# Learning to Immunize Images for Tamper Localization and Self-Recovery

Qichao Ying$^{†}$, *Student Member, IEEE*, Hang Zhou$^{†}$, Zhenxing Qian*, *Member, IEEE*, Sheng Li, *Member, IEEE* and Xinpeng Zhang, *Member, IEEE*

**Abstract**—Digital images are vulnerable to nefarious tampering attacks such as content addition or removal that severely alter the original meaning. It is somehow like a person without protection that is open to various kinds of viruses. Image immunization (Imuge) is a technology of protecting the images by introducing trivial perturbation, so that the protected images are immune to the viruses in that the tampered contents can be auto-recovered. This paper presents Imuge+, an enhanced scheme for image immunization. By observing the invertible relationship between image immunization and the corresponding self-recovery, we employ an invertible neural network to jointly learn image immunization and recovery respectively in the forward and backward pass. We also introduce an efficient attack layer that involves both malicious tamper and benign image post-processing, where a novel distillation-based JPEG simulator is proposed for improved JPEG robustness. Our method achieves promising results in real-world tests where experiments show accurate tamper localization as well as high-fidelity content recovery. Additionally, we show superior performance on tamper localization compared to state-of-the-art schemes based on passive forensics.

**Index Terms**—Image tamper localization; Image immunization; Image recovery; Steganography; Robustness

---◆---

## 1 INTRODUCTION

DIGITAL images have largely replaced conventional photographs from all walks of life. The Online Social Network (OSN) platforms like Instagram and Twitter are designed to amplify the power of image sharing, where users create, curate, and share unique images that spark conversation and speak for themselves. However, digital images can hardly enjoy the credibility of their conventional counterparts. The rapid advancements in digital image processing tools have made it extremely easy to edit images for free, and the modified images can be shared in seconds with social networking services. Although most common image editing in life is benign and unprofitable, maliciously fabricated images can be utilized as a supplement to fake news or criminal investigation to potentially influence public opinion. For example, a critical object can be replaced with an image patch from the same image, which is known as the *copy-move attack*. Or a non-existing object can be added to the scene, which is known as the *splicing attack* . Furthermore, image inpainting techniques [7], [11] have facilitated the removal of unwanted regions without introducing noticeable artifacts. What's worse, the readers are susceptible to well-crafted fake images and they might further circulate these fake images. In that sense, digital images are like *people* without protection yet open to a variety of attacks in the wild, and the social networks will be paralyzed if crowded with *sick people*.

Image tamper localization has aroused extensive research interest to combat forged images and the security threats mentioned above. The research of fake image forensics is to distinguish tampered images, or *sick images*, from non-tampered images, or *healthy images*. The technology is of considerable significance either for academia or industry. With the emergence of deep learning, the capability of image forgery detection is strengthened to a great extent [44], [94], [118]. These schemes are largely built upon detecting noise-level manipulation, where the main component of the images is usually suppressed and the edge information as well as the noise distribution are studied for trace detection. However, the gains fail to meet the expectation in that it is hard to build a universal image tamper detection scheme, considering the enormous ways of image tampering in the wild. In addition, these methods are often less effective on compressed or low-resolution images. Besides, none of the off-the-shelf image tamper detection methods are equipped with self-recovery ability. Once the images are manipulated, it is hard for current techniques to reproduce the original contents. However, barely knowing the tampered regions cannot provide adequate information for image forensics. It is usually hard to infer the intention of the attacks without the reference to the ground-truth image, the intention of the editing, for example, how to distinguish benign image beautification from malicious attacks.

Image immunization is a novel technology for protecting digital images by making them *immune* to malicious attacks, or the *viruses*. If the immunized images are manipulated during image sharing on the OSNs, the tampered contents can be auto-recovered at the recipient's side without reference to any off-the-shelf image forgery detection or image reconstruction schemes. To enable immunization, the image owner is only required to embed some imperceptible perturbations into the image and uploads the immunized version instead of the original one. Since the embedding is trivial, the normal use of the targeted image is not affected, whereas the hidden information can resist common image

- * *Corresponding author (E-mail: zxqian@fudan.edu.cn).* $^{†}$ *The leading two authors contribute equally. Q.Ying, Z.Qian, S. Li and X. Zhang are with the School of Computer Science, Fudan University, Shanghai, China. H. Zhou is with the School of Computer Science, Simon Fraser University, British Columbia, Canada.*

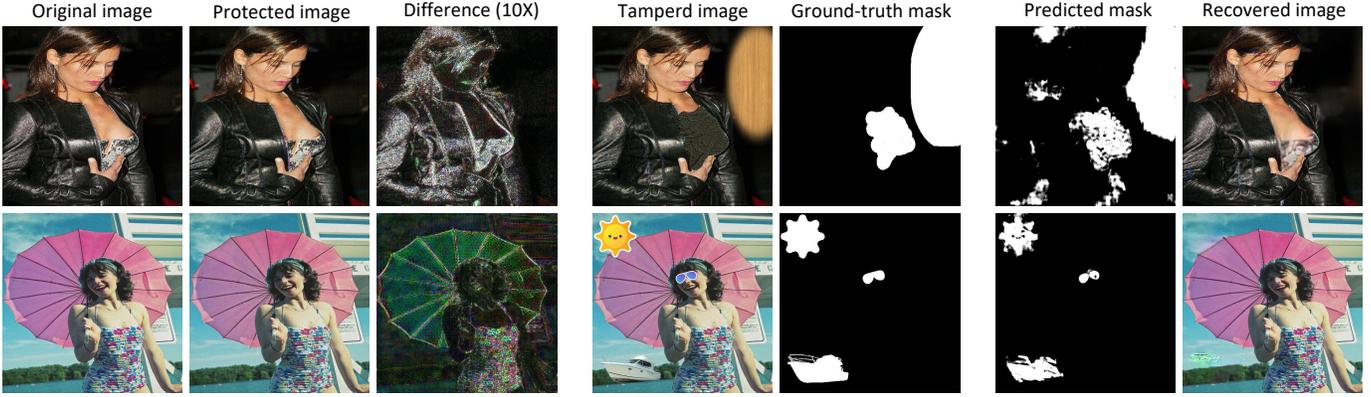| Original image | Protected image | Difference (10X) | Tamperd image | Ground-truth mask | Predicted mask | Recovered image |
|---|---|---|---|---|---|---|

Fig. 1. **Two examples of image tamper localization and content recovery by Imuge+.** Image immunization transforms the unprotected digital images into immunized versions where the residual is close to imperceptible. Malicious attacker deliberately tampers some of the contents and propagates the forged images in order to mislead the audience. Imuge+ localizes the tampered region and successfully recovers the images.

processing attacks and help the recipient localize the tamper and recover the original image faithfully. Image immunization is originally proposed and implemented by a deep-network-based framework named Imuge in our conference paper [74] where we individually employ three networks for image self-embedding, tamper localization and image self-recovery. An attack layer is also proposed to simulate both tampering attacks and image benign processing attacks for robustness training. Nevertheless, there are still many issues to be solved in Imuge. For example, when the critical components of an immunized image are removed, only a blurry and approximate version of them can be reconstructed by Imuge. Besides, the localization accuracy is not satisfactory due to the poor generalization of robustness. Thirdly, Imuge cannot resist image inpainting attacks. It motivates us to develop an improved scheme for image immunization that can be applied on a wilder range of images.

In this paper, we present an enhanced scheme denoted as Imuge+ for image immunization. We introduce trivial perturbation into the original image as immunization and blindly reconstruct the original contents in addition to tamper localization. The network is trained by multi-task learning and contains three modules. The first one is to keep the immunized image consistent with the original image, where the introduced perturbation is imperceptible. The second one is tamper localization by classifying whether a pixel is altered or not. And the third is image self-recovery to encourage the recovered image to resemble the original image. Observing the invertible nature between image immunization and self-recovery, we employ a normalizing-flow-based generator that jointly learns image immunization and self-recovery within a single network. The forward pass transforms an original image as well as its edge map into the corresponding immunized version. On receiving the attacked image, we use a localizer to determine the tampered areas by predicting the tamper mask, and in the backward pass of the generator, the hidden perturbation is transformed into information and we recover the original image as well as the edge map. Three most typical malicious attacks are simulated for effective network training, e.g., copy-move and splicing, and benign attacks (rescaling, blurring). Furthermore, considering that

most attackers deliberately hide their behavior by image filtering or compression, we also concatenate tamper simulation with image post-processing simulation, where a novel knowledge-distillation-based JPEG generator (KD-JPEG) is developed for enhanced JPEG robustness.

We showcase two examples of our scheme in Fig. 1 where some crucial contents containing semantic information in the two protected images are removed by the malicious attacker. The recipient gets the JPEG-compressed tampered images and successfully retrieves the missing contents. Therefore, the recipient can identify the fake images and accordingly block their spreading. We test our scheme by introducing human-participated hybrid digital attacks, including copy-moving, splicing, inpainting and post-processing. The results demonstrate that our scheme can recover the tampered contents with high quality and fidelity. Compared with Imuge [74], Imuge+ can combat more kinds of attacks such as copy-move and inpainting, where more details can be well preserved. Additionally, we show the effectiveness of tamper detection of Imuge+ by comparison with some state-of-the-art passive-based image forensics schemes.

This paper is an extended version of our conference paper [74]. We make the following new contributions:

1) We propose a tailored network named Imuge+ for image immunization, which reconstructs the original contents if the immunized images are manipulated. We develop the scheme upon viewing the image immunization and recovery processes as a pair of inverse problem.
2) We tempt to address the blurry result issue and the poor performance issue over JPEG compression by proposing a novel knowledge-distillation-based JPEG simulator as well as tampering-based data augmentation.
3) We have conducted comprehensive experiments to prove that our network design and training mechanisms remarkably improves the overall performance of image immunization, both in the quality of the recovered image and the accuracy of image tampering localization.

## 2 RELATED WORKS

In this section, we review the related works of Imuge+, namely, image protection using steganography, passive image tampering localization and image immunization.

### 2.1 Image Protection using Steganography

Image steganography aims at hiding secret information into the host images so that the recipient can extract the hidden message for covert communication. In the past decades, many steganography-based schemes for image protection have been proposed, where the extracted secret information is utilized for image tampering localization and fragile self-recovery. For example, He et al. [85] and Zhang et al. [87] respectively embed the Discrete Cosine Transform (DCT) coefficients and the Most Significant Bits (MSB), which is a compressed version of the image, into the Least Significant Bits (LSB). Zhang et al. [88] proposes to embed into the image blocks check-bits and reference-bits, where the former identifies the tampered blocks and the latter can exactly reconstruct the original image. After that, Zhang et al. [90] proposes a reference sharing mechanism, in which the watermark to be embedded is a reference derived from the principal content in different regions and shared by these regions for content recovery. Later, Zhang et al. [91] proposes a watermarking scheme with flexible recovery quality. If the amount of extracted data in the areas without tampering is not enough, the method employs a compressive sensing technique to retrieve the coefficients by exploiting the sparseness in the DCT domain. Besides, Korus et al. [89] theoretically analyzes the reconstruction performance with the use of communication theory.

Some works use steganography to prevent the images from being manipulated by generative models. Khachaturov et al. [72] proposes an adversarial method to attack inpainting systems by forcing them to work abnormally on the targeted images. Similarly, Yin et al [73] proposes a defensive method based on data hiding to defeat Super-Resolution (SR) models. The hidden information mainly resides in higher-band details of the targeted images, which are often analyzed or augmented by many schemes that employ deep networks for image restoration.

Though promising in the presented results, these methods are all fragile and typical image attacks on the modified images can significantly weaken the performance. Besides, the steganography-based schemes against generative tampering are non-blind, where attackers can still modify the protected images using traditional methods such as splicing and copy-move attack.

### 2.2 Passive Image Tampering Localization

Passive image tampering localization schemes aim at finding traces to unveil the behavior of image forgery. Many existing image forensics schemes are specified on detecting typical kinds of attacks, which can be mainly categorized into three groups, i.e., splicing detection [22], [23], copy-move detection [21], [118] and inpainting detection [14], [20]. As manipulating a specific region in a given image inevitably leaves traces between the tampered region and its surrounding, there are also many schemes for universal tampering detection [35], [44], [94] that exploit such noise artifact. Li et al. [14] proposes to implement an FCN's first convolutional layer with trainable high-pass filters and apply their HP-FCN for inpainting detection. Kown et al. [23] propose to model quantized DCT coefficient distribution to trace compression artifacts in splicing attacks. DOA-GAN [118] proposes two attention modules for copy-move detection, where the first is from an affinity matrix based on the extracted feature vectors at every pixel, and the second is to further capture more precise patch inter-dependency.

For universal tampering detection, Mantra-Net [94] uses fully convolutional networks with BayarConv [39] and SRMConv [8] for feature extraction and further uses Z-Pooling layers as well as long short-term memory (LSTM) cells for pixel-wise anomaly detection. In MVSS-Net [44], a system with multi-view feature learning and multi-scale supervision is developed to jointly exploit the noise view and the boundary artifact to learn manipulation detection features. Hu et al. [35] proposes SPAN that models the relationship between image patches at multiple scales by constructing a pyramid of local self-attention blocks.

Despite the existence of these well-designed works, real-world image tampering localization is still an open issue. Besides, attackers can use a chain of image post-processing methods to hide their behaviors. Previous works are reported to have poor generalization against image post-processing or on JPEG-format images [94], [118], where the learned clues can be easily erased. In this paper, we realize robust image tampering localization using image self-embedding.

### 2.3 Image Restoration and Immunization

Image restoration schemes reconstruct an image with higher visual quality. For example, image inpainting schemes [5], [7] restore the contents within missing areas by referring to the ambient regions or the learned deep image prior [115]. Yu et al. [5] proposes a CNN to synthesize novel image structures within the missing areas by explicitly utilizing surrounding image features as references. EdgeConnect [7] improves the visual quality of the generated images by reconstructing the edge information of the missing area ahead of image restoration. LaMa [3] employs fast Fourier convolutions (FFC) to widen the receptive field to grasp more global statistical features for completing large missing areas. However, the results of the inpainting schemes can be natural yet faulty compared to the ground truth in that providing visually pleasing results is the priority other than the reversibility. As a result, image inpainting is more often used to moderately manipulate an image for better layout rather than faithfully recover the image.

Compared to image restoration, image immunization [74] is a recently proposed novel technology that protects images from being illegally tampered. The manipulated contents can be identified and auto-recovered at the recipient's side without reference to sophisticated image forgery detection or image reconstruction schemes. In [74], a U-Net-based encoder [59] is employed to conceal trivial information into the original image, where the hidden data serves as a *vaccine* that helps conduct tampering localization and image self-recovery. A differentiable attack layer is
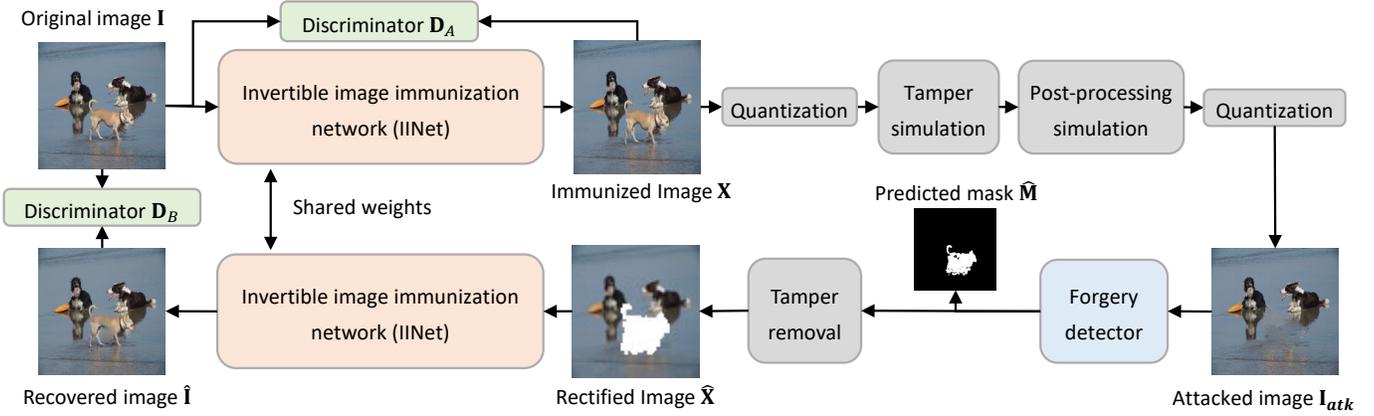
Fig. 2. **Sketch of the pipeline of Imuge+.** IINet is employed to embed slight perturbations into the original images for immunization. Hand-crafted digital attacks in the social medias are decomposed into tampering and post-processing and correspondingly simulated by an attacking layer. Afterwards, a forgery detector helps predicting the tampering mask and removing the tamper from the received image. Finally, the inversed procedures of IINet recovers the original image. We additionally introduce two discriminators to improve the quality of the generated images.

proposed to simulate both tampering attacks and image benign processing attacks for robustness training. We train a forgery detector as well as an image decoder based on U-Net architecture. After the recipient gets the tampered protected image, he can get the original non-tampered version of the received attack image by tampering localization and self-recovery. Nevertheless, the scheme still has several drawbacks. First, Imuge can only approximately recover the original contents within the tampered areas. There is still a big gap toward high-quality and accurate image recovery. Second, the accuracy of tampering prediction is not satisfactory in many cases due to poor generalization. It motivates us to present an enhanced scheme for image immunization to comprehensively address the above issues.

## 3 METHOD

In this section, we elaborate on our method for image immunization and self-recovery. We first present the problem statement of image immunization and self-recovery. Then, we show the details of the network architecture and the learning objectives. Finally, we introduce the training mechanisms.

### 3.1 Approach Overview

**Problem Statement.** Typical malicious image manipulation attack can be generalized by a function $Mani(\cdot)$ in Eq. (1).

$$\mathbf{X}_{atk} = Mani(\mathbf{X}, \mathbf{M}) = IP\left(\mathbf{X} \cdot (1 - \mathbf{M}) + \mathbf{R} \cdot \mathbf{M}\right), \quad (1)$$

where $\cdot$ denotes pixel-wise multiplication. $\mathbf{X}$ and $\mathbf{X}_{atk}$ are respectively the targeted image and its attacked version. $\mathbf{M}$ is the tampering mask, which is generally the region of interest of $\mathbf{X}$. The contents within the mask are replaced with irrelevant contents denoted by $\mathbf{R}$. $IP(\cdot)$ summarizes the image post-processing behaviors such as JPEG compression, rescaling, etc. Image tampering localization is to determine the tampering mask $\mathbf{M}$ from the received attacked image $\mathbf{X}_{atk}$, and image self-recovery is to take an advanced step to reconstruct the destroyed contents, i.e., $\mathbf{X} \cdot \mathbf{M}$. We employ two functions $f_M(\cdot)$, $f_{rec}(\cdot)$ respectively for these two tasks.

$$\hat{\mathbf{X}} = f_{rec}\left(\mathbf{X}_{atk} \cdot \left(1 - \hat{M}\right)\right), \quad (2)$$

$$\hat{M} = f_M(\mathbf{X}_{atk}), \quad (3)$$

where $\hat{\mathbf{X}}$ and $\hat{\mathbf{M}}$ are respectively the reconstructed image and the predicted tampering mask. However, since the non-tampered areas of $\mathbf{X}_{atk}$ do not contain any information of $\mathbf{X} \cdot \mathbf{M}$. As a result, if $\mathbf{X}$ in Eq. (1) is directly the original image $\mathbf{I}$, it leads to a sub-optimal solution that the destroyed contents will be hallucinated using image prior [115] rather than faithfully recovered.

To address the issue, we employ a third function $f_{prt}(\cdot)$ to embed deep representations of an original image into itself, i.e., $\mathbf{X} = f_{prt}(\mathbf{I})$ in Eq. (1), and formulate $\mathbf{X} = f_{prt}(\mathbf{I})$ and $\hat{\mathbf{I}} = f_{rec}(\mathbf{X}_{atk} \cdot (1 - \hat{\mathbf{M}}))$ as a pair of invertible functions, where we expect that $\hat{\mathbf{I}} = \mathbf{I}$ and $\hat{\mathbf{M}} = \mathbf{M}$. In other words, by introducing image self-embedding, we wish to recover $\mathbf{I}$ when the tampered area is detected and removed, even if the self-embedded version $\mathbf{X}$ is attacked with randomized attacking method $IP(\cdot)$ and tampering mask $\mathbf{M}$.

**Network Modeling.** Accordingly, we design four phases for Imuge+, namely, image immunization, image redistribution, forgery detection and image recovery. We use a single INN network called Invertible Immunization Network (IINet) to model $f_{prt}(\cdot)$ and $f_{rec}(\cdot)$ simultaneously. Image tampering and post-processing are implemented by an attacking simulation layer which consists of several differentiable methods. A forgery detector is trained for $f_M(\cdot)$. We also introduce two discriminators denoted as $\mathcal{D}_A$ and $\mathcal{D}_B$ to respectively distinguish $\mathbf{X}$ and $\hat{\mathbf{I}}$ from $\mathbf{I}$. Fig. 2 shows a sketch of the pipeline of our scheme.

Specifically, we transform the original image into the immunized image using the forward pass of IINet. The protected version is uploaded onto the social cloud for daily applications instead of the unprotected original image. To simulate the image redistribution stage, we first perform tampering that evenly varies from splicing, copy-move or inpainting attack. Then, common image post-processing attacks are performed before data storage. On the recipient's side, the forgery detector produces the predicted tampering mask to see which parts of the image are manipulated, and we correspondingly remove the tampered contents. Finally, by inversely running the IINet, we reconstruct the recovered

images. If the generated images are of high quality, they can evade the classification of the two discriminators.

## 3.2 Framework Specification

**Invertible Immunization Network (IINet)** IINet embeds deep representation of the original image into itself, and conducts self-recovery provided with the rectified image. It is empirical that image immunization should not affect normal use of the image. We accepts the invertible U-shaped network proposed by Zhao et al. [106], where the involved wavelet-based blocks can provide empirical biases to decompose the generated features into lower and higher sub-bands. The design helps IINet restrict modifications towards lower sub-bands and conceal the information required by immunization in the higher sub-bands. We additionally generate the edge map $\mathbf{E}$ of the original image using the canny edge detection algorithm as the additional input of IINet. The introduction of the edge map is to aid high-fidelity image recovery, which is inspired by EdgeConnect [7]. The edge map is a one-channel matrix, and therefore the channel number of the input is four. Accordingly, there will be an additional output channel, denoted as $\mathbf{Y}$. We nullify $\mathbf{Y}$ by forcing it to be close to a zero matrix $\mathbf{O}$. In the inversed process, we also feed an extra $\mathbf{O}$ and let IINet output the predicted edge map $\hat{\mathbf{E}}$ along with $\hat{\mathbf{I}}$.

**Attack simulation.** We introduce the malicious attacks before introducing benign post-processing attacks. To begin with, several random areas within $\mathbf{X}$ are selected and manipulated according to the tampering mask $\mathbf{M}$. We denote the tampered image before image post-processing attack as $\mathbf{X}_{tmp}$, where $\mathbf{X}_{tmp} = \mathbf{X} \cdot (1 - \mathbf{M}) + \mathbf{R} \cdot \mathbf{M}$. The generation of $\mathbf{M}$ is borrowed from DeepFillV2 [5]. We restrict the rate of the tampered region as $r_t \in [0.0, 0.5]$. To simulate the copy-move attack, we let $\mathbf{R}$ in Eq. (1) as a spatially-shifted version of $\mathbf{I}_M$, which also contains the embedded information. To simulate the splicing attack, we simply use another clean image for replacement. To simulate the inpainting attacking, we use the open-source model from DeepFillv2 [5] to generate the inpainted result as $\mathbf{R}$. The three kinds of tamper are iteratively and evenly performed for balanced training. Note that minimizing the false alarm rate is also important for computer vision tasks. Accordingly, we force Imuge+ to predict a zero matrix for tampering localization when there is no tampering attack, i.e., $r_t = 0$ and $\mathbf{X}_{tmp} = \mathbf{X}$.

After image tampering, we implement the most-commonly-seen image post-processing attacks to simulate that the attacker wants to conceal the tampering behavior. $\mathbf{X}_{atk} = IP(\mathbf{X}_{tmp})$. During training, we simulate the most common types of image post-processings, including: Addition of White Gaussian noise (AWGN), Gaussian blurring, image rescaling, lossy image compression (PNG and JPEG), median filtering and cropping. Note that if $\mathbf{X}$ is cropped, the recipient only conducts forgery detection and image recovery within the cropped area.

Finally, both $\mathbf{X}$ and $\mathbf{X}_{atk}$ are required to be converted into 8-bit RGB format. We perform differentiable image quantization using Straight-Through Estimator [36] for gradient approximation.

**Knowledge-distillation-based JPEG simulator.** Fig. 3 illustrates the network design of KD-JPEG, which consists of



(a) Knowledge-Distillation-based JPEG simulator (KD-JPEG)



(b) JPEG quality factor (QF) predictor



(c) Student / Teacher network
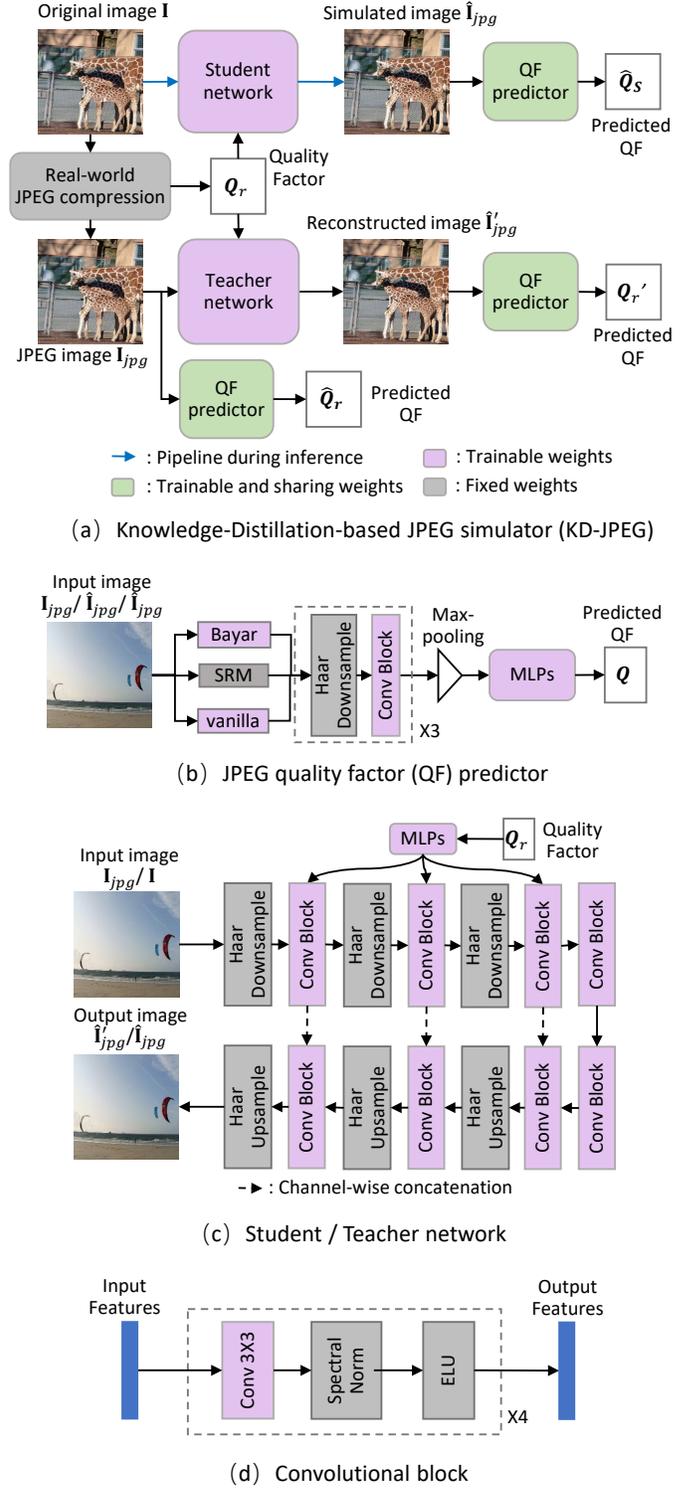


(d) Convolutional block

Fig. 3. **Illustration of KD-JPEG.** The student/teacher network are respectively provided with the plain-text image and its JPEG version using real-world compressor. The two networks output the generated JPEG images. The predictor is pre-trained by classifying the real-world JPEG images and guides the generation of the simulated JPEG images. To save computational complexity, these networks are based on basic architectures. During inference, we only activate the student network.

three parts, namely, a real-world JPEG compressor, a QF predictor, and a pair-wise student and teacher network. Though there are already many schemes that include a carefully-

designed JPEG simulator, e.g., Diff-JPEG [34], MBRS [76] and HiDDeN [56], the real-world JPEG robustness of these schemes is still reported to be limited. It mainly attributes to that neural networks can over-fit fixed compression mode where these works use fixed codes and limited number of quantization tables. In contrast, the quantization table in real-world JPEG images can be customized and more flexibly controlled by the Quality Factor (QF) as well as the image content. We present a novel knowledge-distillation-based JPEG simulator using a novel JPEG generative network to better approximate real-world JPEG compression.

First, we use the FFJPEG library to store the plain-text images $\mathbf{I}$ in JPEG format as the ground-truth JPEG-compressed images $\mathbf{I}_{jpg}$, with ground-truth QF $Q_r$ arbitrarily assigned. We then train the QF predictor to classify $\mathbf{I}_{jpg}$ by their QF. Here we empirically set the labels of the QF classification task as $C_{QF} = \{10, 30, 50, 70, 90, 100\}$, where "100" represents "not compressed". We do not introduce more labels in between such as $\{20, 40, ..., 80, ...\}$ in that when an image is compressed using two close QF ($\delta_{QF} < 10$), the discrepancy is not significant enough for classification. We allow the network to give an imprecise prediction when $Q_r \notin C_{QF}$. Afterwards, we train the pair-wise teacher and student network to produce pseudo-JPEG images. The student network is provided with the plain-text image $\mathbf{I}$ while the teacher network is given with the corresponding real-world compressed image $\mathbf{I}_{jpg}$. They are to produce $\hat{\mathbf{I}}_{jpg}$ and $\hat{\mathbf{I}}'_{jpg}$, which should resemble $\mathbf{I}_{jpg}$. Since reconstruction is much easier than simulation, knowledge distillation is done by employing a feature consisting loss to minimize the distance between the hierarchical feature maps generated by the two networks. We enhance the performance of JPEG simulation by two constraints. First, we encourage the student network to minimize the difference in the intermediate features generated by the last three *Conv* blocks of the two networks. These features are denoted as $\phi_i, i \in [0, 2]$, and they largely affect the final output. The constrain is commonly applied in knowledge-distillation tasks. Second, we encourage both the reconstructed and simulated JPEG image, i.e., $\hat{\mathbf{I}}'_{jpg}$ and $\hat{\mathbf{I}}_{jpg}$, to be classified as $Q_r$ by the QF predictor. The characteristics of JPEG images can be better learned by KD-JPEG if the classification results are consistent with that of the real-world JPEG image.

The student/teacher network shares the same U-shaped Fully Convolutional Network (FCN) architecture [59] with eight *Conv* blocks. To simulate JPEG images with a wide variety of QFs, we respectively employ two individual five-layered Multi-Layer Perceptrons (MLP) to learn the mapping functions that output modulation parameter pairs $a, b$, where the deepest three layers control the standard deviation and mean of the features. The outputs of the leading three *Conv* blocks of are controlled by the MLP as follows.

$$\phi_i = a_i \cdot Conv\_block_i(\phi_{i-1}) + b_i, \tag{4}$$

where $\phi_{i-1}, \phi_i, a_i, b_i$ represent the input and output features, the trainable mean and standard deviation at *Conv* block $i$, with $i \in [0, 2]$. To implement the QF predictor, we consider that the trace for JPEG image classification lies mainly in the higher frequencies, we employ in parallel a vanilla *Conv* layer, an SRM *Conv* layer [8] and a Bayar

*Conv* [39] layer. The latter two are reported to be efficient in depressing the main components of the input image. Thereafter, we use a typical down-sampling FCN followed by three-layered MLP.

After training KD-JPEG, we only activate the student network and feed it with $\mathbf{X}_{tmp}$ for JPEG simulation. Consider that $\mathbf{X}_{tmp}$ contains more higher-band details compared to $\mathbf{I}$ because of information hiding, we additionally introduce some Additive White Gaussian Noise (AWGN) on $\mathbf{I}$ and encourage the removal of higher-frequent details including AWGN during training KD-JPEG.

**Forgery detector and discriminator.** The forgery detector detects the tampered areas within an attacked image. When the simulation of tampering attack is skipped, we force the detector to predict a zero matrix $\mathbf{O}$. For the two discriminators, i.e., $\mathcal{D}_A$ and $\mathcal{D}_B$, the goals are to distinguish the generated images from the original images. High-quality immunization and recovery are expected if the generated images can cause misclassification of the discriminators.

We employ the forgery detector and $\mathcal{D}_B$ respectively using a U-shaped network, which is the same in architecture as the student/teacher network in KD-JPEG (see Fig. 3) except that the MLP is not present. That is to say, $\mathcal{D}_B$ conducts a pixel-wise discrimination [144] on $\hat{\mathbf{I}}$ to predict which parts of the image are not recovered naturally enough. We use a simple Patch-GAN discriminator [61] to implement $\mathcal{D}_A$. The reason we do not employ another U-shaped discriminator is that if $\mathcal{D}_A$ is too strict on $\mathbf{X}$, there will not be enough space for information embedding, resulting in unstable training.

**Implementation details.** In IINet, we apply three Haar down-sampling layers, each appended with four Double-Side Affine Coupling (DSAC) layers proposed in [106], and the last conditional splitting layer is removed. The functions inside each DSAC layer can be represented by arbitrary neural networks by definition. We implement them using a five-layer residual *Conv* block. All down-sampling and up-sampling operations are replaced with Haar down-sampling and up-sampling transformations. Haar down-sampling transformation decomposes each channel of the input into four orthogonal channels, with half the width and height. Haar up-sampling transformation is the exact inverse. Each *Conv* block contains four *Conv* layers appended with a Spectral Normalization (SN) [113] layer and an Exponential Linear Unit (ELU) [4] layer, and we keep the dimension of the inputs and outputs as the same. Exceptions are that the leading *Conv* block of the U-shaped networks transforms the twelve-dimensional inputs into 32-dimensional features, and the last block transforms the 32-dimensional input features back into twelve-dimensional outputs.

### 3.3 Objective Loss Function

The objective functions include the protection loss $\mathcal{L}_{prt}$, the tampering localization loss $\mathcal{L}_{loc}$, the recovery loss $\mathcal{L}_{rec}$ and the nullification loss $\mathcal{L}_{null}$. There is also the simulation loss $\mathcal{L}_{jpg}$ exclusively for KD-JPEG. In the below equations, $\alpha, \beta, \gamma, \theta, \omega$ are the hyper-parameters.

The protection loss $\mathcal{L}_{prt}$ and the recovery loss $\mathcal{L}_{rec}$ respectively encourage $\mathbf{X}$ and $\hat{\mathbf{I}}$ to resemble the original image

**I**. We use the popular $\ell_1$ loss term to measure the distance between images.

$$\mathcal{L}_{prt} = \|\mathbf{I} - \mathbf{X}\|_1, \tag{5}$$

$$\mathcal{L}_{rec} = \|\mathbf{I} - \hat{\mathbf{I}}\|_1 + \|\mathbf{E} - \hat{\mathbf{E}}\|_1. \tag{6}$$

For adversarial training, the discriminators $\mathbf{D}_A$ and $\mathbf{D}_B$ need to respectively distinguish $\mathbf{X}$ and $\hat{\mathbf{I}}$ from $\mathbf{I}$. The adversarial loss $\mathcal{L}_{adv}$ for the main pipeline is as follows.

$$\mathcal{L}_{adv} = \log\left(1 - \mathbf{D}_A(\hat{\mathbf{X}})\right) + \log\left(1 - \mathbf{D}_B(\hat{\mathbf{I}})\right). \tag{7}$$

The loss for the two discriminators are respectively

$$\mathcal{L}_{\mathbf{D}_A} = -\frac{1}{2}(\log \mathbf{D}_A(\mathbf{I}) + \log\left(1 - \mathbf{D}_A(\hat{\mathbf{X}})\right)), \tag{8}$$

$$\mathcal{L}_{\mathbf{D}_B} = -\frac{1}{2}(\log \mathbf{D}_A(\mathbf{I}) + \log\left(1 - \mathbf{D}_A(\hat{\mathbf{I}})\right)). \tag{9}$$

The localization loss $\mathcal{L}_{loc}$ is to improve the accuracy of tampering localization. We minimize the Binary Cross Entropy (BCE) loss between the predicted mask $\hat{\mathbf{M}}$ and the ground-truth mask $\mathbf{M}$.

$$\mathcal{L}_{loc} = -(\mathbf{M}\log\hat{\mathbf{M}} + (1 - \mathbf{M})\log\left(1 - \hat{\mathbf{M}}\right)). \tag{10}$$

The nullification loss $\mathcal{L}_{null}$ is to nullify the additional output $\mathbf{Y}$ of IINet.

$$\mathcal{L}_{null} = \|\mathbf{Y} - \mathbf{O}\|_1. \tag{11}$$

The total loss for the main pipeline of Imuge+ is as follows.

$$\mathcal{L} = \mathcal{L}_{rec} + \alpha \cdot \mathcal{L}_{prt} + \beta \cdot \mathcal{L}_{loc} + \gamma \cdot \mathcal{L}_{null} + \omega \cdot \mathcal{L}_{adv}. \tag{12}$$

KD-JPEG is trained ahead of the whole pipeline. Given a triple $(\mathbf{I}, \mathbf{I}_{jpg}, Q_r)$, we use the Cross-Entropy (CE) loss to train the QF predictor by classifying $\mathbf{I}_{jpg}$.

$$\mathcal{L}_{QF} = CE(\mathbf{Q}_o, \mathbf{Q}_r) = -\sum_{c=1}^{6} y_{o,c}\log(p_{o,c}), \tag{13}$$

where $y$ is the binary indicator if class label $c$ is the correct classification for observation $o$. $p$ is the predicted probability observation $o$ is of class $c$. $\mathbf{Q}_o$ takes the *argmax* of $o$ that maximizes $p$. For the teacher network, we employ the CE loss and the $\ell_1$ loss for reconstructing the real-world JPEG image.

$$\mathcal{L}_{tea} = \|\hat{\mathbf{I}}_{jpg'} - \mathbf{I}_{jpg}\|_1 + \epsilon \cdot CE(\hat{\mathbf{Q}}_r', \mathbf{Q}_r). \tag{14}$$

For the student network, apart from the CE loss and the $\ell_1$ loss, we additionally employ a distillation loss. The total loss for the student network is

$$\mathcal{L}_{stu} = \|\hat{\mathbf{I}}_{jpg} - \mathbf{I}_{jpg}\|_1 + \epsilon \cdot CE(\hat{\mathbf{Q}}_s, \mathbf{Q}_r) + \sum_{i \in [0,2]} \|\phi_i^{stu} - \phi_i^{tea}\|_1. \tag{15}$$



Before DA   After DA    Before DA   After DA    Before DA   After DA

Fig. 4. **Example of data augmentation by image pre-tampering.** We prevent IINet by relying on image prior by deliberately tampering some of the original images using splicing. These add-ons will be tampered again and IINet is required to recover them.

### 3.4 Training Mechanisms

Directly training the network by minimizing $\mathcal{L}$ can hardly achieve satisfying results. The reason is mainly three-folded. First, we observe that simply varying the generated mask is not enough for effective image immunization and self-recovery, in which the network will tend to hallucinate the missing content. Secondly, we find that maintaining balanced performances under different attacks is difficult. Thirdly, a poorly-trained forgery detector will mislead the image recovery process. To address these issues, we propose the following training mechanisms.

**Tampering-based data augmentation**. In many cases, the randomly generated masks are not good enough to cover textured areas or the Regions of Interests (RoIs) within the images. Naturally, Deep Image Prior (DIP) [115] can be learned by networks to recover an approximate version of the original contents with semantic correctness. However, Imuge+ needs to faithfully reproduce the original image without hallucinating the results. Therefore, how to guide the network to correctly recover the image without using DIP is a big issue. We propose a new Data Augmentation (DA) paradigm by modifying some of the original images in the training sets using simulated splicing. Fig. 7 shows two examples of our DA. After image immunization, we exactly tamper the add-ons during tampering simulation, i.e., the mask of the first-round and second-round tampering is the same. In other words, the introduced contents are automatically set as the new RoIs of the images, and since the rest of the image shows no relation with the add-ons, IINet is forced not to use DIP for image recovery but to utilize image immunization to hide information for the recovery. We control the rate of two-round data augmentation by $r_{aug}$, and empirically find that $r_{aug} = 15\%$ provides the best performance. The reason is that DIP can facilitate the efficient encoding of the image representation.

**Asymmetric batch size.** In the majority of previous deep-network-based watermarking schemes [20], [76], the attack layer always arbitrarily and evenly performs one kind of attack on the targeted images. However, we argue that the iterative training strategy might be sub-optimal in that solutions of countering different types of attacks vary. As a result, the network upgrades noisily and unevenly among batches and therefore can be more in favor of providing solutions for trivial attacks. We propose to enlarge the batch size after attack simulation to balance the results among different attacks in each batch. Suppose that the original batch size is $n$, we perform each of the six attacks on the $n$ images, where we get $6 \cdot n$ attacked images. Then, we concatenate these images where the batch size for the backward pass becomes $6 \cdot n$. The technique is explicitly de-

| Original image | Attacked image | Tamper mask | Original image | Attacked image | Tamper mask | Original image | Attacked image | Tamper mask |

(a) Inpainting                                      (b) Copy-move                                      (c) Splicing
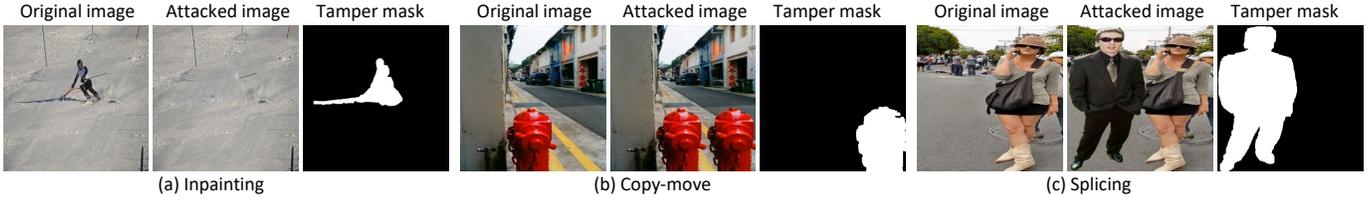
Fig. 5. **Examples of the hand-crafted test set.** The images are first immunized and then manipulated by the volunteers who use Adobe Photoshop and Microsoft Drawing. Malicious attacks include copy-move, inpainting and splicing. Benign attacks include typical image processings such as JPEG compression, rescaling, etc.

TABLE 1
**Composition of the real-world test dataset divided by the tampering rate.** In most cases, people do not modify the images for too much, and the tampering rate is generally less than 0.3. The settings are consistent with that of many off-the-shelf tampering detection datasets.

| Attack | MS-COCO | | | | ILSVRC | | | | CelebA | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | (0,0.1) | (0.1,0.2) | (0.2,0.3) | >0.3 | (0,0.1) | (0.1,0.2) | (0.2,0.3) | >0.3 | (0,0.1) | (0.1,0.2) | (0.2,0.3) | >0.3 |
| Copy-move | 203 | 145 | 104 | 64 | 147 | 85 | 54 | 22 | - | - | - | - |
| Splicing | 168 | 123 | 84 | 66 | 186 | 102 | 75 | 50 | 115 | 73 | 52 | 20 |
| Inpainting | 270 | 147 | 87 | 63 | 152 | 74 | 67 | 33 | 102 | 63 | 56 | 39 |

signed to avoid disparate statistics among results produced by different image post-processing operations.

**Iterative training..** We inherit and improve the task decoupling mechanism from [74], since a wrongly predicted tampering mask in the early training stage will unbalance the invertible function. We divide the training process into two stages. In the first stage, we individually train the forgery detector and the rest of the networks. After the generation of $\mathbf{X}_{atk}$, the forgery detector generates the predicted tampering mask $\hat{\mathbf{M}}$ and only updates itself by minimizing $\mathcal{L}_{cls}$. On the other hand, we provide IINet with the ground-truth rectified image $\hat{\mathbf{X}}_{GT}$, where we assume a perfect tampering prediction. We then update IINet by minimizing the overall loss $\mathcal{L}_{\mathcal{G}}$ with $\alpha = 0$. When $\mathcal{L}_{cls}$ converges to a low level, we move on to the second stage by canceling the task decoupling, and the whole pipeline is thus trained together, where $\hat{\mathbf{X}}$ is influenced by the performance of the forgery localization. It helps IINet to adjust with imperfect localization results.

## 4 EXPERIMENTS

In this section, we conduct experiments to evaluate Imuge+. First, we clarify the experimental setup. Then we provide comprehensive experiments and analysis on tampering localization and image self-recovery using Imuge+. Next, we show the performance comparison with the state-of-the-art methods and the ablation studies. Finally, we showcase three real-world applications where Imuge+ can be applied.

### 4.1 Experimental Setup

**Settings.** We empirically set the hyper-parameters as $\alpha = 3, \beta = 1e\text{-}3, \gamma = 10, \omega = 0.01$ and $\epsilon = 0.1$. Through extensive experiments, we find that the selection of $\alpha$ and *Th* play important role in the ultimate network performance (see Section 4.5), while the rest of the hyper-parameters have less impact. The batch size is set as four, and we use Adam optimizer [78] with the default parameters. The

learning rate is $1 \times 10^{-4}$ with the cosine annealing decay. We binarize the prediction mask by setting the threshold *Th* as 0.2. After binarization, we use image eroding operation with kernel size $k_{ero} = 8$ for noise removal within $\hat{\mathbf{M}}$ and image dilation operation with kernel size $k_{dil} = 16$ to fully cover the tampered areas. We train Imuge+ with four distributed NVIDIA RTX 3090 GPUs. The training finishes in a week.

**Data preparation.** Imuge+ is developed for immunizing natural images from randomized distributions, being it sceneries or facial images. Therefore, during training, the original images $\mathbf{I}$ are prepared by arbitrarily selecting around 10000 images from multiple popular datasets, namely, MS-COCO [64], CelebA [125], Places [126], UCID [119] and ILSVRC [62]. Since convolutions are generally not scale-agnostic, we train different models for some benchmark resolutions, e.g., $512 \times 512$, $256 \times 256$ and $128 \times 128$. Imuge+ can be applied to images with varied resolutions using a proper model. The results on different resolutions are close. Therefore, the following experiments are conducted on images sized $512 \times 512$. The immunized images are saved in BMP format and the attacked images are randomly saved in common formats such as JPEG, BMP or PNG.

**Human-participated real-world testing.** Imuge+ is tested with human-participated real-world attacks where we invite several volunteers to manipulate the immunized images manually using Adobe Photoshop or Microsoft 3D Drawing. The images for testing are from the testing part of the datasets used for training. Note that although there are already many off-the-shelf datasets with manipulated images, such as CASIA [121] and DEFACTO [122], Imuge+ requires image protection ahead of tampering detection. Therefore, we have to first immunize intact images and request volunteers to manipulate them. In Table 1, we summarize the composition of the hand-crafted test set. We group the real-world tampered images according to the tampering rate $r_T$. The rate of the summation of the tampered area versus the whole image is roughly $r_{sum} \in [0.1, 0.5)$, and the rate of the area of the largest tampering region versus the whole

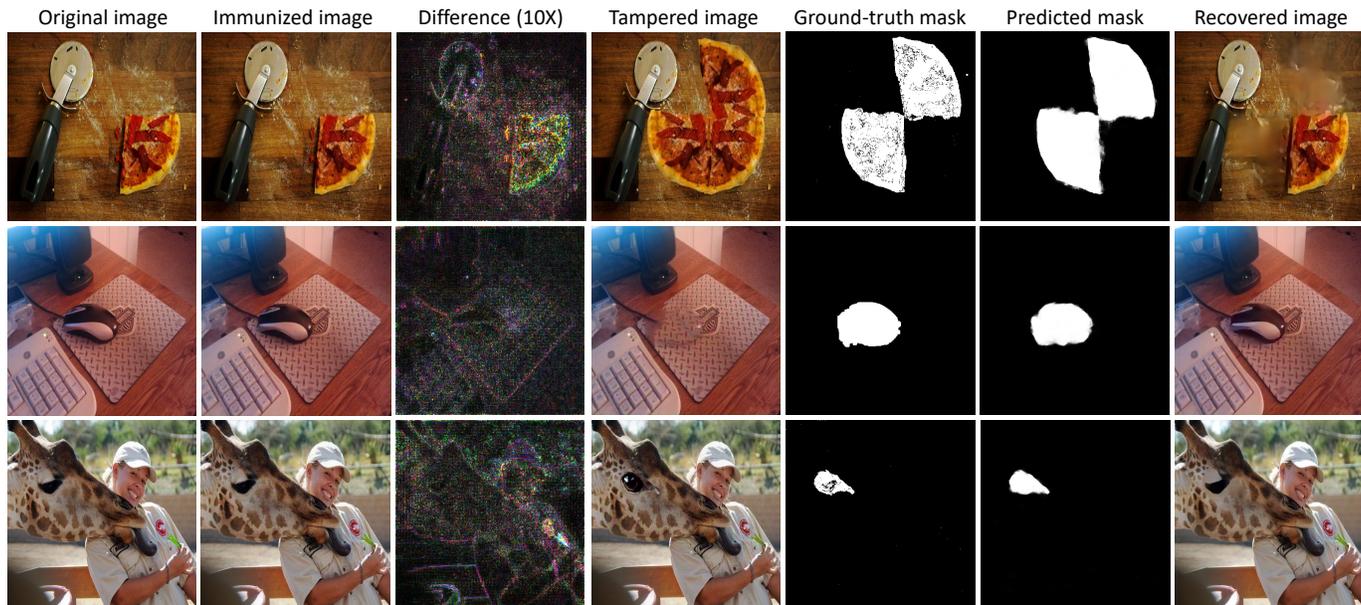| Original image | Immunized image | Difference (10X) | Tampered image | Ground-truth mask | Predicted mask | Recovered image |



Fig. 6. **Performance against combined real-world attacks.** First row: copy-move + JPEG (QF=80). Second row: inpainting + rescaling. Third row: splicing + Gaussian blurring. We successful conduct forgery localization and self-recovery using Imuge+.

TABLE 2
**Comparison of computational complexity. FLOPs: amount of floating point arithmetics. MAdd: amount of multiply-adds. MemR+W: amount of read-write memory.**

| Method | Params | Memory | MAdd | FLOPs | MemR+W |
|---|---|---|---|---|---|
| Imuge+ | 32.0M | 5400MB | 1.53T | 0.76T | 10.92GB |
| Imuge [74] | 13.3M | 1775MB | 0.14T | 0.28T | 3.68GB |
| MantraNet [94] | 3.84M | 4706MB | 2.01T | 1.01T | 8.18GB |
| MVSS-Net [44] | 142.7M | 1377MB | 0.32T | 0.16T | 3.32GB |

TABLE 3
**Average PSNR and SSIM between the protected images and the original images under different resolutions on MS-COCO.**

| Dataset | $512 \times 512$ | | $256 \times 256$ | | $128 \times 128$ | |
|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| COCO | 33.31 | 0.945 | 33.96 | 0.950 | 34.26 | 0.959 |
| ILSVRC | 33.48 | 0.944 | 34.13 | 0.951 | 34.47 | 0.960 |
| Places | 33.15 | 0.940 | 33.77 | 0.949 | 34.02 | 0.953 |
| UCID | 32.96 | 0.937 | 33.45 | 0.945 | 33.92 | 0.952 |

image $r_{max} \in [0.1, 0.25)$, which is generally in line with that of CASIA [121] and DEFACTO [122]. The testing set is composed of 3091 images in total. Fig. 5 shows three examples.

**Evaluation metrics.** We employ the peak signal-to-noise ratio (PSNR), the structural similarity (SSIM) [66] to evaluate the image quality, and the F1 score to measure the accuracy of tampering localization.

**Benchmark.** We compare Imuge+ with Imuge [74] to validate the improvement in the performance of image immunization. Additionally, there are several off-the-shelf schemes for image tampering localization. We employ three state-of-the-art schemes, which detect universal image manipulations, namely, MVSS-Net [44], SPAN [35] and ManTra-Net [94].

**Computational complexity.** In Table 2, we analyze the computational complexity of Imuge+ and compare it with Imuge, MVSS-Net and Mantra-Net. First, Imuge+ requires three times and eight times more parameters than Imuge and Mantra-Net, but much less than MVSS-Net and many other Transformer-based vision pretraining models such as Swin Transformer [2]. Besides, the efficiency of Imuge+ is much improved compared to Imuge, where the memory cost and the amount of floating point arithmetics are

comparable with Mantra-Net and MVSS-Net. Therefore, the computational complexity of Imuge+ is affordable.

### 4.2 Real-world performances

In Fig. 6, we randomly select three test images from MS-COCO test set. We first immunize the images and respectively invite volunteers to conduct combined attacks on them. Then Imuge+ locates the tampered areas and recovers the original image.
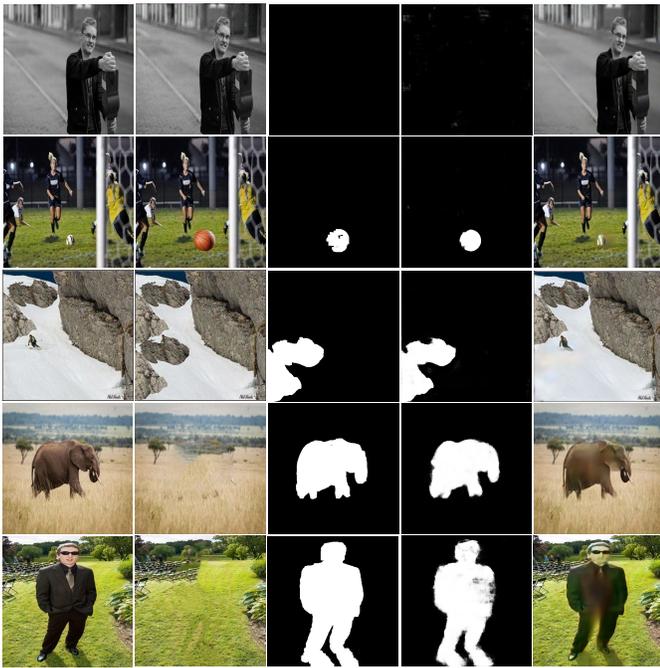
**Image quality of immunized images.** We can observe that the quality of the immunized images is satisfactory where the differences before and after image immunization are close to imperceptible to human visual systems. We have conducted more embedding experiments and the results are reported in Table 3, where stronger perturbations are required for datasets with textured images. For example, for MS-COCO, the average PSNR and SSIM are respectively $33.51dB$. UCID contains way more textured images, and the PSNR and SSIM slightly drop to $32.96dB$ and 0.955. In Table 3, we also conduct experiments on images with several typical resolutions. Images with smaller size enjoy a higher PSNR after image immunization, and we believe the reason is that less information is required to be self-embedded.

TABLE 4
**Performance of tampering localization and image recovery tested by real-world image tampering attack.** The image tampering localization and image recovery are generally satisfactory and there is no significant performance drop against image post-processing attacks.

| Dataset | Index | No Attack | JPEG | | | Scaling | | | Crop | Blurring | | AWGN | Drop-out |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | QF=90 | QF=70 | QF=50 | 150% | 70% | 50% | | Gaussian | Median | | |
| COCO | F1 | 0.918 | 0.894 | 0.873 | 0.827 | 0.915 | 0.874 | 0.853 | 0.864 | 0.852 | 0.792 | 0.844 | 0.802 |
| | PSNR | 29.93 | 28.87 | 28.33 | 27.82 | 29.92 | 28.45 | 27.44 | 27.61 | 28.63 | 25.65 | 27.68 | 27.13 |
| | SSIM | 0.916 | 0.905 | 0.887 | 0.873 | 0.915 | 0.864 | 0.847 | 0.854 | 0.880 | 0.787 | 0.848 | 0.847 |
| ILSVRC | F1 | 0.903 | 0.889 | 0.863 | 0.834 | 0.855 | 0.870 | 0.834 | 0.818 | 0.826 | 0.784 | 0.815 | 0.742 |
| | PSNR | 30.67 | 29.45 | 28.47 | 27.63 | 29.05 | 29.33 | 28.47 | 27.54 | 28.87 | 26.45 | 27.73 | 26.92 |
| | SSIM | 0.933 | 0.902 | 0.890 | 0.858 | 0.882 | 0.892 | 0.858 | 0.850 | 0.878 | 0.813 | 0.856 | 0.841 |
| CelebA | F1 | 0.925 | 0.904 | 0.874 | 0.822 | 0.912 | 0.873 | 0.844 | 0.838 | 0.824 | 0.855 | 0.840 | 0.788 |
| | PSNR | 31.23 | 29.90 | 28.73 | 27.84 | 30.15 | 29.53 | 29.02 | 28.79 | 28.62 | 27.47 | 28.55 | 29.19 |
| | SSIM | 0.934 | 0.916 | 0.886 | 0.859 | 0.895 | 0.877 | 0.874 | 0.868 | 0.870 | 0.855 | 0.865 | 0.877 |



Immunized image   Attacked image   Ground-truth mask   Predicted mask   Recovered image

Fig. 7. **Performance under various tampering rate and post-processing.** The involved post-processing attacks from row one to five is respectively AWGN, Gaussian Blur, resizing, JPEG compression (QF=80) and JPEG compression (QF=90).
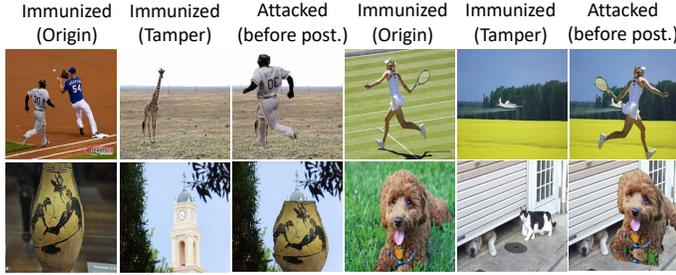
TABLE 5
**Average performance under different tampering rate on MS-COCO.** The data before and after each slash report the performance under no attack and JPEG compression (QF=70).
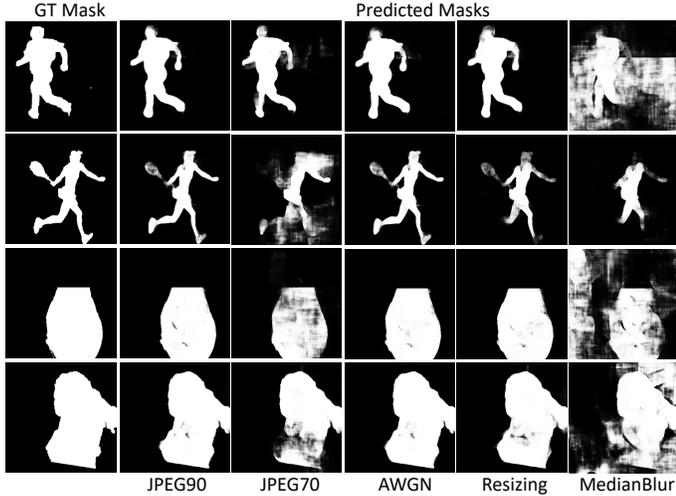
| $r_T$ | Index | Copy-move | Splicing | Inpainting |
|---|---|---|---|---|
| [0,0.1] | PSNR | 30.24 / 28.78 | 30.47 / 28.74 | 30.03 / 28.94 |
| | SSIM | 0.917 / 0.882 | 0.924 / 0.890 | 0.916 / 0.886 |
| | F1 | 0.941 / 0.881 | 0.920 / 0.881 | 0.915 / 0.880 |
| [0.1,0.2] | PSNR | 29.79 / 28.30 | 29.87 / 28.60 | 29.90 / 28.51 |
| | SSIM | 0.910 / 0.879 | 0.917 / 0.883 | 0.915 / 0.880 |
| | F1 | 0.918 / 0.874 | 0.916 / 0.873 | 0.910 / 0.871 |
| [0.2,0.3] | PSNR | 29.57 / 28.02 | 29.23 / 28.32 | 29.53 / 27.91 |
| | SSIM | 0.908 / 0.871 | 0.898 / 0.867 | 0.910 / 0.871 |
| | F1 | 0.910 / 0.861 | 0.883 / 0.861 | 0.890 / 0.843 |
| [0.3, 0.5] | PSNR | 28.97 / 27.67 | 28.69 / 27.73 | 28.56 / 27.36 |
| | SSIM | 0.892 / 0.858 | 0.889 / 0.855 | 0.897 / 0.858 |
| | F1 | 0.883 / 0.844 | 0.875 / 0.854 | 0.875 / 0.822 |

In Fig. 7, we showcase more examples where the tampering rate $r_T$ varies from zero to near 0.5. Also, these images have gone through different image post-processing methods, which are marked below the last row. When the rate is zero, we test the false alarm rate of Imuge+ where the network should identically output the provided image. The predicted mask shows that Imuge+ can evade predicting non-tampering pixels as positive even though additive AWGN is added. Besides, in the subsequent tests, the volunteers conducted the copy-move, splicing or inpainting attack respectively which result in the removal of certain non-trivial objects in the original image. From the results, we see that the tampered areas are correctly classified from the whole image plane where the borders are largely consistent. Even provided with extreme cases where $r_T > 0.4$, Imuge+ can still recover the missing objects, though some higher-band details are lost.

**Quantitative analysis.** In Table 4, we clarify the performance of Imuge+ with the presence of different image post-processing attacks. In Table 5, we show the average performance for different tampering rate and different image post-processing attacks. Here for inpainting, the volunteers use online demo sites [11], free tools [5] or open-source models [7]. The tampering rate ranges from zero to 0.5

**Qualitative analysis.** As a successful localization of tampered areas is a prerequisite of successful image recovery, Fig. 6 further shows the satisfactory results of tampering localization as well as image recovery. The predicted masks are close to the ground truth and the tampered contents are correctly reconstructed. In these examples, the added pizza is the result of copy-moving, and the missing mouse and eye of the giraffe is respectively removed by image inpainting and splicing. Besides, the condition of image post-processing applied in the three examples varies, which involve the famous compression, blurring and rescaling operations. However, the diversity of the hybrid attacks still cannot prevent Imuge+ from recovering the original contents, though some details might be lost. It proves the robustness of our scheme.

(a) Examples of coincident forgery, i.e., both the origin and the tampering source are immunized images.
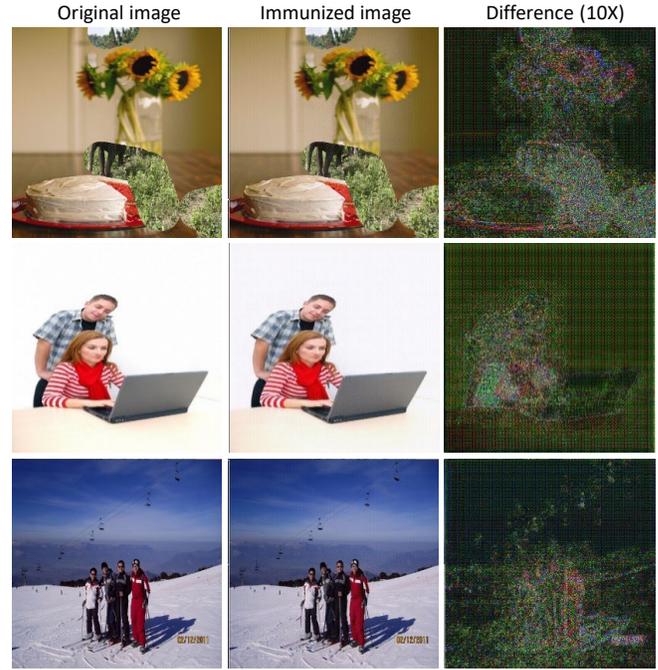


(b) Localization results on multiple types of post-processed versions of the manipulated images in (a).

Fig. 8. **Performance of Imuge+ on localization against coincident forgery.** Owing to the fact that the added items from the tampering source in most real-world application will be resized/rotated/flipped for realistic forgery, the immunization signals on them are therefore largely weakened, resulting in successful localization.

whose distribution follows that reported in Table 1. Concerning image cropping attack, Imuge+ only works within the cropped region.

From the results, we can observe that the performance of both image tampering localization and image self-recovery are successful, where the F1 score and the SSIM score are both above 0.9. Besides, Imuge+ shows strong robustness against common image post-processing behaviors in that the performances do not drop significantly despite the presence of heavy attacks such as Median blurring or AWGN addition. For example, the F1 score and the SSIM score of the worst averaged performance is near 0.8, where median blurring filters out the higher-band details of the images. In most cases, the F1 scores are between 0.8 and 0.9 and the SSIM score between 0.85 and 0.9. Third, it is proven by the provided results that Imuge+ can effectively work on images from different distributions, where the test images are from the most-commonly-used image datasets. Interestingly, the performances against copy-move and splicing are slightly better than that against inpainting. The reason might be that the simulation for copy-move and splicing might be more effective than that for inpainting. Besides, the performances do not drop significantly with the increase of the tampering rate.

**Coincident forgery.** We additionally conduct experiments



Fig. 9. **Analysis of the composition of immunized images.** There are two types of hidden patterns. The chess-board pattern is believed to be responsible for tampering detection, and the other is the compressed version of the original image for image self-recovery.

to test the performance of tampering localization where the tampering source used for splicing is coincidentally another immunized image. We find that in the most common situations of image manipulation, the attacker always either resizes, flips, rotates or spatially transforms the forgery content before adding them onto the victim image, in order to produce plausible manipulation. We simulate this manipulation using two immunized images. Figure 8 provides four examples, where we observe that the coincident forgery can be detected, even though we have not included coincident forgery in the attack simulation stage. Notice that even if the tampering rates in the last two rows exceed 50%, the added items are not mistakenly predicted as original by the forgery detector. The reason is that the above-mentioned distortion operations will inevitably weaken or destroy the immunized signal inside the forgery content.

### 4.3  What is in the immunized images?

To analyze how Imuge+ locates the tampered areas as well as recovers the received image, we show three more examples of the immunized images in Fig. 9 and especially have a closer observation of the augmented residual images. First, we find that in the smooth or empty areas, a checkerboard pattern is introduced by IINet, which is not so remarkable in the textured areas. We believe this chess-board pattern is embedded mainly for tampering localization, since when the immunized images are tampered, the learnt pattern is very unique and out of the distribution of natural images, which can hardly be forged or constructed by tampering attacks. Even for copy-move attacks, it will result in pattern inconsistency inside the local areas. Besides, the chess-board pattern can somehow survive JPEG compression, resizing
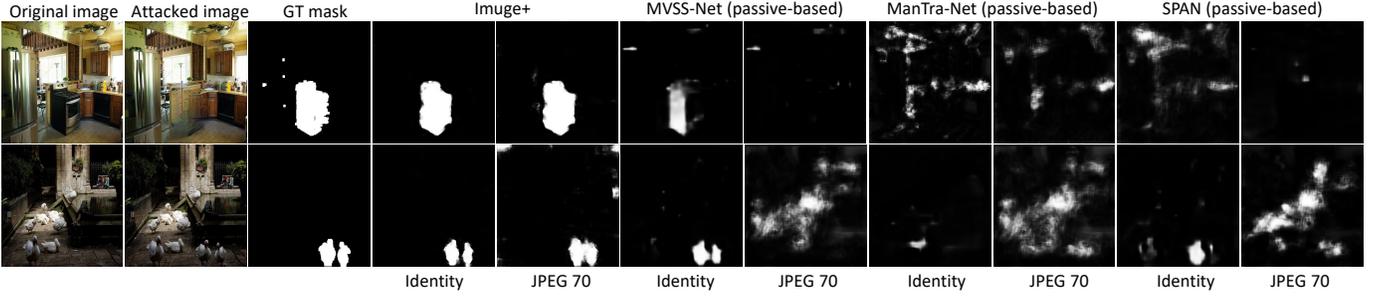
Fig. 10. **Comparison of tampering localization among Imuge+ and several state-of-the-art schemes.** Upper: inpainting. Lower: copy-move. Imuge+ can accurately localize the tampered areas even with the presence of post-processing attack. In contrast, many state-of-the-art schemes are reported to not have robustness.
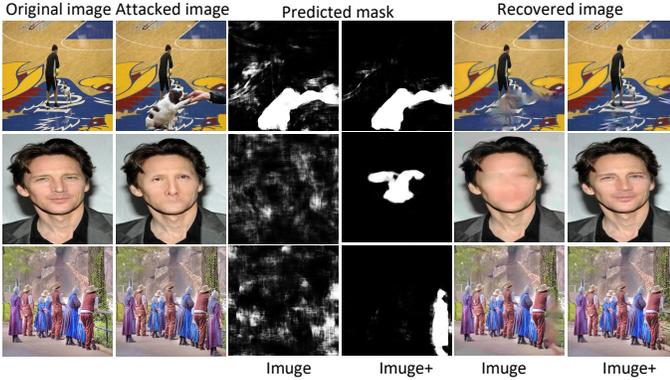


Fig. 11. **Performance comparison with Imuge.** We can observe a noticeable performance boost in all three attacks. The attacked images are stored in JPEG format with QF=70.

TABLE 6
**Performance comparison with Imuge on tampering localization and image recovery**. The performance of [74] is in the brackets.

| Attack | Index | Copy-move | Splicing | Inpainting |
|---|---|---|---|---|
| No | PSNR | 29.67 (27.72) | 30.13 (28,87) | 29.06 (24.33) |
| Attack | F1 | 0.924 (0.627) | 0.903 (0.754) | 0.912 (0.431) |
| Gaussian | PSNR | 28.24 (24.49) | 28.77 (25.23) | 28.04 (23.52) |
| Blur | F1 | 0.843 (0.531) | 0.866 (0.637) | 0.853 (0.333) |
| JPEG | PSNR | 28.42 (25.37) | 28.51 (25.68) | 28.03 (23.75) |
| QF=70 | F1 | 0.854 (0.582) | 0.881 (0.654) | 0.837 (0.295) |

and several kinds of blurring. Therefore, the forgery detector can detect the existence of such a pattern to determine which part of the image is forged. In contrast, state-of-the-art passive forensics schemes have to find a universally-present tampering trace, which is much harder than detecting an embedded tailored pattern. Besides, we can clearly observe that the embedded patterns near the border of the objects are different from the pattern for localization. For example, the pizza in the first example and the mouse pad in the second example are compressed into residual information and scattered around the nearby areas. We believe that on recovering the image, Imuge+ must have learned to check the surrounding area for residual information. Finally, these two kinds of patterns can harmoniously co-exist, in that we can also localize the tampered areas though the second kind of pattern is stronger in textured areas, and we believe only if the second kind of pattern is weak, will Imuge+ embed the first kind of pattern into the image.

### 4.4 Comparison

**Content recovery within tampered areas.** We compare Imuge+ with [74] to verify the performance boost brought by our enhanced network design and novel training mechanisms. The averaged PSNRs of the two methods between $\mathbf{I}$ and $\mathbf{X}$ are kept close for a fair comparison. Besides, in order to compare the overall quality of image recovery, assume that the tampering localization is correct in measuring the performance of image recovery.

In Fig. 11, we show three groups of experimental comparison between Imuge and Imuge+, where we perform splicing attack in the first row, inpainting attack in the second and copy-move attack in the third row. In Table 6, we provide the averaged results of the comparison. Notice that the face of the man in the second example is completely removed, but Imuge+ can successfully recover the original face. Besides, Imuge is not trained against image inpainting attack and therefore it performs poorly in the second example. From the averaged results, where we see that the performance boost is noticeable. First, Imuge+ is more accurate in localizing the tampers. Second, Imuge+ can preserve more detail of the missing objects. Third, Imuge+ is stronger in overall robustness against various kinds of image post-processing attacks.

**Image tampering localization.** Passive image forensics schemes do not rely on hiding additional information, but they usually find tiny traces that are vulnerable to traditional image post-processing, so the performances can deteriorate in real-world OSN applications. Though tampering localization in Imuge+ requires an additional and mandatory procedure of information hiding, our scheme aims at replacing original images which are prone to a variety of hybrid tampering and post-processing attacks with their immunized versions. Therefore, by restricting the magnitude of embedded perturbation, we wish the immunized images to be viewed as *original images* in the future.

To better evaluate the accuracy of tampering localization of our scheme, here we briefly compare the performance with those of the state-of-the-art passive image tampering localization schemes. Note that for [23], [44], [94], we conduct fair comparisons by performing the same attacks on the original images, not the immunized images. Fig. 10 shows

TABLE 7
**F1 score comparison for tampering detection among our scheme and the state-of-the-art methods.** Our scheme ranks first in all the tests and leads by a large margin in robustness.

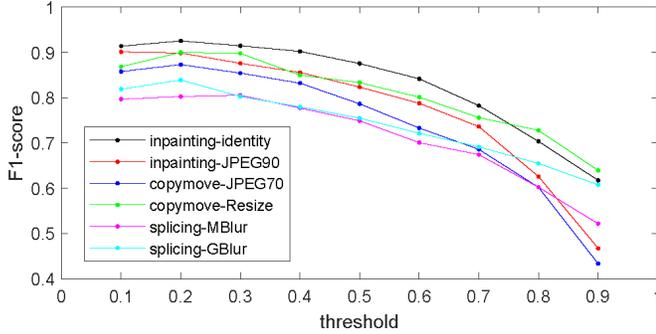| Method | Inpainting | | | | Splicing | | | | Copy-move | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | NoAttack | JPEG | Blurring | Scaling | NoAttack | JPEG | Blurring | Scaling | NoAttack | JPEG | Blurring | Scaling |
| Imuge [74] | 0.372 | 0.253 | 0.219 | 0.327 | 0.754 | 0.712 | 0.679 | 0.747 | 0.642 | 0.463 | 0.453 | 0.566 |
| Mantra-Net [94] | 0.584 | 0.540 | 0.517 | 0.559 | 0.678 | 0.563 | 0.504 | 0.539 | 0.553 | 0.477 | 0.423 | 0.511 |
| MVSS-Net [44] | 0.673 | 0.487 | 0.445 | 0.570 | 0.787 | 0.659 | 0.535 | 0.672 | 0.701 | 0.627 | 0.551 | 0.675 |
| SPAN [35] | 0.701 | 0.624 | 0.599 | 0.683 | 0.732 | 0.657 | 0.574 | 0.665 | 0.685 | 0.614 | 0.565 | 0.648 |
| Imuge+ | **0.917** | **0.872** | **0.854** | **0.885** | **0.923** | **0.878** | **0.849** | **0.873** | **0.903** | **0.844** | **0.829** | **0.837** |



Fig. 12. **Performance curves of F1 score versus threshold.** Generally, F1 reaches the highest peak when the threshold is $0.2$.
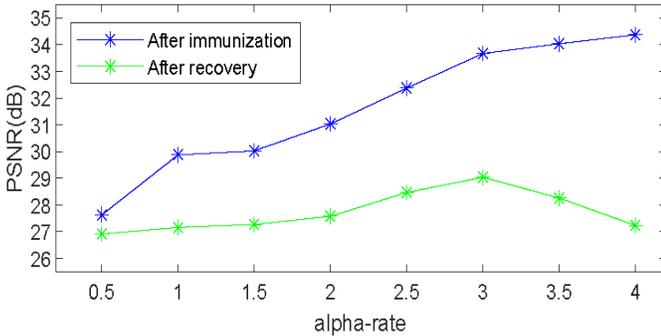


Fig. 13. **Performance curves of PSNR versus $\alpha$.** We observe that the trade-off between imperceptibility and network effectiveness can be best balanced when $\alpha$ equals to $3$.

TABLE 8
**Ablation study of Imuge+ using varied partial settings.** The tests are done under JPEG attack (QF=70) and $r_T \in [0.1, 0.4]$. 'IP': image protection (forward pass of IINet). 'E': including edge as additional input and supervision. 'KD': using KD-JPEG as JPEG simulator. 'TDA': tampering-based DA. 'IT': iterative training. 'AB': asymmetric batch size. $\Delta^1$: using Diff-JPEG [34] as JPEG simulator. $\Delta^2$: using two separate U-Net [59] to replace IINet. '-': failed to train steadily.

| Network components | | | | | | Index | | |
|---|---|---|---|---|---|---|---|---|
| IP | E | KD | TDA | IT | AB | F1 | PSNR | SSIM |
| | ✓ | ✓ | ✓ | ✓ | ✓ | 0.435 | 18.13 | 0.585 |
| $\Delta^1$ | ✓ | ✓ | ✓ | ✓ | ✓ | 0.933 | 21.54 | 0.713 |
| ✓ | | ✓ | ✓ | ✓ | ✓ | 0.854 | 24.57 | 0.784 |
| ✓ | ✓ | $\Delta^2$ | ✓ | ✓ | ✓ | 0.962 | 23.78 | 0.755 |
| ✓ | ✓ | ✓ | | ✓ | ✓ | 0.435 | 24.63 | 0.744 |
| ✓ | ✓ | ✓ | ✓ | | ✓ | - | - | - |
| ✓ | ✓ | ✓ | ✓ | ✓ | | 0.913 | 25.59 | 0.810 |
| ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | **0.965** | **26.41** | **0.822** |

### 4.5 Ablation Study

We discuss the selection of the hyper-parameters and the threshold for binarizing the prediction mask, respectively in Fig. 12 and Fig. 13. Besides, we explore the influence of the key components in Imuge+ by evaluating the performance of the model with varied and partial setups. In each test, we remove a single component and train the models from scratch till they converge. Fig. 14 visualizes two groups of performance comparison under hybrid attacks. We also summarize the average result on partial setups in Table 8 by performing the same image-processing attacks and the same tampering attacks in each test.

**Selection of the critical hyper-parameters.** In Fig. 12, we study the selection of proper threshold for predicted mask binarization $Th$ and rate of data augmentation $r_{aug}$. From the curves, we find that when $Th$ is $0.2$, we can generally acquire the best performance. Besides, we also empirically find that the performance can be further enhanced when $r_{aug}$ is set as $15\%$. In Fig. 13, we also study the trade-off between the imperceptibility of image protection and the effectiveness of image recovery. Generally, when the PSNR surpasses $33dB$, there will be little easily-noticeable artifacts caused by the immunization. According to both the human visual system and the figure, $\alpha$ is set as three for the best equilibrium. To begin with, we verify that directly conducting forgery detection and image recovery on the attacked *unimmunized* images leads to sub-optimal results.

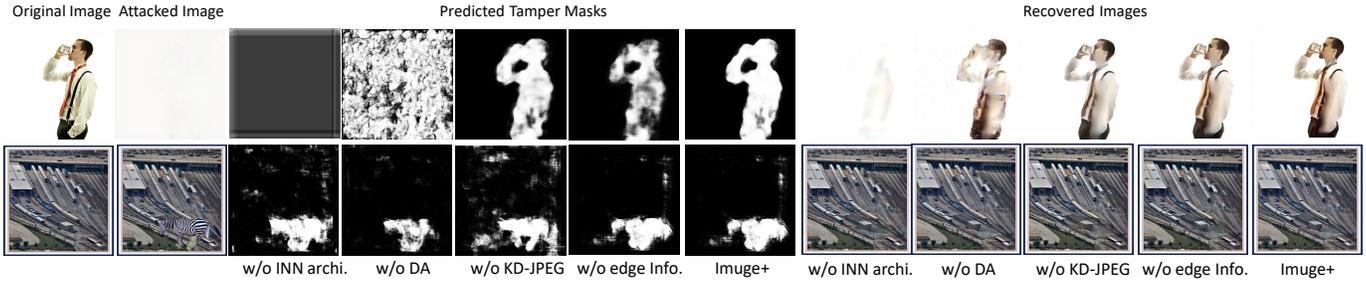**Influence of the JPEG simulator.** Fig 15 shows the com-

two comparison results where we respectively involve inpainting and copy-move attacks.

According to the results, Imuge+ provides the leading performance which is robust to the JPEG compression attack with QF=70. In contrast, passive methods encounter a significant performance drop. In Table 7, we further show the averaged accuracy of the above-mentioned schemes. Generally, the F1 scores of Imuge+ are above 0.8, which shows high resilience against JPEG compression, blurring and scaling. The overall performance is significantly improved compared to [74] which performs only fair against copy-move attack and poorly against inpainting. The testing results of [23], [44], [94] are consistent with those recorded in these papers. The average F1 score of these methods under no attack is within [0.55,0.8], and the performances under image post-processing attacks drop by 0.1 to 0.2.

Fig. 14. **Ablation study of Imuge+.** In each group, the performed attack is the same, with JPEG QF=80. Among all of the setups, Imuge+ with full setup performs the best in both tampering localization and image recovery.
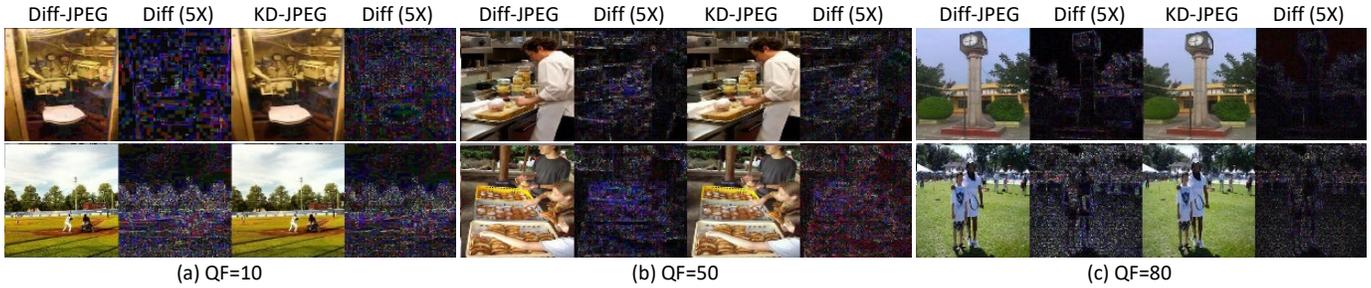


Fig. 15. **Comparison on the fidelity of the simulated JPEG images.** Owing to the flexibility of the generative network, KD-JPEG in average contains less distortion compared to Diff-JPEG [34].

parison between KD-JPEG and Diff-JPEG [34] for JPEG simulation. In each group, we respectively compare the absolute difference between the generated JPEG images and the real-world JPEG images. Our generative method gives closer results to the real-world JPEG in that the differences are smaller and the characteristics of JPEG compression are better studied. The checkerboard artifact of JPEG compression can also be found in the generated JPEG images. On average, the PSNR between the generated JPEG images by KD-JPEG and the real-world JPEG images is $29.64dB$ when QF is 10, and that of Diff-JPEG is $28.03dB$. Moreover, the PSNR of KD-JPEG is $31.17dB$ when QF is 50 and $34.42dB$ when QF is 90. In comparison, Diff-JPEG shows less flexibility in the simulation where the PSNR is $30.22dB$ when QF is 50 and $33.50dB$ when QF is 90. Besides, the QF classification accuracy of images generated by KD-JPEG is $95.47\%$, while that of the differentiable method Diff-JPEG [34] $72.48\%$. It indicates that more feature representations of the JPEG images can be learned by KD-JPEG. Fig. 14 shows that applying Diff-JPEG leads to more blurry results.

**Influence of the INN-based architecture.** The benefit of introducing INN-based architectures for invertible function modeling has been well studied by many researches [103], [105]. These networks learn deterministic and invertible distribution mapping, where the forward and back propagation operations are in the same network. A typical alternative is to model image protection and image recovery independently using the well-known "Encoder-Decoder" architecture. From the results, we observe that INN-based Imuge+ provides much better performance. The reason is that INN has much fewer hyper-parameters and therefore the training process is stable compared to GAN training.

**Influence of the data augmentation.** We find that without using our tampering-based data augmentation, the network

during training still tends to use DIP for image recovery. The reason is that the randomly generated masks may catch the plain-text areas where the contents inside are trivial compared to the surroundings. In the first example of Fig. 14, we see that without the novel data augmentation, Imuge+ can still somehow recover the removed person, suggesting that even if without disabling semantic hints, Imuge+ still reconstructs the images using the hidden information. However, we see a much lower quality and fidelity compared to the full implementation. As we introduce irrelevant information and force the network to tamper and recover them, Imuge+ cannot always rely on DIP and tries to embed more essential information for reliable recovery.

**Influence of other components.** Previous work [7] has proven that the recovery of intermediate representation, like edges and the gray-scaled version, can boost the performance of image reconstruction. Compared to [74], Imuge+ additionally embeds and recovers the edge information to enforce that the recovered results are semantically and trustfully correct. According to the ablation studies, without introducing the edge supervision, Imuge+ tends to produce more blurry recovered images, while there is no noticeable effect on the tampering localization. Besides, we find that without using iterative training, we cannot steadily train Imuge+ where poorly predicted tampering masks heavily disrupts the image recovery stage. Also, performing the asymmetric-batch-size technique can help noticeable promoting the overall performance.

## 4.6 Applications

Image immunization is tested on thousands of hand-crafted tampered images. It is designed for real-world automatic image tampering localization and self-recovery. In Fig. 16,

Original image  Attacked image  Predicted mask  Recovered image

(a) Prevention of copy-right information removal

(b) Prevention of facial image manipulation by DeepFake

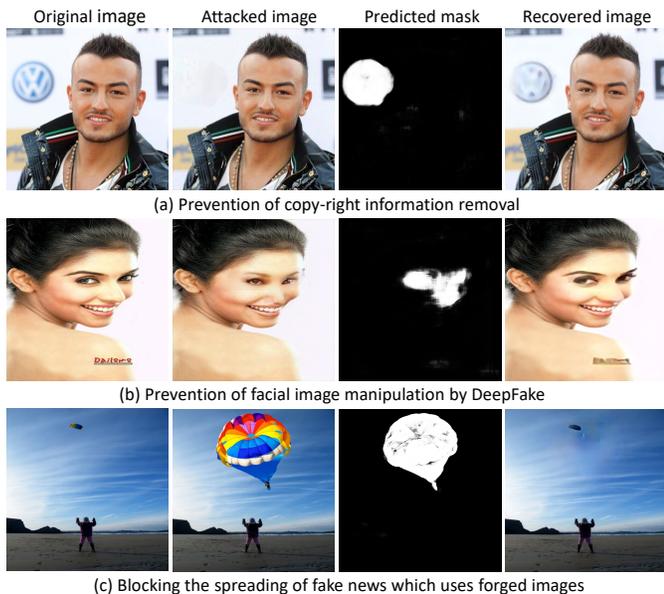(c) Blocking the spreading of fake news which uses forged images

Fig. 16. **Applications of Imuge+ against fake images.** The proposed method is effective in both localizing the tampered area and recovering the original content.

we exemplify three applications of Imuge+ for different purposes in image forensics.

**Prevention of copy-right information removal.** Several image tampering attacks are based on the removal of critical information, such as signatures, logos, etc. Such behaviors are a serious violation of copy-right protection and image fidelity. Using image immunization, Imuge+ can benefit the prevention of their removal. For example, in the first row of the figure, the logo of Volkswagen was removed by the attacker and later identified as well as recovered in high quality by Imuge+.

**Prevention of facial image manipulation by DeepFake.** The second example is to combat image manipulation by DeepFake. Facial images are also prone to modifications that illegally shift the identity of the person involved. Such tampering can be easily accomplished by GAN-based technologies and result in severe security issues. In the second row, Imuge+ predicts that the face was forged by image inpainting technique and successfully reconstructs the identifiable original face. In that sense, Imuge+ also blocks the spreading of fake news which uses forged images as supplementary multimedia.

**Reversible image editing.** Normally, after we edit an image, a copy of the original version has to be stored if we wish to revert the unwanted editing later. However, storing every version as a backup after each tiny modification might be expensive. In the third example, the data owner removes the added parachute, and uses Imuge+ to revert the editing without having stored the original image. In the recovered image, the missing tiny parachute is also recovered. Therefore, the users are no longer required to store the original version of an image after modification. We leave it a future work to develop a repeatable image immunization that can revert each modification individually with only one immunized version stored.

## 5 CONCLUSIONS AND FUTURE WORKS

In this paper, we present a novel generative scheme called Imuge+, which is an image tamper resilient generative scheme for image self-recovery. We transform the original images into immunized images, where the tamper attacks can be accurately localized and the original content within the tampered areas can be recovered. To boost the performance, we form the invertible function for image immunization and employ an INN architecture for implementation. Besides, we propose a novel JPEG simulator as well as an enhanced attack layer for greater resilience against common image post-processing attacks. We conduct comprehensive experiments on several popular datasets and invite several volunteers to manipulate the immunized images, and the results prove the effectiveness of Imuge+ in both tamper localization and content recovery against splicing, copy-move and inpainting attacks.

There are still some remaining issues to be addressed in future works. First, the imperceptibility of the information embedded required by image immunization can be further improved. We find that we cannot immunize an immunized image, otherwise the corresponding generated image will be disastrous in quality. Second, though we have made steady improvements in the overall quality of the recovered images, the blurry issue still exists in many cases, e.g., the attacked images are heavily compressed, or a single tampered region is too big resulting in the center area not recoverable. Third, IMUGE+ is still largely at black-bos level, and therefore, we wish to conduct more theoretical analysis, especially on the upper bound for image immunization against common types of OSN attack. We hope that in the near future, the above-mentioned issues can be addressed well, and Imuge+ can be integrated into cameras, so that image immunization can be introduced into the image signal processing pipeline, thereby changing the current situation where digital images can be freely edited.

## REFERENCES

[1] I. Yerushalmy and H. Hel-Or, "Digital image forgery detection based on lens and sensor aberration," *International Journal of Computer Vision (IJCV)*, vol. 92, no. 1, pp. 71–91, 2011.

[2] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 012–10 022.

[3] R. Suvorov, E. Logacheva, A. Mashikhin, A. Remizova, A. Ashukha, A. Silvestrov, N. Kong, H. Goka, K. Park, and V. Lempitsky, "Resolution-robust large mask inpainting with fourier convolutions," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2022, pp. 2149–2159.

[4] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *arXiv preprint arXiv:1511.07289*, 2015.

[5] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Generative image inpainting with contextual attention," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5505–5514.

[6] N. Wang, J. Li, L. Zhang, and B. Du, "Musical: Multi-scale image contextual attention learning for inpainting." in *IJCAI*, 2019, pp. 3748–3754.

[7] K. Nazeri, E. Ng, T. Joseph, F. Qureshi, and M. Ebrahimi, "Edgeconnect: Structure guided image inpainting using edge prediction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops (ICCVW)*, 2019, pp. 0–0.

[8] P. Zhou, X. Han, V. I. Morariu, and L. S. Davis, "Learning rich features for image manipulation detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1053–1061.

[9] ——, "Learning rich features for image manipulation detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 1053–1061.

[10] MATLAB, *9.7.0.1190202 (R2019b)*. Natick, Massachusetts: The MathWorks Inc., 2018.

[11] J. Yu, Z. Lin, J. Yang, X. Shen, X. Lu, and T. S. Huang, "Free-form image inpainting with gated convolution," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 4471–4480.

[12] T. Wang, H. Ouyang, and Q. Chen, "Image inpainting with external-internal learning and monochromic bottleneck," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 5120–5129.

[13] P. Zhuang, H. Li, S. Tan, B. Li, and J. Huang, "Image tampering localization using a dense fully convolutional network," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 16, pp. 2986–2999, 2021.

[14] H. Li and J. Huang, "Localization of deep inpainting using high-pass fully convolutional network," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 8301–8310.

[15] Y. Jo, S. Y. Chun, and J. Choi, "Rethinking deep image prior for denoising," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 5087–5096.

[16] H. Guan, M. Kozak, E. Robertson, Y. Lee, A. N. Yates, A. Delgado, D. Zhou, T. Kheyrkhah, J. Smith, and J. Fiscus, "Mfc datasets: Large-scale benchmark datasets for media forensic challenge evaluation," in *2019 IEEE Winter Applications of Computer Vision Workshops (WACVW)*. IEEE, 2019, pp. 63–72.

[17] Y.-F. Hsu and S.-F. Chang, "Detecting image splicing using geometry invariants and camera characteristics consistency," in *2006 IEEE International Conference on Multimedia and Expo*. IEEE, 2006, pp. 549–552.

[18] X. Liao, K. Li, X. Zhu, and K. R. Liu, "Robust detection of image operator chain with two-stream convolutional neural network," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 5, pp. 955–968, 2020.

[19] C. Yang, H. Li, F. Lin, B. Jiang, and H. Zhao, "Constrained r-cnn: A general image manipulation detection model," in *2020 IEEE International Conference on Multimedia and Expo (ICME)*. IEEE, 2020, pp. 1–6.

[20] X. Zhu, Y. Qian, X. Zhao, B. Sun, and Y. Sun, "A deep learning approach to patch-based image inpainting forensics," *Signal Processing: Image Communication*, vol. 67, pp. 90–99, 2018.

[21] Y. Li and J. Zhou, "Fast and effective image copy-move forgery detection via hierarchical feature point matching," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 14, no. 5, pp. 1307–1322, 2018.

[22] R. Salloum, Y. Ren, and C.-C. J. Kuo, "Image splicing localization using a multi-task fully convolutional network (mfcn)," *Journal of Visual Communication and Image Representation*, vol. 51, pp. 201–209, 2018.

[23] M.-J. Kwon, I.-J. Yu, S.-H. Nam, and H.-K. Lee, "Cat-net: Compression artifact tracing network for detection and localization of image splicing," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, 2021, pp. 375–384.

[24] Y. Wu, W. Abd-Almageed, and P. Natarajan, "Deep matching and validation network: An end-to-end solution to constrained image splicing localization and detection," in *Proceedings of the 25th ACM international conference on Multimedia*, 2017, pp. 1480–1502.

[25] Y. Rao and J. Ni, "A deep learning approach to detection of splicing and copy-move forgeries in images," in *2016 IEEE International Workshop on Information Forensics and Security (WIFS)*. IEEE, 2016, pp. 1–6.

[26] D. Cozzolino, G. Poggi, and L. Verdoliva, "Efficient dense-field copy–move forgery detection," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 10, no. 11, pp. 2284–2297, 2015.

[27] T. F. van der Ouderaa and D. E. Worrall, "Reversible gans for memory-efficient image-to-image translation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4720–4728.

[28] Y. Wang, M. Xiao, C. Liu, S. Zheng, and T.-Y. Liu, "Modeling lost information in lossy image compression," *arXiv preprint arXiv:2006.11999*, 2020.

[29] A. Lugmayr, M. Danelljan, L. Van Gool, and R. Timofte, "Srflow: Learning the super-resolution space with normalizing flow," in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 715–732.

[30] L. Ardizzone, C. Lüth, J. Kruse, C. Rother, and U. Köthe, "Guided image generation with conditional invertible neural networks," *arXiv preprint arXiv:1907.02392*, 2019.

[31] B. Boehm, "Stegexpose-a tool for detecting lsb steganography," *arXiv preprint arXiv:1410.6656*, 2014.

[32] J.-H. Jacobsen, A. Smeulders, and E. Oyallon, "i-revnet: Deep invertible networks," *arXiv preprint arXiv:1802.07088*, 2018.

[33] J. Behrmann, W. Grathwohl, R. T. Chen, D. Duvenaud, and J.-H. Jacobsen, "Invertible residual networks," in *International Conference on Machine Learning (ICML)*, 2019, pp. 573–582.

[34] R. Shin and D. Song, "Jpeg-resistant adversarial images," in *NIPS 2017 Workshop on Machine Learning and Computer Security*, vol. 1, 2017.

[35] X. Hu, Z. Zhang, Z. Jiang, S. Chaudhuri, Z. Yang, and R. Nevatia, "Span: Spatial pyramid attention network for image manipulation localization," in *European Conference on Computer Vision (ECCV)*. Springer, 2020, pp. 312–328.

[36] Y. Bengio, N. Léonard, and A. Courville, "Estimating or propagating gradients through stochastic neurons for conditional computation," *arXiv preprint arXiv:1308.3432*, 2013.

[37] C. Zhang, A. Karjauv, P. Benz, and I. S. Kweon, "Towards robust deep hiding under non-differentiable distortions for practical blind watermarking," in *Proceedings of the 29th ACM International Conference on Multimedia*, 2021, pp. 5158–5166.

[38] D.-A. Clevert, T. Unterthiner, and S. Hochreiter, "Fast and accurate deep network learning by exponential linear units (elus)," *arXiv preprint arXiv:1511.07289*, 2015.

[39] B. Bayar and M. C. Stamm, "Constrained convolutional neural networks: A new approach towards general purpose image manipulation detection," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 13, no. 11, pp. 2691–2706, 2018.

[40] J. Tao, S. Li, X. Zhang, and Z. Wang, "Towards robust image steganography," *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, vol. 29, no. 2, pp. 594–600, 2018.

[41] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille, "Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 40, no. 4, pp. 834–848, 2017.

[42] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 3431–3440.

[43] V. Sedighi, R. Cogranne, and J. Fridrich, "Content-adaptive steganography by minimizing statistical detectability," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 11, no. 2, pp. 221–234, 2015.

[44] X. Chen, C. Dong, J. Ji, J. Cao, and X. Li, "Image manipulation detection by multi-view multi-scale supervision," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2021, pp. 14 185–14 193.

[45] K. Chen, H. Zhou, W. Zhou, W. Zhang, and N. Yu, "Defining cost functions for adaptive jpeg steganography at the microscale," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 14, no. 4, pp. 1052–1066, 2018.

[46] H. Fang, W. Zhang, H. Zhou, H. Cui, and N. Yu, "Screen-shooting resilient watermarking," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 14, no. 6, pp. 1403–1418, 2018.

[47] N. Agarwal, A. K. Singh, and P. K. Singh, "Survey of robust and imperceptible watermarking," *Multimedia Tools and Applications*, vol. 78, no. 7, pp. 8603–8633, 2019.

[48] M. Asikuzzaman and M. R. Pickering, "An overview of digital video watermarking," *IEEE Transactions on Circuits and Systems for Video Technology (TCSVT)*, vol. 28, no. 9, pp. 2131–2153, 2017.

[49] Y.-Q. Shi, X. Li, X. Zhang, H.-T. Wu, and B. Ma, "Reversible data hiding: advances in the past two decades," *IEEE access*, vol. 4, pp. 3210–3237, 2016.

[50] S. Baluja, "Hiding images in plain sight: Deep steganography," *Advances in Neural Information Processing Systems (NIPS)*, vol. 30, pp. 2069–2079, 2017.

[51] ——, "Hiding images within images," *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019.

[52] X. Duan, K. Jia, B. Li, D. Guo, E. Zhang, and C. Qin, "Reversible image steganography scheme based on a u-net structure," *IEEE Access*, vol. 7, pp. 9314–9323, 2019.

[53] R. Rahim, S. Nadeem *et al.*, "End-to-end trained cnn encoder-decoder networks for image steganography," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 0–0.

[54] S.-M. Mun, S.-H. Nam, H. Jang, D. Kim, and H.-K. Lee, "Finding robust domain from attacks: A learning framework for blind watermarking," *Neurocomputing*, vol. 337, pp. 191–202, 2019.

[55] H. Kandi, D. Mishra, and S. R. S. Gorthi, "Exploring the learning capabilities of convolutional neural networks for robust image watermarking," *Computers & Security*, vol. 65, pp. 247–268, 2017.

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

[57] X. Luo, R. Zhan, H. Chang, F. Yang, and P. Milanfar, "Distortion agnostic deep watermarking," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 548–13 557.

[58] X. Bi, Y. Liu, B. Xiao, W. Li, C.-M. Pun, G. Wang, and X. Gao, "D-unet: A dual-encoder u-net for image splicing forgery detection and localization," *arXiv preprint arXiv:2012.01821*, 2020.

[59] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*. Springer, 2015, pp. 234–241.

[60] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 2818–2826.

[61] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1125–1134.

[62] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European Conference on Computer Vision (ECCV)*. Springer, 2016, pp. 694–711.

[63] L. I. Rudin, S. Osher, and E. Fatemi, "Nonlinear total variation based noise removal algorithms," *Physica D: nonlinear phenomena*, vol. 60, no. 1-4, pp. 259–268, 1992.

[64] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision (ECCV)*. Springer, 2014, pp. 740–755.

[65] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. Ieee, 2009, pp. 248–255.

[66] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing (TIP)*, vol. 13, no. 4, pp. 600–612, 2004.

[67] C. Yu, "Attention based data hiding with generative adversarial networks," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 01, 2020, pp. 1120–1128.

[68] X. Mao, Q. Li, H. Xie, R. Y. Lau, Z. Wang, and S. Paul Smolley, "Least squares generative adversarial networks," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 2794–2802.

[69] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International Conference on Machine Learning (ICML)*, 2015, pp. 448–456.

[70] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in *2017 ieee symposium on security and privacy (sp)*. IEEE, 2017, pp. 39–57.

[71] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.

[72] D. Khachaturov, I. Shumailov, Y. Zhao, N. Papernot, and R. Anderson, "Markpainting: Adversarial machine learning meets inpainting," in *International Conference on Machine Learning (ICML)*, 2021, pp. 5409–5419.

[73] M. Yin, Y. Zhang, X. Li, and S. Wang, "When deep fool meets deep prior: Adversarial attack on super-resolution network," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 1930–1938.

[74] Q. Ying, Z. Qian, H. Zhou, H. Xu, X. Zhang, and S. Li, "From image to imuge: Immunized image generation," in *Proceedings of the 29th ACM international conference on Multimedia*, 2021, pp. 1–9.

[75] K. Liu, D. Chen, J. Liao, W. Zhang, H. Zhou, J. Zhang, W. Zhou, and N. Yu, "Jpeg robust invertible grayscale," *IEEE Transactions on Visualization and Computer Graphics*, 2021.

[76] Z. Jia, H. Fang, and W. Zhang, "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*, 2021, pp. 41–49.

[77] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," *arXiv preprint arXiv:1511.07122*, 2015.

[78] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[79] H. Cheng and X. Shi, "A simple and effective histogram equalization approach to image enhancement," *Digital signal processing*, vol. 14, no. 2, pp. 158–170, 2004.

[80] A. Bruhn, J. Weickert, and C. Schnörr, "Lucas/kanade meets horn/schunck: Combining local and global optic flow methods," *International Journal of Computer Vision (IJCV)*, vol. 61, no. 3, pp. 211–231, 2005.

[81] H. Bay, T. Tuytelaars, and L. Van Gool, "Surf: Speeded up robust features," in *European Conference on Computer Vision (ECCV)*. Springer, 2006, pp. 404–417.

[82] D. G. Lowe, "Object recognition from local scale-invariant features," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, vol. 2. Ieee, 1999, pp. 1150–1157.

[83] S.-Y. Wang, O. Wang, R. Zhang, A. Owens, and A. A. Efros, "Cnn-generated images are surprisingly easy to spot... for now," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 8695–8704.

[84] H. Lu, R. Shen, and F.-L. Chung, "Fragile watermarking scheme for image authentication," *Electronics Letters*, vol. 39, no. 12, pp. 898–900, 2003.

[85] H. He, J. Zhang, and H.-M. Tai, "A wavelet-based fragile watermarking scheme for secure image authentication," in *International Workshop on Digital Watermarking*. Springer, 2006, pp. 422–432.

[86] X. Zhang, S. Wang, Z. Qian, and G. Feng, "Self-embedding watermark with flexible restoration quality," *Multimedia Tools and Applications*, vol. 54, no. 2, pp. 385–395, 2011.

[87] X. Zhang and S. Wang, "Fragile watermarking with error-free restoration capability," *IEEE Transactions on Multimedia*, vol. 10, no. 8, pp. 1490–1499, 2008.

[88] X. Zhang, "Fragile watermarking scheme using a hierarchical mechanism," *Signal processing*, vol. 89, no. 4, pp. 675–679, 2009.

[89] P. Korus and A. Dziech, "Efficient method for content reconstruction with self-embedding," *IEEE Transactions on Image Processing*, vol. 22, no. 3, pp. 1134–1147, 2012.

[90] X. Zhang, Z. Qian, and G. Feng, "Reference sharing mechanism for watermark self-embedding," *IEEE Transactions on Image Processing (TIP)*, vol. 20, no. 2, pp. 485–495, 2010.

[91] X. Zhang, Z. Qian, Y. Ren, and G. Feng, "Watermarking with flexible self-recovery quality based on compressive sensing and compositive reconstruction," *IEEE Transactions on Information Forensics and Security (TIFS)*, vol. 6, no. 4, pp. 1223–1232, 2011.

[92] L. Verdoliva, "Media forensics and deepfakes: an overview," *IEEE Journal of Selected Topics in Signal Processing*, vol. 14, no. 5, pp. 910–932, 2020.

[93] X. Bi, Y. Wei, B. Xiao, and W. Li, "Rru-net: The ringed residual u-net for image splicing forgery detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019, pp. 0–0.

[94] Y. Wu, W. AbdAlmageed, and P. Natarajan, "Mantra-net: Manipulation tracing network for detection and localization of image forgeries with anomalous features," in *Proceedings of the*

IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 9543–9552.

[95] F. Chen, H. He, and Y. Huo, "Self-embedding watermarking scheme against jpeg compression with superior imperceptibility," Multimedia Tools and Applications, vol. 76, no. 7, pp. 9681–9712, 2017.

[96] R. Preda and D. Vizireanu, "Watermarking-based image authentication robust to jpeg compression," Electronics Letters, vol. 51, no. 23, pp. 1873–1875, 2015.

[97] M.-J. Tsai and C.-C. Chien, "Authentication and recovery for wavelet-based semifragile watermarking," Optical Engineering, vol. 47, no. 6, p. 067005, 2008.

[98] L. Li, J. Bao, T. Zhang, H. Yang, D. Chen, F. Wen, and B. Guo, "Face x-ray for more general face forgery detection," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2020, pp. 5001–5010.

[99] A. R. Bruna, G. Messina, and S. Battiato, "Crop detection through blocking artefacts analysis," in International Conference on Image Analysis and Processing. Springer, 2011, pp. 650–659.

[100] M. Fanfani, M. Iuliani, F. Bellavia, C. Colombo, and A. Piva, "A vision-based fully automated approach to robust image cropping detection," Signal Processing: Image Communication, vol. 80, p. 115629, 2020.

[101] B. Van Hoorick and C. Vondrick, "Dissecting image crops," arXiv preprint arXiv:2011.11831, 2020.

[102] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," arXiv preprint arXiv:1406.2661, 2014.

[103] L. Dinh, D. Krueger, and Y. Bengio, "Nice: Non-linear independent components estimation," arXiv preprint arXiv:1410.8516, 2014.

[104] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using real nvp," arXiv preprint arXiv:1605.08803, 2016.

[105] D. P. Kingma and P. Dhariwal, "Glow: Generative flow with invertible 1x1 convolutions," Advances in Neural Information Processing Systems (NIPS), vol. 31, 2018.

[106] R. Zhao, T. Liu, J. Xiao, D. P. Lun, and K.-M. Lam, "Invertible image decolorization," IEEE Transactions on Image Processing (TIP), vol. 30, pp. 6081–6095, 2021.

[107] M. Xiao, S. Zheng, C. Liu, Y. Wang, D. He, G. Ke, J. Bian, Z. Lin, and T.-Y. Liu, "Invertible image rescaling," in European Conference on Computer Vision (ECCV). Springer, 2020, pp. 126–144.

[108] Y. Xing, Z. Qian, and Q. Chen, "Invertible image signal processing," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 6287–6296.

[109] S.-P. Lu, R. Wang, T. Zhong, and P. L. Rosin, "Large-capacity image steganography based on invertible neural networks," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 10 816–10 825.

[110] R. O. Park E, Liu W, "Ilsvrc-2017," 2017. [Online]. Available: http://www.image-net.org/challenges/LSVRC/2017

[111] T. Karras, T. Aila, S. Laine, and J. Lehtinen, "Progressive growing of gans for improved quality, stability, and variation," arXiv preprint arXiv:1710.10196, 2017.

[112] Q. Ying, X. Hu, X. Zhang, Z. Qian, S. Li, and X. Zhang, "Rwn: Robust watermarking network for image cropping localization," in 2022 IEEE International Conference on Image Processing (ICIP). IEEE, 2022, pp. 301–305.

[113] T. Miyato, T. Kataoka, M. Koyama, and Y. Yoshida, "Spectral normalization for generative adversarial networks," arXiv preprint arXiv:1802.05957, 2018.

[114] S. Iizuka, E. Simo-Serra, and H. Ishikawa, "Globally and locally consistent image completion," ACM Transactions on Graphics (ToG), vol. 36, no. 4, pp. 1–14, 2017.

[115] D. Ulyanov, A. Vedaldi, and V. Lempitsky, "Deep image prior," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2018, pp. 9446–9454.

[116] N. Otsu, "A threshold selection method from gray-level histograms," IEEE transactions on systems, man, and cybernetics, vol. 9, no. 1, pp. 62–66, 1979.

[117] J. Fridrich, T. Pevnỳ, and J. Kodovskỳ, "Statistically undetectable jpeg steganography: dead ends challenges, and opportunities," in Proceedings of the 9th workshop on Multimedia & security, 2007, pp. 3–14.

[118] A. Islam, C. Long, A. Basharat, and A. Hoogs, "Doa-gan: Dual-order attentive generative adversarial network for image copy-move forgery detection and localization," in Proceedings of the

[119] G. Schaefer and M. Stich, "Ucid: An uncompressed color image database," in Storage and Retrieval Methods and Applications for Multimedia 2004, vol. 5307. International Society for Optics and Photonics, 2003, pp. 472–480.

[120] E. Agustsson and R. Timofte, "Ntire 2017 challenge on single image super-resolution: Dataset and study," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), 2017, pp. 126–135.

[121] J. Dong, W. Wang, and T. Tan, "Casia image tampering detection evaluation database," in 2013 IEEE China Summit and International Conference on Signal and Information Processing. IEEE, 2013, pp. 422–426.

[122] G. Mahfoudi, B. Tajini, F. Retraint, F. Morain-Nicolier, J. L. Dugelay, and P. Marc, "Defacto: Image and face manipulation dataset," in 2019 27th European Signal Processing Conference (EUSIPCO). IEEE, 2019, pp. 1–5.

[123] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017, pp. 2961–2969.

[124] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," Advances in Neural Information Processing Systems (NIPS), vol. 28, pp. 91–99, 2015.

[125] Z. Liu, P. Luo, X. Wang, and X. Tang, "Large-scale celebfaces attributes (celeba) dataset," Retrieved August, vol. 15, no. 2018, p. 11, 2018.

[126] B. Zhou, A. Lapedriza, A. Khosla, A. Oliva, and A. Torralba, "Places: A 10 million image database for scene recognition," IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI), vol. 40, no. 6, pp. 1452–1464, 2017.

[127] W. Grathwohl, R. T. Chen, J. Bettencourt, I. Sutskever, and D. Duvenaud, "Ffjord: Free-form continuous dynamics for scalable reversible generative models," arXiv preprint arXiv:1810.01367, 2018.

[128] L. Ardizzone, J. Kruse, S. Wirkert, D. Rahner, E. W. Pellegrini, R. S. Klessen, L. Maier-Hein, C. Rother, and U. Köthe, "Analyzing inverse problems with invertible neural networks," arXiv preprint arXiv:1808.04730, 2018.

[129] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in Proceedings of the IEEE International Conference on Computer Vision (ICCV), 2017, pp. 2223–2232.

[130] T. Pevnỳ, T. Filler, and P. Bas, "Using high-dimensional image models to perform highly undetectable steganography," in International Workshop on Information Hiding. Springer, 2010, pp. 161–177.

[131] A. Cheddad, J. Condell, K. Curran, and P. Mc Kevitt, "Digital image steganography: Survey and analysis of current methods," Signal processing, vol. 90, no. 3, pp. 727–752, 2010.

[132] I. Cox, M. Miller, J. Bloom, J. Fridrich, and T. Kalker, Digital watermarking and steganography. Morgan kaufmann, 2007.

[133] J. Hayes and G. Danezis, "Generating steganographic images via adversarial training," arXiv preprint arXiv:1703.00371, 2017.

[134] W. Tang, S. Tan, B. Li, and J. Huang, "Automatic steganographic distortion learning using a generative adversarial network," IEEE Signal Processing Letters, vol. 24, no. 10, pp. 1547–1551, 2017.

[135] P. Wu, Y. Yang, and X. Li, "Stegnet: Mega image steganography capacity with deep convolutional network," Future Internet, vol. 10, no. 6, p. 54, Jun 2018. [Online]. Available: http://dx.doi.org/10.3390/fi10060054

[136] K. A. Zhang, A. Cuesta-Infante, and K. Veeramachaneni, "Steganogan: Pushing the limits of image steganography," 2019.

[137] E. Wengrowski and K. Dana, "Light field messaging with deep photographic steganography," in Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2019, pp. 1515–1524.

[138] C. W. Park, Y. H. Moon, and I. K. Eom, "Image tampering localization using demosaicing patterns and singular value based prediction residue," IEEE Access, vol. 9, pp. 91 921–91 933, 2021.

[139] N. Le and F. Retraint, "An improved algorithm for digital image authentication and forgery localization using demosaicing artifacts," IEEE Access, vol. 7, pp. 125 038–125 053, 2019.

[140] J.-H. Choi, H. Zhang, J.-H. Kim, C.-J. Hsieh, and J.-S. Lee, "Evaluating robustness of deep image super-resolution against

adversarial attacks," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 303–311.

[141] P. Neekhara, B. Dolhansky, J. Bitton, and C. C. Ferrer, "Adversarial threats to deepfake detection: A practical perspective," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 923–932.

[142] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4401–4410.

[143] Y. Wang, L. Wang, J. Yang, W. An, and Y. Guo, "Flickr1024: A large-scale dataset for stereo image super-resolution," in *The IEEE International Conference on Computer Vision Workshops (ICCVW)*, Oct 2019, pp. 3852–3857.

[144] E. Schonfeld, B. Schiele, and A. Khoreva, "A u-net based discriminator for generative adversarial networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 8207–8216.

[145] L. Wang, Y. Wang, Z. Liang, Z. Lin, J. Yang, W. An, and Y. Guo, "Learning parallax attention for stereo image super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, p. 12250–12259.

**Sheng Li** received the Ph.D. degree at the School of Electrical and Electronic Engineering, Nanyang Technological University, Singapore, in 2013. From 2013 to 2016, he was a research fellow in Rapid Rich Object Search (ROSE) Lab, Nanyang Technological University. He is currently an Associate Professor with the School of Computer Science, Fudan University, China. His research interests include biometric template protection, pattern recognition, multimedia forensics and security. He is the recipient of the IEEE WIFS Best Student Paper Silver Award.



**Qichao Ying** received the B.S. and M.S. degree from the School of Communication and Information Engineering, Shanghai University, China, respectively in 2017 and 2020. He is currently a doctoral candidate in the School of Computer Science, Fudan University, China. His research interests include multimedia forensics, data hiding and fake news detection.



**Hang Zhou** received his B.S. degree in 2015 from Shanghai University and a Ph.D. degree in 2020 from the University of Science and Technology of China. Currently, he is a postdoctoral researcher at Simon Fraser University. His research interests include computer graphics, multimedia security and deep learning.



**Xinpeng Zhang** received the B.S. degree in computational mathematics from Jilin University, China, in 1995, and the M.E. and Ph.D. degrees in communication and information system from Shanghai University, China, in 2001 and 2004, respectively, where he has been with the faculty of the School of Communication and Information Engineering, since 2004, and is currently a Professor. His research interests include information hiding, image processing, and digital forensics. He has published over 200 papers in these areas.



**Zhenxing Qian** received the B.S. and Ph.D. degrees from the University of Science and Technology of China (USTC), in 2003 and 2007, respectively. He is currently a Professor with the School of Computer Science, Fudan University. He has published over 100 peer-reviewed papers on international journals and conferences. His research interests include information hiding, image processing, and multimedia security.