# Visible Watermark Removal via Self-calibrated Localization and Background Refinement

Jing Liang
MoE Key Lab of Artificial Intelligence,
Shanghai Jiao Tong University
leungjing@sjtu.edu.cn

Li Niu*
MoE Key Lab of Artificial Intelligence,
Shanghai Jiao Tong University
ustcnewly@sjtu.edu.cn

Fengjun Guo
INTSIG
fengjun_guo@intsig.net

Teng Long
INTSIG
mike_long@intsig.net

Liqing Zhang
MoE Key Lab of Artificial Intelligence,
Shanghai Jiao Tong University
zhang-lq@cs.sjtu.edu.cn

## ABSTRACT

Superimposing visible watermarks on images provides a powerful weapon to cope with the copyright issue. Watermark removal techniques, which can strengthen the robustness of visible watermarks in an adversarial way, have attracted increasing research interest. Modern watermark removal methods perform watermark localization and background restoration simultaneously, which could be viewed as a multi-task learning problem. However, existing approaches suffer from incomplete detected watermark and degraded texture quality of restored background. Therefore, we design a two-stage multi-task network to address the above issues. The coarse stage consists of a watermark branch and a background branch, in which the watermark branch self-calibrates the roughly estimated mask and passes the calibrated mask to background branch to reconstruct the watermarked area. In the refinement stage, we integrate multi-level features to improve the texture quality of watermarked area. Extensive experiments on two datasets demonstrate the effectiveness of our proposed method.

## CCS CONCEPTS

• **Computing methodologies → Image processing**.

## KEYWORDS

watermark removal; multi-task learning; two-stage network

*Corresponding Author

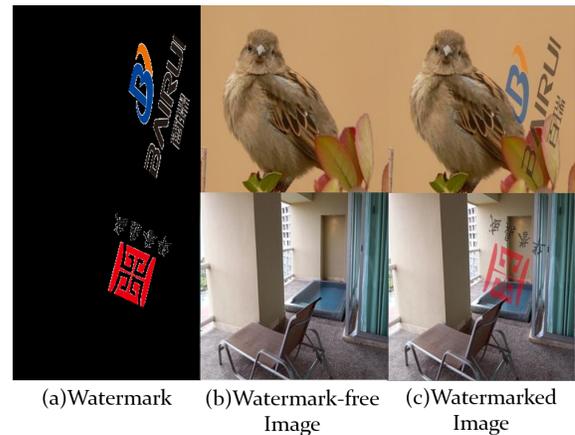| (a)Watermark | (b)Watermark-free Image | (c)Watermarked Image |

**Figure 1: The watermarked image (c) is acquired by superimposing a watermark (a) on the background image (b) via alpha blending. Given a watermarked image (c), watermark removal task aims to reconstruct the watermark-free image (b) without knowing the watermark mask.**

## 1 INTRODUCTION

With the surge of social media, images become the most prevailing carriers for recording and conveying information. To protect the copyright or claim the ownership, various types of visible watermarks are designed and overlaid on background images via alpha blending. Superimposing visible watermark is considered as an efficient and effective approach to combat against attackers. However, watermarked images are likely to be converted back to watermark-free images by virtue of modern watermark removal techniques. To evaluate and strengthen the robustness of visible watermarks in an adversarial way, watermark removal task has raised the research interest in recent years [2, 4, 6, 9, 19, 24].

Watermark removal, which aims to reconstruct background images based on watermarked images, is an open and challenging problem. Watermarks can be overlaid at any position of a background image with different sizes, shapes, colors, and transparencies. Besides, the watermarks often contain complex patterns like warped symbols, thin line, shadow effects, *etc.* The above reasons render the watermark removal task dramatically difficult when no

prior knowledge is provided. An example of watermark, watermark-free image, and watermarked image is shown in Figure 1. In the remainder of this paper, we use two terms "background image" and "watermark-free image" exchangably.

In some pioneer works, the location of watermarked area is required. Guided by watermark mask, watermark removal is similar to image inpainting [20] or feature matching problem [29, 31]. Nevertheless, manually annotating the watermark mask for each image is extremely time-consuming and cost-expensive. Noticing the fact that multiple images are often marked with the same watermark, watermark could be detected and removed in a more effective way [9, 15]. Unfortunately, the assumption in [9, 15] limits their application to real-world scenarios. Recently, researchers [2, 4, 6, 19, 24, 28] attempt to solve the blind watermark removal problem in an end-to-end manner with deep learning approaches. Some works [2, 24] formulated the watermark removal problem as an image-to-image translation task without localizing the watermark. On the contrary, other works realized that watermark should be localized and removed sequentially [4] or simultaneously [6, 19, 28]. Despite the great success of these emerging methods, they are still struggling to localize the watermark precisely and completely, especially when the watermark has complex patterns, diverse colors, or isolated fragments. The inaccurate watermark mask will interfere the reconstruction of background image. Moreover, the reconstructed images are suffering from quality issues like blur, artifacts, and distorted structures, which awaits further improvement.

In this paper, we propose a novel watermark removal network via **S**elf-calibrated **L**ocalization and **B**ackground **R**efinement (**SLBR**), which consists of a coarse stage and a refinement stage. In the coarse stage, we consider watermark localization and watermark removal as two tasks in a multi-task learning framework. Specifically, we employ a U-Net [33] structure, in which two tasks share the same encoder but have two separate decoders. The mask decoder branch predicts multi-scale watermark masks, which provides guidance for the background decoder branch via Mask-guided Background Enhancement (MBE) module to better reconstruct watermark-free images. Considering that the watermarks in various images are considerably different in many aspects, we design a Self-calibrated Mask Refinement (SMR) module, in which the watermark feature is propagated to the whole feature map to better handle image-specific watermark. In the refinement stage, we take the predicted watermark mask and watermark-free image in the coarse stage as input, to produce a refined watermark-free image. To fully exploit the useful information in the coarse stage, we add skip-stage connections between the background decoder branch in the coarse stage and the encoder in the refinement stage. Considering that different levels of features capture the structure information or texture details, we repeatedly use Cross-level Feature Fusion (CFF) modules to aggregate multi-level encoder features in the refinement stage. The output image from the refinement stage is the final recovered background image. Our main contributions could be summarized as follows,

- We propose a novel two-stage multi-task network named SLBR with cross-stage and cross-task information propagation for watermark removal task.

- In the coarse stage, we devise a novel Self-calibrated Mask Refinement (SMR) module to calibrate the watermark mask and a novel Mask-guided Background Enhancement (MBE) module to enhance the background representation.
- In the refinement stage, we propose a novel Cross-level Feature Fusion (CFF) module, which is repeatedly used to get the refined watermark-free image.
- Experiments on two datasets demonstrate the effectiveness of our proposed method.

## 2 RELATED WORKS

In this section, we first introduce a broad range of image content removal applications, and then describe the existing watermark removal methods. Besides, since our network involves multi-level feature fusion, we also briefly review the related methods.

**Image Content Removal:** Similar to watermark removal task, some existing tasks also focus on removing some undesirable content from an image, for example, deraining [14, 32, 36, 40], blind shadow removal [7, 11, 37], dehazing [1, 5, 12, 18, 39, 41, 43], and so on. However, these removed contents (*e.g.*, rain, shadow, haze) often consist of repeated patterns and monotonous colors. Different from the above tasks, watermark removal task targets at removing the watermarks which have diverse shapes and colors. Therefore, watermark removal task is a unique and challenging task.

**Visible Watermark Removal:** Visible watermark provides a powerful weapon for protecting the copyright. To evaluate and improve the robustness of visible watermarks, watermark removal techniques are proposed and gradually draw attentions from the security community. In the earlier explorations [20, 29, 31], they generally interacted with users to indicate the watermark locations for the following background recovery, which limits its practical usage. Since acquiring each of image location is ineffective, [9, 15] assumed that multiple images have the same watermark pattern, in which multiple images are processed simultaneously to remove the common watermark pattern. However, the assumption in [9, 15] is too stringent and unpractical, which weakens its potential in real-world applications.

The development of deep learning techniques have greatly advanced the watermark removal task. Some methods [2, 24] formulated the watermark removal as an image-to-image translation task. Other methods [6, 19, 28] performed watermark localization and removal tasks at the same time. In [6, 19, 28], watermark localization and watermark removal were wrapped up in a multi-task learning framework. Nevertheless, the above methods [2, 6, 19, 24, 28] are still struggling to achieve satisfactory performance in localizing watermark and restoring the watermark-free images.

**Multi-level Feature Fusion:** Multi-level feature fusion has been widely used in various computer vision tasks [7, 12, 21, 27, 45] for boosting network performance. Aggregation strategies could vary from task to task, but most of them fall into the following classical scopes: dense connection [42], top-down and/or bottom up feature integration [25, 27], feature concatenation [7, 45, 47], weighted element-wise summation [3, 44]. Although these methods are capable of merging multi-level features, how to propagate multi-level information properly and efficiently in watermark removal task is still unsolved. In watermark removal approaches [19, 28], Hertz et

al. [19] only considered the skip connection from encoder; Liu et al. [28] further passed the shallowest decoder feature from coarse stage to refinement stage. Nevertheless, these methods overlook the potential capacity of multi-level features integration. Thus, we propose to bridge the coarse stage and refinement stage by multi-level feature propagation, and further perform cross-level feature interweaving for better background reconstruction.

# 3 OUR METHOD

Given a watermarked image J which is obtained by superimposing a watermark on the background image I, the goal of watermark removal is recovering the watermark-free image I based on the watermarked image J. Because the watermark mask M is usually unknown, we need to perform two tasks simultaneously: watermark localization and watermark removal, which can be accommodated under a multi-task learning framework. As exemplified in Figure 2, our whole network is designed in a coarse-to-fine manner, which comprises of a coarse stage and a refinement stage. In the coarse stage, similar to previous multi-task learning methods [19, 28], we employ one shared encoder and two split decoders, in which two decoders account for localizing the watermark (mask decoder branch) and restoring the background image (background decoder branch) respectively. In the mask decoder branch, we design a Self-calibrated Mask Refinement (SMR) module to promote the quality of predicted watermark mask. To ease the information flow from the mask decoder branch to the background decoder branch, we employ a Mask-guided Background Enhancement (MBE) module to enhance the background decoder features. In the refinement stage, we build skip-stage connections between the decoder features in the coarse stage and the encoder features in the refinement stage to facilitate information propagation from coarse stage to refinement stage. To better recover the structure and texture of background image, we also devise a Cross-level Feature Fusion (CFF) module to aggregate multi-level encoder features iteratively in the refinement stage. Next, we will elaborate on the coarse stage in Section 3.1 and the refinement stage in Section 3.2.

## 3.1 Coarse Stage

In the coarse stage, we adopt the U-Net [33] architecture with skip links connecting encoder and decoder features as shown in Figure 2. Specifically, we employ the structure of encoder block and decoder block in [19]. Watermark localization and watermark removal are treated as two tasks, which share all five encoder blocks and the first decoder block. But they have three separate decoder blocks, which form the mask decoder branch and background decoder branch separately. In the mask decoder branch, it is equipped with our designed Self-calibrated Mask Refinement (SMR) module and assigned to indicate watermark position. Apart from the predicted mask from the last decoder block, we also predict side output masks based on the features in the other two decoder blocks. In the background decoder branch, it is composed of Mask-guided Background Enhancement (MBE) module and assigned to recover the corrupt background area overlaid with watermark. SMR and MBE block will be detailed next.

**Self-calibrated Mask Refinement (SMR) module:** When predicting the watermark mask, we observe that the predicted masks are often incomplete. One possible reason is that the watermarks in different images have diverse shapes, colors, patterns, and transparencies, so one global predictor can hardly localize all various types of watermarks. Thus, we consider calibrating the mask predictor according to the watermark characteristics in each image, to improve the quality of predicted watermark mask. By taking the last decoder block in the mask decoder branch as an example, as shown in Figure 3, we concatenate the features from previous decoder block and skip connection, followed by stacked residual blocks [28]. We denote that $X^m$ is the feature map used to predict the watermark mask $\hat{M}$. Following [19, 28], we use binary cross-entropy loss to enforce $\hat{M}$ to be close to the ground-truth watermark mask M:

$$\mathcal{L}_{mask} = -\sum_{i,j} \left( M_{i,j} \log \hat{M}_{i,j} + (1 - M_{i,j}) \log(1 - \hat{M}_{i,j}) \right), \quad (1)$$

where $M_{i,j}$ (*resp.*, $\hat{M}_{i,j}$) is the $(i, j)$-th entry in M(*resp.*, $\hat{M}$). We first apply this roughly estimated mask $\hat{M}$ to the feature map $X^m$ to pool the averaged feature vector $x^m$. Although the estimated mask $\hat{M}$ has missed detection and false alarms, watermarked pixels still dominate the estimated mask and thus the averaged feature vector can roughly represent the watermark characteristics. After obtaining the averaged watermark feature $x^m$, we tend to compare all pixel-level features in $X^m$ with $x^m$. Specifically, we first employ a $1 \times 1$ conv layer (*resp.*, fully-connected layer) to project $X^m$ (*resp.*, $x^m$) to $\tilde{X}^m$ (*resp.*, $\tilde{x}^m$). In the projected space, we expect that the averaged watermark feature is close to the watermarked pixels but far away from the unmasked pixels. Then, we spatially replicate $\tilde{x}^m$ to the same size as $\tilde{X}^m$, giving rise to $\bar{X}^m$. We concatenate $\tilde{X}^m$ and $\bar{X}^m$, followed by a $1 \times 1$ conv layer to predict a binary affinity map, in which 1 (*resp.*, 0) indicates that this pixel-level feature is similar (*resp.*, dissimilar) to the averaged watermark feature. Apparently, the ground-truth affinity map should be identical with the ground-truth watermark mask. Therefore, we can apply the same loss as Eqn. (1) to supervise the affinity map. By using $\hat{M}'$ to denote the predicted affinity map, the loss can be expressed as

$$\mathcal{L}'_{mask} = -\sum_{i,j} \left( M_{i,j} \log \hat{M}'_{i,j} + (1 - M_{i,j}) \log(1 - \hat{M}'_{i,j}) \right), \quad (2)$$

in which $\hat{M}'_{i,j}$ is the $(i, j)$-th entry in $\hat{M}'$. By comparing all pixel-level features with the averaged watermark feature, the predicted affinity map can identify some missed detection and erase some false alarms. Because $\hat{M}'$ is refined $\hat{M}$, we use $\hat{M}'$ as input for the background decoder branch and the refinement stage. We refer to the above module as Self-calibrated Mask Refinement (SMR) module and replace the original decoder blocks [19] in the mask decoder branch by our SMR modules.

**Mask-guided Background Enhancement (MBE) module:** In the coarse stage, watermark localization and watermark removal are two closely related tasks under a multi-task learning framework. According to [6, 19, 28], knowing the watermark area will offer strong guidance for the watermark removal task. In previous multi-task learning works [8, 16, 26, 48], myriads of strategies have been proposed to encourage the information sharing and propagation across different tasks. In our problem, we conjecture that mask localization would provide more benefit for watermark removal
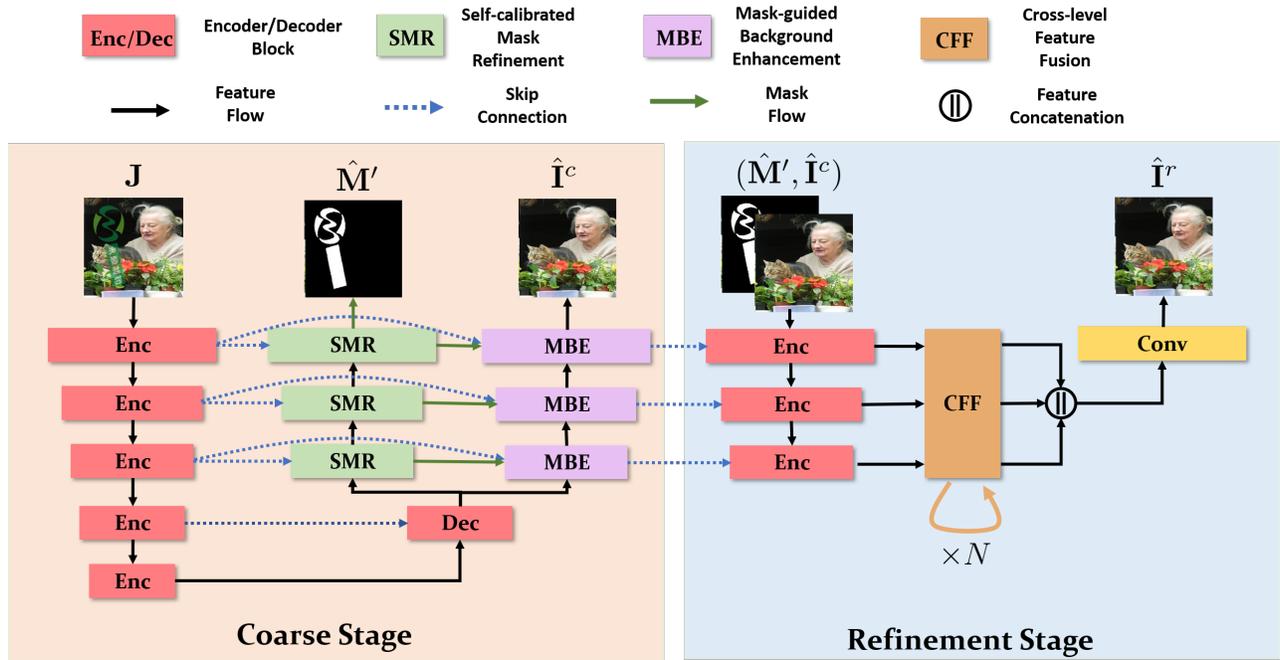
Figure 2: The illustration of our SLBR network which consists of a coarse stage and a refinement stage. The coarse stage contains one shared encoder and two separate decoder branches, which accounts for watermark localization and watermark-free image reconstruction respectively. The refinement stage takes the predicted watermark mask and watermark-free image from the coarse stage, producing the refined watermark-free image. We omit the side output masks in this figure for clarity.
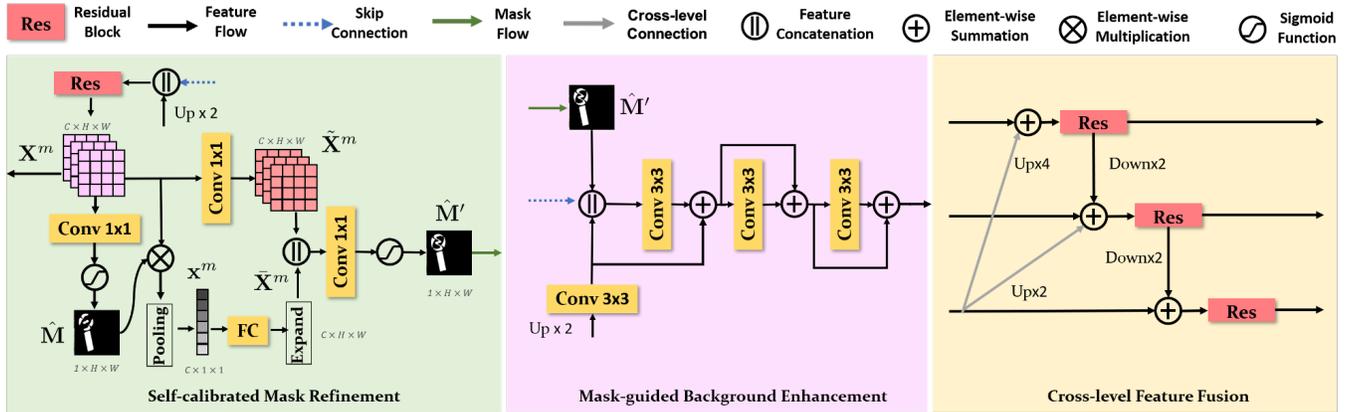


Figure 3: The illustration of our Self-calibrated Mask Refinement (SMR), Mask-guided Background Enhancement (MBE), Cross-level Feature Fusion (CFF) modules. "Pooling" means average pooling. "FC" means fully-connected layer. "Expand" means spatial replication.

than the other way around. Furthermore, our main goal is recovering the watermark-free image. Therefore, we design a Mask-guided Background Enhancement (MBE) module to guide the information flow from mask decoder branch to background decoder branch.

As shown in Figure 3, in each MBE module, we concatenate the output mask $\hat{M}'$ from the corresponding SMR module with the features from previous decoder block and skip connection. Then,

we apply a $3 \times 3$ conv layer to the concatenated feature to generate a feature residue, which is added back to the input feature. Following [19], We repeat this residual process for three times to produce the enhanced background decoder feature, which is fed into the next decoder block.

We notice that in some previous multi-task learning networks [17, 35], the features in one decoder branch are also appended to the

features in the other decoder branch. Different from them, our MBE module incorporates the predicted mask and learns residual information to boost the capacity of background representation. Here, we denote generated background image as $\hat{I}^c$, which is expected to be close to the ground-truth watermark-free image I using $L_1$ loss:

$$\mathcal{L}_{bg-L_1}^c = \|I - \hat{I}^c\|_1. \tag{3}$$

## 3.2 Refinement Stage

We observe that the restored watermark-free image $\hat{I}^c$ in the coarse stage may suffer from some quality issues like blur, artifact, and distorted structure, which calls for further improvement. Thus, we additionally attach a refinement stage to the coarse stage. We concatenate the coarse watermark-free image $\hat{I}^c$ and predicted watermark mask $\hat{M}'$ as the input for the refinement stage. First, we employ three encoder blocks [19] to extract multi-level features. To fully exploit the repaired content information in the coarse stage, we add skip-stage connections between the decoder features in the coarse stage and the encoder features in the refinement stage. Although WDNet [28] also uses the coarse-stage feature in the refinement stage, they simply append the last feature map in the coarse stage to the input for the refinement stage. Distinctively, we connect each background decoder feature in the coarse stage to its corresponding encoder feature with the same spatial size in the refinement stage in a symmetrical way, yielding enhanced multi-level encoder features in the refinement stage. Compared with [28], our skip-stage connections can integrate the content information in the coarse stage and refinement stage more thoroughly.

**Cross-level Feature Fusion (CFF) module:** Generally, we assume that the low-level encoder features with larger spatial size encode the texture details, while the high-level encoder features with smaller spatial size encode the structure information. To recover clear and coherent texture and structure for watermark-free image, we need to leverage multi-level encoder features in a better way. Thus, we design a Cross-level Feature Fusion (CFF) module, which is repeatedly used after the initial multi-level encoder features. As shown in Figure 3, in each CFF module, we upsample the high-level encoder feature to the same size of different low-level encoder features. After concatenating the upsampled high-level encoder feature with each low-level encoder feature, we also apply stacked residual blocks [28] to all encoder features including the high-level encoder feature. Besides this sparse connection fashion (*i.e.*, only propagating the high-level feature to the other levels of features), we have also tried dense connection fashion (*i.e.*, propagating all levels of features to the other levels of features) as in [45]. However, we observe that sparse connection is able to achieve comparable or even better results than dense connection. Thus, we adopt sparse connection in our CFF module for efficiency. We stack CFF module for $N$ times ($N = 3$ in our experiments).

Finally, based on the multi-level encoder features output from the last CFF module, we resize the encoder features of all levels to the target image size and aggregate them to obtain the final feature map. A $1 \times 1$ conv layer is applied to the final feature map to generate the refined watermark-free image $\hat{I}^r$. Similar to Eqn. (3), we employ $L_1$ loss to enforce the refined watermark-free image

to approach the ground-truth one:

$$\mathcal{L}_{bg-L_1}^r = \|I - \hat{I}^r\|_1. \tag{4}$$

To further ensure the quality of generated watermark-free image, we additionally employ perception loss [22, 46] based on VGG16 [34] pretrained on ImageNet [10]. The perception loss can be written as

$$\mathcal{L}_{bg-vgg} = \sum_{k \in 1,2,3} \|\Phi_{vgg}^k(\hat{I}^r) - \Phi_{vgg}^k(I)\|_1, \tag{5}$$

in which $\Phi_{vgg}^k(\cdot)$ means the activation map of $k$-th layer in VGG16.

Finally, we collect the losses in the coarse stage and the refinement stage, leading to the total loss:

$$\mathcal{L}_{all} = \mathcal{L}_{bg-L_1}^c + \mathcal{L}_{bg-L_1}^r + \lambda_{vgg}\mathcal{L}_{bg-vgg} + \\ \lambda_{mask}(\mathcal{L}_{mask} + \mathcal{L}_{mask}'), \tag{6}$$

in which $\lambda_{vgg}$ and $\lambda_{mask}$ are trade-off parameters. The whole network including the coarse stage and refinement stage can be trained in an end-to-end manner. In the testing stage, given a watermarked input image J, we use the output image $\hat{I}^r$ from the refinement stage as the final result.

## 4 EXPERIMENTS

In this section, we first introduce our used datasets, implementation details, and evaluation metrics. Then, we compare our SLBR method with existing watermark removal methods and image content removal methods. We also provide visualization results of all methods to demonstrate the effectiveness of our method. Moreover, we conduct comprehensive ablation studies to investigate the benefit of each stage and each module in our network.

## 4.1 Datasets and Implementation Details

Following [28], we conduct experiments on two large-scale benchmark datasets for watermark removal: Large-scale Visible Watermark Dataset (LVW) [4] and Colored Large-scale Watermark Dataset (CLWD) [28]. LVW mainly contains gray-scale watermarks, which have monotonous patterns and limited shapes. To overcome the shortcoming of LVW, the recent work [28] contributed a large-scale dataset CLWD with colored and diverse watermarks, which is more realistic and challenging than LVW dataset.

**LVW [4]:** LVW contains 48,000 watermarked images made of 64 gray-scale watermarks for training and 12,000 watermarked images made of 16 gray-scale watermarks for testing. The background images used in the training and test sets are randomly chosen from the train/val and test sets in PASCAL VOC2012 dataset [13] respectively.

**CLWD [28]:** CLWD contains 60,000 watermarked images made of 160 colored watermarks for training and 10,000 watermarked images made of 40 colored watermarks for testing. In CLWD, the watermarks are collected from open-sourced logo images websites. The original images used in training set and test sets are randomly chosen from PASCAL VOC2012 [13] training and test dataset respectively. When making watermarked image, the transparency is set in the range of (0.3, 0.7). Besides, the size, locations, rotation angle, and transparency of each watermark is randomly set in different images.

| Method | LVW | | | | CLWD | | | |
|---|---|---|---|---|---|---|---|---|
| | PSNR ↑ | SSIM ↑ | RMSE ↓ | RMSEw ↓ | PSNR ↑ | SSIM ↑ | RMSE ↓ | RMSEw ↓ |
| U-Net [33] | 30.33 | 0.9517 | 7.11 | 42.18 | 23.21 | 0.8567 | 19.35 | 48.43 |
| Qian *et al.* [32] | 39.92 | 0.9902 | 3.31 | 21.40 | 34.60 | 0.9694 | 5.40 | 19.34 |
| Cun *et al.* [7] | 40.68 | 0.9949 | 2.62 | 17.29 | 35.29 | 0.9712 | 5.28 | 18.25 |
| Li *et al.* [24] | 33.57 | 0.9690 | 5.84 | 34.71 | 27.96 | 0.9161 | 12.63 | 46.80 |
| Cao *et al.* [2] | 34.16 | 0.9714 | 5.51 | 33.42 | 29.04 | 0.9363 | 10.36 | 41.21 |
| WDNet [28] | 42.45 | 0.9954 | 2.39 | 12.75 | 35.53 | 0.9738 | 5.11 | 17.27 |
| BVMR [19] | 40.14 | 0.9910 | 3.24 | 18.57 | 35.89 | 0.9734 | 5.02 | 18.71 |
| SplitNet [6] | 43.16 | 0.9946 | 2.28 | 14.06 | 37.41 | 0.9787 | 4.23 | 15.25 |
| **SLBR (Ours)** | **43.48** | **0.9959** | **2.15** | **12.14** | **38.28** | **0.9814** | **3.76** | **14.07** |

Table 1: The results of different methods on LVW [4] and CLWD [28] datasets. The best results are denoted in boldface.
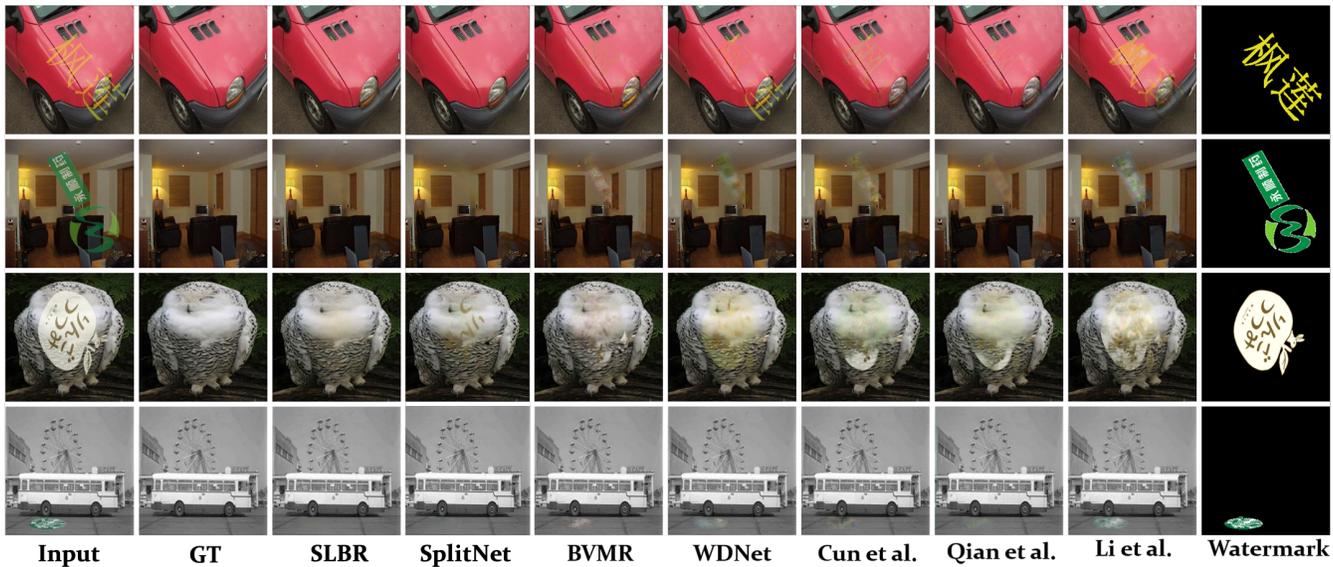


Figure 4: Visualization results of different methods on CLWD [28] dataset. Input is the watermarked image, GT is the ground-truth watermark-free image.

We implement our method using Pytorch [30]. We conduct all the experiments on the above two datasets. We set the input image size as $256 \times 256$. We choose Adam [23] optimizer with the initial learning rate 0.001, batch size 8, and momentum parameters $\beta_1 = 0.5, \beta_2 = 0.999$. The hyper-parameters $\lambda_{\text{vgg}}$ and $\lambda_{\text{mask}}$ in (6) are empirically set as 0.001 and 1 respectively, after a few trials by observing the quality of predicted masks and reconstructed images.

## 4.2 Baselines

To the best of our knowledge, there are only a few deep learning methods specifically designed for watermark removal: conditional GAN based watermark removal method Li *et al.* [24], self-attention model Cao *et al.* [2], blind visual motif removal method (BVMR) [19], split and refine network(SplitNet) [6],watermark-decomposition network (WDNet) [28]. We compare with these methods as the first group of baselines. Following [28], we also consider some image content removal methods and general image-to-image translation

methods as the second group of baselines. Concretely, we compare with attentive recurrent network Qian *et al.* [32] for deraining, attention-guided dual hierarchical aggregation network Cun *et al.* [7] for shadow removal, and U-Net [33] for general image-to-image translation.

## 4.3 Evaluation Metrics

Following [28], we adopt Peak Signal-to-Noise Radio (PSNR), Structural Similarity (SSIM) [38], Root-Mean-Square (RMSE) distance, weighted Root-Mean-Square distance (RMSE$_w$) as evaluation metrics. The difference between RMSE and RMSE$_w$ lies in that RMSE$_w$ is only computed within the watermarked area.

## 4.4 Experimental Results

The results of all methods on two datasets are summarized in Table 1. We reproduce the baseline results using their released code [7, 19, 28, 32, 33] or our own implementation [2, 24]. One

| # | SMR | MBE | CFF | Skip-stage | Evaluation Metrics | | | |
|---|-----|-----|-----|-----------|--------|--------|--------|---------|
| | | | | | PSNR ↑ | SSIM ↑ | RMSE ↓ | RMSEw ↓ |
| 1 | ∘ | ∘ | - | - | 35.99 | 0.9708 | 5.01 | 18.84 |
| 2 | ×1 | ∘ | - | - | 36.38 | 0.9740 | 4.87 | 17.43 |
| 3 | ×3 | ∘ | - | - | 36.50 | 0.9754 | 4.67 | 17.16 |
| 4 | ×3 | ×1 | - | - | 36.77 | 0.9759 | 4.53 | 16.92 |
| 5 | ×3 | ×3 | - | - | 36.90 | 0.9761 | 4.48 | 16.31 |
| 6 | ×3 | ×3 | ×0 | - | 37.19 | 0.9771 | 4.39 | 15.90 |
| 7 | ×3 | ×3 | ×1 | - | 37.27 | 0.9774 | 4.31 | 15.72 |
| 8 | ×3 | ×3 | ×2 | - | 37.35 | 0.9780 | 4.28 | 15.59 |
| 9 | ×3 | ×3 | ×3 | - | 37.42 | 0.9785 | 4.24 | 15.37 |
| 10 | ×3 | ×3 | ×3 | ×1 | 37.84 | 0.9797 | 3.94 | 14.59 |
| 11 | ×3 | ×3 | ×3 | ×2 | 38.02 | 0.9801 | 3.84 | 14.27 |
| 12 | ×3 | ×3 | ×3 | ×3 | **38.28** | **0.9814** | **3.76** | **14.07** |
| 13 | ×3 | ×3 | ∗ | ×3 | 37.65 | 0.9791 | 4.07 | 14.87 |

Table 2: Ablation studies of our method on CLWD [28] dataset. ∘ means using original decoder block. − means not using certain module or connection. ∗ means replacing CFF modules with original decoder blocks. ×N means the times of using certain module or connection. For ×1 (*resp.,* ×2), we replace the module or add the skip-stage connection in the shallowest one (*resp.,* two) layer(s). The best results are denoted in boldface.

may notice that our reported results are different from those reported in [28], especially the result of WDNet which is much worse than that in [28]. The performance degradation is attributed to a bug[1] in their released evaluation code. After fixing this bug, we re-evaluate and report the results of WDNet trained from scratch using their released code. One observation is that results on LVW dataset are much better than those on CLWD dataset, because LVW dataset only contains gray-scale watermarks and is much easier than CLWD dataset. Another observation is that the image content removal methods [7, 32] and watermark removal methods [6, 19, 28] based multi-task learning outperform image-to-image translation method [24] by a large margin, which verifies the effectiveness and necessity of predicting watermark mask. Moreover, baselines SplitNet [6], BVMR [19] and WDNet [28] specifically designed for watermark removal perform more favorably on two datasets than image content removal methods [7, 32].

Our SLBR method outperforms all baselines and achieves the best results on two datasets, which demonstrates the effectiveness of cross-task cross-stage information sharing and our devised modules. Our performance gain on LVW dataset [4] is not so obvious as that on CLWD dataset, which is again due to the simplicity of LVW dataset. In particular, gray-scale watermark removal task is much easier and we observe that the baseline methods can also capture the key pattern in LVW dataset within several training epochs. Therefore, the results on CLWD dataset can better justify the advantage of our proposed method.

For qualitative comparison, we show the visualization results of our method as well as baselines [6, 7, 19, 24, 28, 32] in Figure 4. In each row, from left to right, we show the input watermarked image, the ground-truth watermark-free image, the watermark-free images generated by different methods, and the watermark. It can be seen that our method can reconstruct the structure information

| Method | Evaluation Metrics | |
|--------|--------|---------|
| | F1 | IoU (%) |
| BVMR [19] | 0.7871 | 70.21 |
| WDNet [28] | 0.7240 | 61.20 |
| SplitNet [6] | 0.8027 | 71.96 |
| SLBR ($\hat{M}$) | 0.8107 | 73.10 |
| **SLBR** ($\hat{M}'$) | **0.8234** | **74.63** |

Table 3: Quantitative evaluation of watermark masks predicted by our method and baselines on CLWD [28] dataset.

and texture details of background more clearly and coherently, which shows the advantage of our proposed method for watermark removal task. For example, in the first row, baseline methods are capable of removing the main part of watermark, but there are still some remaining watermark, especially at the car light. In the second row, baseline methods suffer from color inconsistency and noticeable artifacts. In contrast, our method can generally erase the entire watermark and reconstruct the background image with clear texture.

## 4.5 Ablation Studies

In this section, we perform ablation studies to investigate the effectiveness of each module and each stage in our network. We start from a simple coarse stage network and gradually build up our full model. First, we only use the coarse stage and discard the refinement stage. Besides, we replace SMR and MBE modules with original decoder blocks [19] as mentioned in Section 3.2. In this case, we obtain a standard U-Net structure except two separate decoder branches for watermark localization and watermark removal respectively. The results of this simplest case are reported in row 1 in Table 2.

---

[1]They ignore the fact that return format of "imread" function in OpenCV is unsigned, which will raise a numeric overflow issue when subtracting images.
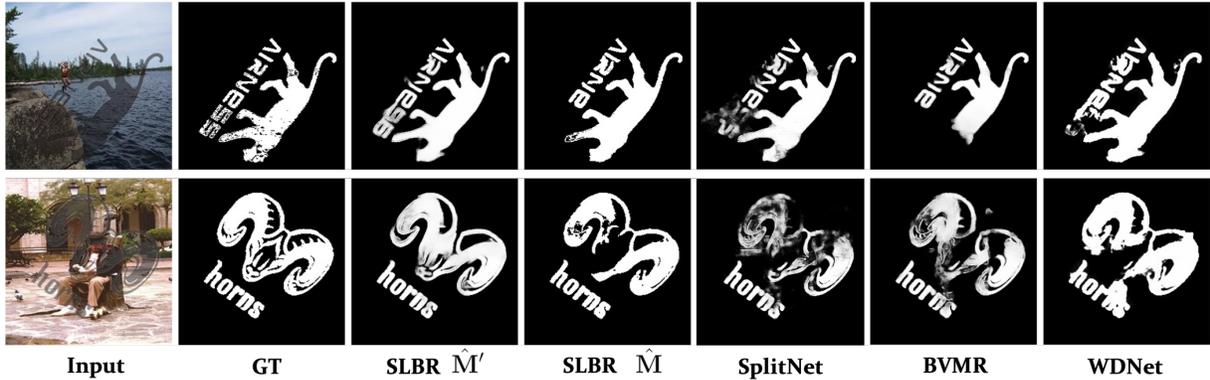
**Figure 5: Watermark localization results. From left to right, we show the input watermarked image, ground-truth watermark mask, the predicted results of our $\hat{M}'$, $\hat{M}$, and baselines.**

Then, we replace the last decoder block in the mask decoder branch with SMR, which corresponds to row 2 in Table 2. ×1 means that we only use one SMR module. Furthermore, we replace all the decoder blocks in the mask decoder branch with SMR, corresponding to row 3 in Table 2. By comparing the first three rows in Table 2, it is evident that our SMR is better than original decoder block and able to predict the watermark mask more accurately. Besides, using three SMR works better than only using one SMR, which implies that learning better side output masks can contribute to better intermediate decoder features.

Based on row 3, we replace the last decoder block in the background decoder branch with our MBE module, which import the output mask from mask decoder branch to enhance the last decoder feature, leading to the results in row 4. Furthermore, we replace all the decoder blocks in the background decoder branch with MBE, which utilizes all side output masks to enhance all the decoder features in the background decoder branch. The results using all three MBE modules are reported in row 5. By comparing row 3-5 in Table 2, we observe that our MBE is better than original decoder block and able to recover the background image better. Besides, using all three MBE performs better than only using one MBE, which implies that the side output masks from mask decoder branch can also benefit the reconstruction of watermark-free images.

Based on row 5, we introduce the refinement stage which only uses three encoder blocks and the final $1 \times 1$ conv layer without CFF block, resulting in row 6 in Table 2. Then, we gradually increase the number of CFF blocks, leading to row 7-9 in Table 2. By comparing row 6-9, we can draw a conclusion that using CFF to aggregate multi-level features is necessary and using more CFF leads to better results.

Based on row 9, we further bridge the coarse stage and the refinement stage by adding skip-stage connections. In row 10, we only connect the last decoder feature in the coarse stage and the first encoder feature in the refinement stage. In row 11-12, we link the last two (*resp.* three) decoder features and the first two (*resp.* three) encoder features using skip-stage connections, gradually yielding our full-fledged model. By comparing row 9-12, we can observe that the information propagation through skip-stage connection is

beneficial and more skip-stage connections can bring larger performance improvement. Finally, we replace CCF modules with decoder blocks [19], making the refinement network a U-Net structure. The results are listed in row 13, based on which our design of refinement stage performs more favorably than a U-Net network structure.

### 4.6 Watermark Localization

In this section, we evaluate the quality of our predicted watermark masks $\hat{M}$ and $\hat{M}'$. We also compare with SplitNet [6], BVMR [19], WDNet [28], which can also predict watermark mask as a byproduct. In terms of quantitative comparison, we calculate $F_1$ and IoU score based on the predicted mask and the ground-truth mask, where we simply use 0.5 as the threshold in all the experiments. The results are recorded in Table 3, which shows that $\hat{M}'$ indeed improves $\hat{M}$ and also outperforms the baselines [6, 19, 28] by a large margin.

We also show the predicted masks and ground-truth masks in Figure 5 for qualitative comparison. From Figure 5, we can see that $\hat{M}'$ is more complete and accurate. For example, in the first row, some texts in the rough estimation $\hat{M}$ are missing. Thanks to our SMR block, the final result $\hat{M}'$ is capable of predicting a complete mask, while other methods are struggling with the missed detection issue.

## 5 CONCLUSION

In this paper, we have studied watermark removal task and developed a two-stage multi-task network with novel MBE, SMR, and CCF modules, which can localize the watermark and recover the watermark-free image simultaneously. Extensive experiments on two datasets have verified the superiority of our proposed network.

# REFERENCES

[1] Bolun Cai, Xiangmin Xu, Kui Jia, Chunmei Qing, and Dacheng Tao. 2016. Dehazenet: An end-to-end system for single image haze removal. *IEEE Transactions on Image Processing* 25, 11 (2016), 5187–5198.

[2] Zhiyi Cao, Shaozhang Niu, Jiwei Zhang, and Xinyi Wang. 2019. Generative adversarial networks model for visible watermark removal. *IET Image Processing* 13, 10 (2019), 1783–1789.

[3] Dongdong Chen, Mingming He, Qingnan Fan, Jing Liao, Liheng Zhang, Dongdong Hou, Lu Yuan, and Gang Hua. 2019. Gated context aggregation network for image dehazing and deraining. In *IEEE winter conference on applications of computer vision*. 1375–1383.

[4] Danni Cheng, Xiang Li, Wei-Hong Li, Chan Lu, Fake Li, Hua Zhao, and Wei-Shi Zheng. 2018. Large-scale visible watermark detection and removal with deep convolutional networks. In *Chinese Conference on Pattern Recognition and Computer Vision*. 27–40.

[5] Xiaofeng Cong, Jie Gui, Kai-Chao Miao, Jun Zhang, Bing Wang, and Peng Chen. 2020. Discrete Haze Level Dehazing Network. In *Proceedings of the 28th ACM International Conference on Multimedia*. 1828–1836.

[6] Xiaodong Cun and Chi-Man Pun. 2020. Split then Refine: Stacked Attention-guided ResUNets for Blind Single Image Visible Watermark Removal. *arXiv preprint arXiv:2012.07007* (2020).

[7] Xiaodong Cun, Chi-Man Pun, and Cheng Shi. 2020. Towards ghost-free shadow removal via dual hierarchical aggregation network and shadow matting GAN. In *Proceedings of the AAAI Conference on Artificial Intelligence*. 10680–10687.

[8] Jifeng Dai, Kaiming He, and Jian Sun. 2016. Instance-aware semantic segmentation via multi-task network cascades. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3150–3158.

[9] Tali Dekel, Michael Rubinstein, Ce Liu, and William T Freeman. 2017. On the effectiveness of visible watermarks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2146–2154.

[10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. 248–255.

[11] Bin Ding, Chengjiang Long, Ling Zhang, and Chunxia Xiao. 2019. Argan: Attentive recurrent generative adversarial network for shadow detection and removal. In *Proceedings of the IEEE International Conference on Computer Vision*. 10213–10222.

[12] Hang Dong, Jinshan Pan, Lei Xiang, Zhe Hu, Xinyi Zhang, Fei Wang, and Ming-Hsuan Yang. 2020. Multi-scale boosted dehazing network with dense feature fusion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2157–2167.

[13] Mark Everingham, SM Ali Eslami, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2015. The pascal visual object classes challenge: A retrospective. *International journal of computer vision* 111, 1 (2015), 98–136.

[14] Zhiwen Fan, Huafeng Wu, Xueyang Fu, Yue Huang, and Xinghao Ding. 2018. Residual-guide network for single image deraining. In *Proceedings of the 26th ACM international conference on Multimedia*. 1751–1759.

[15] Yosef Gandelsman, Assaf Shocher, and Michal Irani. 2019. " Double-DIP": Unsupervised Image Decomposition via Coupled Deep-Image-Priors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 11026–11035.

[16] Yuan Gao, Jiayi Ma, Mingbo Zhao, Wei Liu, and Alan L Yuille. 2019. Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3205–3214.

[17] Zhangxuan Gu, Li Niu, Haohua Zhao, and Liqing Zhang. 2020. Hard Pixel Mining for Depth Privileged Semantic Segmentation. *IEEE Transactions on Multimedia* (2020).

[18] Kaiming He, Jian Sun, and Xiaoou Tang. 2010. Single image haze removal using dark channel prior. *IEEE transactions on pattern analysis and machine intelligence* 33, 12 (2010), 2341–2353.

[19] Amir Hertz, Sharon Fogel, Rana Hanocka, Raja Giryes, and Daniel Cohen-Or. 2019. Blind visual motif removal from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6858–6867.

[20] Chun-Hsiang Huang and Ja-Ling Wu. 2004. Attacking visible watermarking schemes. *IEEE transactions on multimedia* 6, 1 (2004), 16–30.

[21] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Baojin Huang, Yimin Luo, Jiayi Ma, and Junjun Jiang. 2020. Multi-scale progressive fusion network for single image deraining. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8346–8355.

[22] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. 2016. Perceptual losses for real-time style transfer and super-resolution. In *European conference on computer vision*. 694–711.

[23] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[24] Xiang Li, Chan Lu, Danni Cheng, Wei-Hong Li, Mei Cao, Bo Liu, Jiechao Ma, and Wei-Shi Zheng. 2019. Towards Photo-Realistic Visible Watermark Removal with Conditional Generative Adversarial Networks. In *International Conference on Image and Graphics*. 345–356.

[25] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. 2017. Feature pyramid networks for object detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2117–2125.

[26] Shikun Liu, Edward Johns, and Andrew J Davison. 2019. End-to-end multi-task learning with attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1871–1880.

[27] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. 2018. Path aggregation network for instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 8759–8768.

[28] Yang Liu, Zhen Zhu, and Xiang Bai. 2021. WDNet: Watermark-Decomposition Network for Visible Watermark Removal. In *Proceedings of the IEEE Winter Conference on Applications of Computer Vision*. 3685–3693.

[29] Jaesik Park, Yu-Wing Tai, and In So Kweon. 2012. Identigram/watermark removal using cross-channel correlation. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*. 446–453.

[30] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. 2019. Pytorch: An imperative style, high-performance deep learning library. *arXiv preprint arXiv:1912.01703* (2019).

[31] Soo-Chang Pei and Yi-Chong Zeng. 2006. A novel image recovery algorithm for visible watermarked images. *IEEE Transactions on Information Forensics and Security* 1, 4 (2006), 543–550.

[32] Rui Qian, Robby T Tan, Wenhan Yang, Jiajun Su, and Jiaying Liu. 2018. Attentive generative adversarial network for raindrop removal from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2482–2491.

[33] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. 2015. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*. 234–241.

[34] Karen Simonyan and Andrew Zisserman. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* (2014).

[35] Yi-Hsuan Tsai, Xiaohui Shen, Zhe Lin, Kalyan Sunkavalli, Xin Lu, and Ming-Hsuan Yang. 2017. Deep image harmonization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3789–3797.

[36] Cong Wang, Yutong Wu, Zhixun Su, and Junyang Chen. 2020. Joint self-attention and scale-aggregation for self-calibrated deraining network. In *Proceedings of the 28th ACM International Conference on Multimedia*. 2517–2525.

[37] Jifeng Wang, Xiang Li, and Jian Yang. 2018. Stacked conditional generative adversarial networks for jointly learning shadow detection and shadow removal. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1788–1797.

[38] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13, 4 (2004), 600–612.

[39] Dong Yang and Jian Sun. 2018. Proximal dehaze-net: A prior learning-based deep network for single image dehazing. In *Proceedings of the European Conference on Computer Vision*. 702–717.

[40] Youzhao Yang and Hong Lu. 2019. Single image deraining via recurrent hierarchy enhancement network. In *Proceedings of the 27th ACM International Conference on Multimedia*. 1814–1822.

[41] He Zhang and Vishal M Patel. 2018. Densely connected pyramid dehazing network. In *Proceedings of the IEEE Conference on Computer Vision and pattern recognition*. 3194–3203.

[42] He Zhang, Vishwanath Sindagi, and Vishal M Patel. 2018. Multi-scale single image dehazing using perceptual pyramid deep network. In *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*. 902–911.

[43] Jing Zhang, Yang Cao, Zheng-Jun Zha, and Dacheng Tao. 2020. Nighttime dehazing with a synthetic benchmark. In *Proceedings of the 28th ACM International Conference on Multimedia*. 2355–2363.

[44] Lu Zhang, Ju Dai, Huchuan Lu, You He, and Gang Wang. 2018. A bi-directional message passing model for salient object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1741–1750.

[45] Pingping Zhang, Dong Wang, Huchuan Lu, Hongyu Wang, and Xiang Ruan. 2017. Amulet: Aggregating multi-level convolutional features for salient object detection. In *Proceedings of the IEEE International Conference on Computer Vision*. 202–211.

[46] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. 2018. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 586–595.

[47] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. 2018. Residual dense network for image super-resolution. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2472–2481.

[48] Xiangyun Zhao, Haoxiang Li, Xiaohui Shen, Xiaodan Liang, and Ying Wu. 2018. A modulation module for multi-task learning with applications in image retrieval. In *Proceedings of the European Conference on Computer Vision*. 401–416.