m} + \lambda_{RI_m} \mathcal{L}_{RI_m} \quad (13)$$

4 EXPERIMENTS

4.1 Baseline

The baseline model (denote as *CIN**) contains only the invertible part. Using the method proposed in the article [50] to concatenate each bit watermark after duplication with the image channels. And the channel squeezing method proposed in the article [9] is used to output the watermarked image, as shown in the left part of Fig. 7. The specific architecture of *CIN** is given in the Appendix.

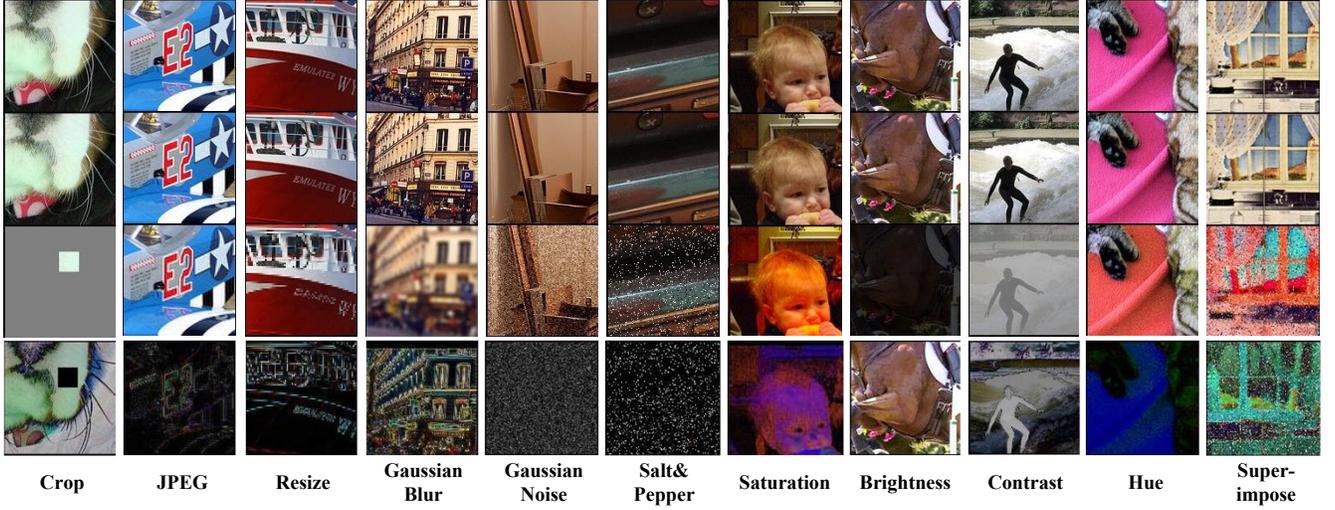


Figure 8: Visual comparisons of the experimental results under different traditional noises. Each column corresponds to a type of noise. Top: input image I_m . Second row: watermarked image WI_m . Third row: noised image I_m^{noised} . Bottom: the magnified difference $|I_m^{noised} - WI_m|$.

The results of the reference models are from either the published results or the open-source works and are partially quoted from the article [48]. Our experimental setup is consistent with the reference method.

4.2 Datasets.

To verify the robustness and imperceptibility of the proposed *CIN*, we utilize the real-world acquired COCO dataset [25] for training and evaluation. We also evaluate the transform performance of the model on the super-resolution dataset DIV2K dataset [1]. For the COCO dataset, 10000 images are collected for training, and evaluation is performed on the other 5000 images. For the DIV2K dataset, we use 100 images from the validation set for evaluation. For each input image, there is a corresponding watermarking message which is randomly sampled from the binary distribution $W_m \sim \{0, 1\}^L$.

4.3 Evaluation Metrics

To objectively evaluate the robustness and imperceptibility of our proposed watermarking framework, we apply a series of quantitative metrics. To validate the robustness, we evaluate the accuracy (*Acc*) between the extracted RW_m and the embedded W_m . For each input image $I_m(x_i)$, its corresponding watermark embedded and restored are $W_m(x_i)$ and $RW_m(x_i)$, respectively. Bit Error Ratio (*BER*) is also listed below, and the corresponding *Acc*(%) is $(1 - BER)$.

$$BER(\%) = \left(\frac{1}{L} \times \sum_{k=1}^L |RW_m(x_i) - W_m(x_i)| \right) \times 100\% \quad (14)$$

For imperceptibility of watermarked images, we adopt the Peak Signal-to-Noise Ratio (*PSNR*) and Structural Similarity (*SSIM*) for evaluation.

$$PSNR(I_m(x_i), RI_m(x_i)) = 20 \times \log_{10} \frac{MAX(I_m(x_i), RI_m(x_i)) - 1}{MSE(I_m(x_i), RI_m(x_i))} \quad (15)$$

$$SSIM(I_m(x_i), RI_m(x_i)) = \frac{(2\mu_x\mu_y + C_1)(2\sigma_{xy} + C_2)}{(\mu_x^2 + \mu_y^2 + C_1)(\sigma_x^2 + \sigma_y^2 + C_2)} \quad (16)$$

where $MAX(\cdot)$ is the maximum pixel value of images, and $MSE(\cdot)$ represents the Mean Squared Error. Symbol σ , μ_x and σ_{xy} represent the average, variances and covariance of images, respectively. C_1 and C_2 are two constants for preventing unstable results.

4.4 Implementation Details

To keep a fair comparison, we adopt exactly the same settings with the reference methods. Images are resized to 128×128 for all models, and the watermark length is 30 or 64. For our model, training with Nvidia 3080 graphics cards, the batch size is set to 32, and the Adam optimizer [23] with default hyperparameters is adopted. In the implementation, we train and evaluate the model under Specific Noise and Combined Noise, respectively. In the Specific Noise, all training and evaluations are performed only for one noise. In the Combined Noise, each mini-batch randomly samples a specific noise from N_{pool} , N_{pool}^{cp1} or N_{pool}^{cp2} . In the evaluation stage, we utilize the trained model to test the performance of each noise in turn. Throughout the training phase, we first trained the model in the noise-free case, at which point the loss weights are set to $\lambda_{WI_m} = 1$, $\lambda_{RW_m} = 0.001$ and $\lambda_{RI_m} = 1$, respectively. Next, the model is trained to resist different noise. We load the trained noise-free model and subsequently set the loss weights to $\lambda_{WI_m} = 1$, $\lambda_{RW_m} = 0.01$ and $\lambda_{RI_m} = 0$, respectively. In training combined noise N_{pool} , N_{pool}^{cp1} and N_{pool}^{cp2} , the loss weight are set to $\lambda_{WI_m} = 1$,

Table 1: Results of robustness and imperceptibility against various distortions. $PSNR_1$ and $PSNR_2$ denote the similarity between input image I_m and watermarked image WI_m , WI_m and noised image I_m^{noised} , respectively. Pre shows pre-trained accuracy without noise attack during the training stage. Specified and Combined mean the performance of specified noise and combined noises N_{pool} , respectively.

Noise	Factor	Specified				Combined	
		$PSNR_1$ dB	$PSNR_2$ dB	Pre (%)	Acc (%)	$PSNR_1$ dB	Acc (%)
<i>Identity</i>	-	67.66	-	-	99.99	39.28	99.99
<i>Dropout</i>	$p = 30\%$	61.39	62.93	99.43	99.99	39.29	99.99
<i>GausBlur</i>	$k = 7$	52.44	21.64	50.21	99.94	39.28	99.99
<i>Resize</i>	$p = 50\%$	53.56	21.06	59.10	99.97	39.29	99.99
<i>GausNoise</i>	$\sigma = 25$	61.50	18.65	99.90	99.99	39.29	99.99
<i>SaltPepper</i>	$p = 10\%$	53.83	14.80	81.24	99.96	39.29	99.99
<i>Cropout</i>	$p = 30\%$	62.30	62.65	99.90	99.99	39.29	99.99
<i>Crop</i>	$p = 3.5\%$	41.62	11.06	60.32	99.70	39.29	99.94
<i>RealJpeg</i>	$Q = 50$	42.70	27.13	50.32	99.11	39.29	95.80
<i>Brightness</i>	$f = 2$	46.83	11.12	89.07	99.13	39.28	99.70
<i>Contrast</i>	$f = 2$	51.70	17.87	89.87	99.58	39.29	99.99
<i>Saturation</i>	$f = 2$	56.91	23.18	95.95	99.92	39.29	99.98
<i>Hue</i>	$f = 0.1$	58.85	27.73	96.60	99.98	39.29	99.99
<i>Superimpose</i>	-	46.43	13.92	50.24	98.54	-	-
Average	-	54.12	25.67	78.62	99.69	39.28	99.64

$\lambda_{RWm} = 1$ and $\lambda_{RI_m} = 0$, respectively. More experimental details can be found in the Appendix.

4.5 Visualization Results

The results of our model *CIN* against various noises are visualized in Fig. 8. Each column indicates the result against a specific noise. And we omit the results of noise *Identity*, *Cropout*, and *Dropout* since they appear almost identical to the input image. The first two rows are the input image I_m and the output watermarked image WI_m , respectively. We can find that the image WI_m and I_m are almost indistinguishable visually, which indicates that our model has excellent imperceptibility. The third row is the noised image I_m^{noised} attacked by the specific noise. The bottom row shows the magnified difference $|I_m^{noised} - WI_m|$ between the watermarked image WI_m and the noised image I_m^{noised} , which indirectly indicates the intensity of the noise. In the Appendix, we present detailed noise parameters and experimental results against rotation, affine, and combinatorial attacks.

The detailed experimental performance corresponding to Fig. 8 is listed in Table 1. In the *Identity* case, the $PSNR$ reaches 67.66dB with BER less than 10^{-4} , which demonstrates the high imperceptibility of our framework. The results are given for the two mechanisms of specific and combined noise. For specific noise, we list the $PSNR_1$ between I_m and WI_m , the $PSNR_2$ between WI_m and I_m^{noised} , the *Pretraining* (*Pre*) accuracy tested on the *Identity* model and the accuracy (*Acc*) tested on the specific noise model. The average $PSNR$ and *Acc* with specific noise reach 54.76dB and 99.78% respectively, which indicates that our framework has great potential against

Table 2: Comparison results of combined noise on COCO dataset. The message length of all models is 30. Red represents the top accuracy value, blue takes the second place, and underlining indicates equal accuracy. Adjusting the $PSNR$ to 38.51 (dB) by the strength factor.

Models	Imp $PSNR$	Robustness (Acc%)				
		<i>Cropout</i> $p = 30\%$	<i>Dropout</i> $p = 30\%$	<i>Resize</i> $p = 50\%$	<i>Jpeg</i> $Q = 50$	<i>Ave</i> -
HiDDeN	33.5	75.96	76.89	82.72	84.09	79.915
ABDH	32.8	74.82	75.31	80.23	82.62	78.245
DA	33.7	78.58	77.13	81.72	82.82	80.06
IGA	32.8	79.33	77.51	81.44	87.35	81.40
ReDMark	-	92.5	92.00	94.1	74.6	86.36
TSDL	33.5	97.30	97.40	92.80	76.20	90.92
MBRS	33.5	<u>99.99</u>	<u>99.99</u>	-	<u>95.51</u>	<u>98.49</u>
CIN*	<u>36.2</u>	97.41	99.33	<u>99.64</u>	77.49	93.46
CIN	38.51	<u>99.99</u>	<u>99.99</u>	<u>99.99</u>	<u>99.24</u>	<u>99.80</u>

specific noise we have no test. For combined noise, we list the *Acc* of the restored watermark RWm and the $PSNR$ between I_m^{noised} and WI_m . Moreover, the mean values of $PSNR$ and *Acc* with the combined noises N_{pool} reach 39.28dB and 99.64%, respectively.

As shown in Fig. 9, we visualize the watermark patterns that the model tends to embed for different noises. The watermarked pixel with *Identity* is relatively concentrated in areas with texture information where watermarks can be easily embedded, while the region of the combined model and *Crop* are more globally embedded to resist multiple noise and random cropping. *RealJpeg* model is embedded in a way that can resist quantization loss.

4.6 Comparison against SOTA methods

In this work, we compare with several outstanding methods, such as HiDDeN [50], DA [30], ABDH [46], IGA [48], TSDL [26], ReDMark [2], and MBRS [21]. To evaluate the performance of our model compared with other methods, we conduct experiments using noise N_{pool}^{cp1} in Table 2. Our model not only has a higher $PSNR$ than other methods but also achieves the best results in terms of robustness.

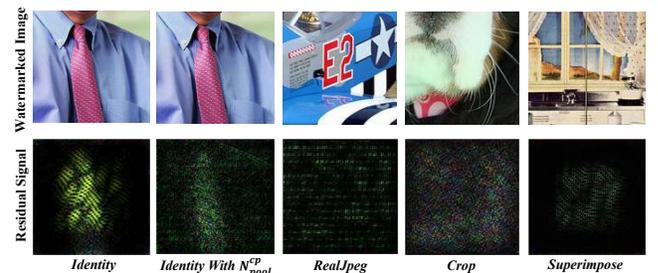


Figure 9: Visually compare watermarks embedded in images. Top: watermarked image. Bottom: the magnified difference $|I_m - WI_m|$ between input image and watermarked image.

Table 3: Compare the robustness of combined noise on the COCO and DIV2K dataset. We use the model trained with noise N_{pool} to evaluation on both datasets. The PSNR of CIN is adjusted to 37.28 dB on the COCO dataset and 40.08 dB on DIV2K. The PSNR of the other references is 33.5 dB.

Dataset	Methods	Robustness (Acc%)					
		Crop $p = 1\%$	Salt&Pepper $p = 10\%$	GauNoise $\sigma = 25$	GauBlur $k = 3$		
COCO	TSDL	75.3	90.9	74.4	99.1		
	CIN*	77.31	99.82	99.88	99.58		
	CIN	98.81	99.99	99.99	99.97		
DIV2K		Crop	Cropout	Dropout	Resize	Jpeg	
		$p = 3.5\%$	$p = 30\%$	$p = 30\%$	$p = 50\%$	$Q = 50$	
		HiDDeN	68.24	60.92	63.78	66.28	66.37
		ABDH	62.24	59.71	58.72	60.83	63.44
		DA	77.32	77.11	74.55	71.01	82.35
		IGA	77.39	60.93	76.63	72.19	82.90
		CIN*	83.70	97.00	99.29	99.72	75.31
CIN	99.99	99.84	99.99	99.99	97.52		

In Table 3, all models are trained and evaluated with watermark length $L = 30$. We find that our model can resist more kinds of noise (N_{pool} has 12 types of noise) while achieving optimal robustness and a much higher PSNR than other methods.

4.7 Ablation Study

We conduct experiments with noise pool N_{pool}^{cp2} in Table 4. At watermark length $L=30$, the PSNR is considerably higher than the reference method, and the BER is lower. At $L=64$, our approach is significantly more robust to cropping than MBRS and has a slightly higher PSNR. Our method achieves excellent results in terms of robustness and imperceptibility compared to the SOTA MBRS [21].

In Fig. 10, the comparison experiments with the model MBRS show that our framework achieves higher PSNR and SSIM at lower BER. Meanwhile, in the experiments with higher Jpeg compression strength, as shown in the left part for $Q=10$, our BER is significantly lower than MBRS. In addition, as shown in the right part, our SSIM is also noticeably higher than the reference at the strength factor of 1.4.

Through the ablation experiments in Table 5, we can find that when only the IM module (ICN*) is available, the Acc against RealJpeg is 77.49%, and the watermark intensity cannot be flexibly adjusted. After adding the DEM and FSM modules, the Acc of the

Table 4: Comparison to MBRS with combined noise N_{pool}^{cp2} on COCO dataset. The PSNR (dB) and BER(%) are given in the table. The watermark lengths are 30 and 64.

Models	L=30			L=64			
	PSNR	Crop	Jpeg	PSNR	Crop	Cropout	Jpeg
MBRS	33.5	4.15	4.48	33.5	45.86	32.86	4.14
CIN	38.51	0.09	2.6	34.22	13.40	13.27	6.77

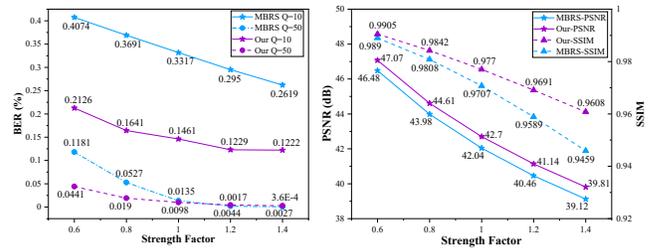


Figure 10: Compared to the methods MBRS with RealJpeg noise. The left part gives the BER with a quality factor of 10 and 50, respectively. The right part provides the PSNR and SSIM with strength factors, respectively. CIN achieve outstanding performance not only in BER but also in PSNR and SSIM.

watermark improves by 7.2%, and the PSNR improves by 2.38%. Finally, after employing NIAM and NSM modules, the model's accuracy against RealJpeg improves by 16.9%, and the PSNR and the Acc of resistance to multiple noises are also enhanced.

Table 5: Ablation experiments. The average Acc and PSNR with N_{pool} and RealJpeg are given in the table, respectively. S indicates whether the watermark strength is adjustable.

Modules				Acc (%)		PSNR (dB)	
IM	DEM&FSM	NIAM&NSM	S	N_{pool}	RealJpeg	N_{pool}	RealJpeg
✓			×	91.23	77.49	36.21	36.30
✓	✓		✓	98.43	78.90	38.59	38.50
✓	✓	✓	✓	99.64	95.80	39.28	39.29

5 CONCLUSIONS

We propose a CIN framework that learns a joint representation between watermark embedding and extraction, which effectively improve the imperceptibility of watermarking against traditional noise. To resist the non-differentiable lossy compression noise, we introduce a NIAM to improve the decoder's performance against non-additive quantization noise. In addition, we present a DEM to embed and extract watermark with high robustness. Finally, the NSM enables the appropriate decoder for compression or other noises. Extensive experiments on COCO and DIV2K datasets show that our method performs better in imperceptibility and robustness.

ACKNOWLEDGMENTS

This work was supported by National Key R&D Program of China 2021ZD0109802 and National Science Foundation of China 61971047.

REFERENCES

- [1] Eirikur Agustsson and Radu Timofte. 2017. Ntire 2017 challenge on single image super-resolution: Dataset and study. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 126–135.
- [2] Mahdi Ahmadi, Alireza Norouzi, Nader Karimi, Shadrokh Samavi, and Ali Emami. 2020. ReDMark: Framework for residual diffusion watermarking based on deep networks. *Expert Systems with Applications* (2020), 113157.

- [3] Reem A Alotaibi and Lamiaa A Elrefaie. 2019. Text-image watermarking based on integer wavelet transform (IWT) and discrete cosine transform (DCT). *Applied Computing and Informatics* 15, 2 (2019), 191–202.
- [4] Lynton Ardizzone, Jakob Kruse, Sebastian Wirkert, Daniel Rahner, Eric W Pellegrini, Ralf S Klessen, Lena Maier-Hein, Carsten Rother, and Ullrich Köthe. 2018. Analyzing inverse problems with invertible neural networks. *arXiv preprint arXiv:1808.04730* (2018).
- [5] Lynton Ardizzone, Carsten Lüth, Jakob Kruse, Carsten Rother, and Ullrich Köthe. 2019. Guided image generation with conditional invertible neural networks. *arXiv preprint arXiv:1907.02392* (2019).
- [6] Jens Behrmann, Will Grathwohl, Ricky TQ Chen, David Duvenaud, and Jörn-Henrik Jacobsen. 2019. Invertible residual networks. In *International Conference on Machine Learning*. PMLR, 573–582.
- [7] Olivia Byrnes, Wendy La, Hu Wang, Congbo Ma, Minhui Xue, and Qi Wu. 2021. Data hiding with deep learning: A survey unifying digital watermarking and steganography. *arXiv preprint arXiv:2107.09287* (2021).
- [8] Ricky TQ Chen, Jens Behrmann, David K Duvenaud, and Jörn-Henrik Jacobsen. 2019. Residual flows for invertible generative modeling. *Advances in Neural Information Processing Systems* (2019).
- [9] Ka Leong Cheng, Yueqi Xie, and Qifeng Chen. 2021. IICNet: A Generic Framework for Reversible Image Conversion. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1991–2000.
- [10] Laurent Dinh, David Krueger, and Yoshua Bengio. 2014. Nice: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516* (2014).
- [11] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. 2016. Density estimation using real nvp. *arXiv preprint arXiv:1605.08803* (2016).
- [12] Xintao Duan, Kai Jia, Baoxia Li, Daidou Guo, En Zhang, and Chuan Qin. 2019. Reversible image steganography scheme based on a U-Net structure. *IEEE Access* 7 (2019), 9314–9323.
- [13] Will Grathwohl, Ricky TQ Chen, Jesse Bettencourt, Ilya Sutskever, and David Duvenaud. 2018. Ffjord: Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint arXiv:1810.01367* (2018).
- [14] Zhenyu Guan, Junpeng Jing, Xin Deng, Mai Xu, Lai Jiang, Zhou Zhang, and Yipeng Li. 2022. DeepMIH: Deep Invertible Network for Multiple Image Hiding. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2022).
- [15] Mohamed Hamidi, Mohamed El Haziti, Hocine Cherifi, and Mohammed El Hassouni. 2018. Hybrid blind robust image watermarking technique based on DFT-DCT and Arnold transform. *Multimedia Tools and Applications* 77, 20 (2018), 27181–27214.
- [16] Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. 2019. Flow++: Improving flow-based generative models with variational dequantization and architecture design. In *International Conference on Machine Learning*. PMLR, 2722–2730.
- [17] Hai-tao Hu, Ya-dong Zhang, Chao Shao, and Quan Ju. 2014. Orthogonal moments based on exponent functions: Exponent-Fourier moments. *Pattern Recognition* 47, 8 (2014), 2596–2606.
- [18] Ming-Kuei Hu. 1962. Visual pattern recognition by moment invariants. *IRE transactions on information theory* 8, 2 (1962), 179–187.
- [19] Aapo Hyvärinen and Petteri Pajunen. 1999. Nonlinear independent component analysis: Existence and uniqueness results. *Neural networks* 12, 3 (1999), 429–439.
- [20] Jörn-Henrik Jacobsen, Arnold Smeulders, and Edouard Oyallon. 2018. i-irvnet: Deep invertible networks. *arXiv preprint arXiv:1802.07088* (2018).
- [21] Zhaoyang Jia, Han Fang, and Weiming Zhang. 2021. Mbrs: Enhancing robustness of dnn-based watermarking by mini-batch of real and simulated jpeg compression. In *Proceedings of the 29th ACM International Conference on Multimedia*. 41–49.
- [22] Haribabu Kandi, Deepak Mishra, and Subrahmanyam RK Sai Gorthi. 2017. Exploring the learning capabilities of convolutional neural networks for robust image watermarking. *Computers & Security* 65 (2017), 247–268.
- [23] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [24] Durk P Kingma and Prafulla Dhariwal. 2018. Glow: Generative flow with invertible 1x1 convolutions. *Advances in neural information processing systems* 31 (2018).
- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.
- [26] Yang Liu, Mengxi Guo, Jian Zhang, Yuesheng Zhu, and Xiaodong Xie. 2019. A novel two-stage separable deep learning framework for practical blind watermarking. In *Proceedings of the 27th ACM International Conference on Multimedia*. 1509–1517.
- [27] Yang Liu, Zhenyue Qin, Saeed Anwar, Pan Ji, Dongwoo Kim, Sabrina Caldwell, and Tom Gedeon. 2021. Invertible denoising network: A light solution for real noise removal. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13365–13374.
- [28] Shao-Ping Lu, Rong Wang, Tao Zhong, and Paul L Rosin. 2021. Large-capacity image steganography based on invertible neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10816–10825.
- [29] Andreas Lugmayr, Martin Danelljan, Luc Van Gool, and Radu Timofte. 2020. SrfFlow: Learning the super-resolution space with normalizing flow. In *European conference on computer vision*. Springer, 715–732.
- [30] Xiyang Luo, Ruohan Zhan, Huiwen Chang, Feng Yang, and Peyman Milanfar. 2020. Distortion agnostic deep watermarking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13548–13557.
- [31] Rajesh Mehta, Navin Rajpal, and Virendra P Vishwakarma. 2016. LWT-QR decomposition based robust and efficient image watermarking scheme using Lagrangian SVR. *Multimedia Tools and Applications* 75, 7 (2016), 4129–4150.
- [32] Seung-Min Mun, Seung-Hun Nam, Haneol Jang, Dongkyu Kim, and Heung-Kyu Lee. 2019. Finding robust domain from attacks: A learning framework for blind watermarking. *Neurocomputing* 337 (2019), 191–202.
- [33] Zahra Pakdaman, Saeid Saryzadi, and Hossein Nezamabadi-Pour. 2017. A prediction based reversible image watermarking in Hadamard domain. *Multimedia Tools and Applications* 76, 6 (2017), 8517–8545.
- [34] Albert Pumarola, Stefan Popov, Francesc Moreno-Noguer, and Vittorio Ferrari. 2020. C-flow: Conditional generative flow models for images and 3d point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 7949–7958.
- [35] Kartik Sharma, Ashutosh Aggarwal, Tanay Singhania, Deepak Gupta, and Ashish Khanna. 2019. Hiding data in images using cryptography and deep neural network. *arXiv preprint arXiv:1912.10413* (2019).
- [36] Yang Song, Chenlin Meng, and Stefano Ermon. 2019. Mintnet: Building invertible neural networks with masked convolutions. *Advances in Neural Information Processing Systems* 32 (2019).
- [37] Abdallah Soualmi, Adel Alti, and Lamri Laouamer. 2018. Schur and DCT decomposition based medical images watermarking. In *2018 Sixth International Conference on Enterprise Systems (ES)*. IEEE, 204–210.
- [38] Qingtang Su, Yugang Niu, Hailin Zou, Yongsheng Zhao, and Tao Yao. 2014. A blind double color image watermarking algorithm based on QR decomposition. *Multimedia tools and applications* 72, 1 (2014), 987–1009.
- [39] Teycho FA van der Ouderaa and Daniel E Worrall. 2019. Reversible gans for memory-efficient image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4720–4728.
- [40] Ron G Van Schyndel, Andrew Z Tirkel, and Charles F Osborne. 1994. A digital watermark. In *Proceedings of 1st international conference on image processing*, Vol. 2. IEEE, 86–90.
- [41] Yaolong Wang, Mingqing Xiao, Chang Liu, Shuxin Zheng, and Tie-Yan Liu. 2020. Modeling lost information in lossy image compression. *arXiv preprint arXiv:2006.11999* (2020).
- [42] Bingyang Wen and Sergul Aydore. 2019. Romark: A robust watermarking system using adversarial training. *arXiv preprint arXiv:1910.01221* (2019).
- [43] Eric Wengrowski and Kristin Dana. 2019. Light field messaging with deep photographic steganography. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1515–1524.
- [44] Mingqing Xiao, Shuxin Zheng, Chang Liu, Yaolong Wang, Di He, Guolin Ke, Jiang Bian, Zhouchen Lin, and Tie-Yan Liu. 2020. Invertible image rescaling. In *European Conference on Computer Vision*. Springer, 126–144.
- [45] Yazhou Xing, Zian Qian, and Qifeng Chen. 2021. Invertible image signal processing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 6287–6296.
- [46] Chong Yu. 2020. Attention based data hiding with generative adversarial networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 1120–1128.
- [47] Chaoning Zhang, Chenguo Lin, Philipp Benz, Kejiang Chen, Weiming Zhang, and In So Kweon. 2021. A brief survey on deep learning based data hiding, steganography and watermarking. *arXiv preprint arXiv:2103.01607* (2021).
- [48] Honglei Zhang, Hu Wang, Yuanzhouhan Cao, Chunhua Shen, and Yidong Li. 2020. Robust Data Hiding Using Inverse Gradient Attention. *arXiv preprint arXiv:2011.10850* (2020).
- [49] Kevin Alex Zhang, Alfredo Cuesta-Infante, Lei Xu, and Kalyan Veeramachaneni. 2019. SteganoGAN: high capacity image steganography with GANs. *arXiv preprint arXiv:1901.03892* (2019).
- [50] Jiren Zhu, Russell Kaplan, Justin Johnson, and Li Fei-Fei. 2018. Hidden: Hiding data with deep networks. In *Proceedings of the European conference on computer vision (ECCV)*. 657–672.
- [51] Xiaobin Zhu, Zhuangzi Li, Xiao-Yu Zhang, Changsheng Li, Yaqi Liu, and Ziyu Xue. 2019. Residual invertible spatio-temporal network for video super-resolution. In *Proceedings of the AAAI conference on artificial intelligence*. 5981–5988.