# Split then Refine: Stacked Attention-guided ResUNets for Blind Single Image Visible Watermark Removal

**Xiaodong Cun and Chi-Man Pun***

University of Macau, Macau, China
yb87432@umac.mo, cmpun@umac.mo

## Abstract

Digital watermark is a commonly used technique to protect the copyright of medias. Simultaneously, to increase the robustness of watermark, attacking technique, such as watermark removal, also gets the attention from the community. Previous watermark removal methods require to gain the watermark location from users or train a multi-task network to recover the background indiscriminately. However, when jointly learning, the network performs better on watermark detection than recovering the texture. Inspired by this observation and to erase the visible watermarks blindly, we propose a novel two-stage framework with a stacked attention-guided ResUNets to simulate the process of detection, removal and refinement. In the first stage, we design a multi-task network called SplitNet. It learns the basis features for three sub tasks altogether while the task-specific features separately use multiple channel attentions. Then, with the predicted mask and coarser restored image, we design RefineNet to smooth the watermarked region with a mask-guided spatial attention. Besides network structure, the proposed algorithm also combines multiple perceptual losses for better quality both visually and numerically. We extensively evaluate our algorithm over four different datasets under various settings and the experiments show that our approach outperforms other state-of-the-art methods by a large margin. The code is available at: http://github.com/vinthony/deep-blind-watermark-removal.

## Introduction

Sharing rich contents such as images, audios or videos on social media has become a significant trend recently. Thus, digital watermarks, especially the visible ones, are used to credit the affiliation of the digital media. In general, the media vendor or the users assume the embedded watermarks are robust to various attacks, such as JPEG compression and the removal techniques. Also, the watermarks should have minimal visual influences on recognizing the original media contents. The robustness of the watermarking systems is a very important multimedia security issue for protecting the copyright or ownership of the digital media. Therefore, many research work have been studying the watermark removal methods to verify and improve the resilience of digital watermarks.
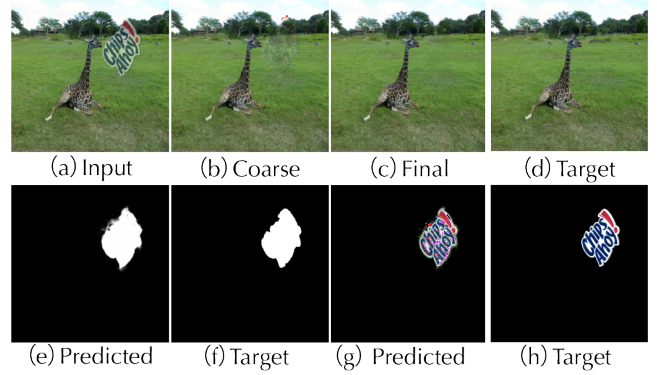


Figure 1: We propose an end-to-end, two-stage method to remove the visible watermark from a single image without any prior knowledge or user input. In the first stage, we use a multi-task attention-guided network to predict the coarser background (b), the location of watermark (e), and the recovered watermark (g). Then, the second-stage network uses (b) and (e) as input and learns the final results (c).

Since the invention of the watermark, removal techniques (Huang and Wu 2004; Park, Tai, and Kweon 2012; Pei and Zeng 2006) have been developed to attack them synchronously. In some previous work, detecting the location is a precondition for removing because it is relatively easier to match the features with the known mask. After getting the location, the watermark can be removed by image inpainting (Huang and Wu 2004) or features matching (Park, Tai, and Kweon 2012; Pei and Zeng 2006). On the other hand, since the single image-based watermark removal is difficult, multiple images with the same watermarks (Dekel et al. 2017; Gandelsman, Shocher, and Irani 2019) supply a strong prior knowledge for removing. However, erasing the watermark interactively and using multiple images restrict the application range of previous methods. Recently, Hertz et al. (2019) use a multi-task network to remove the watermark with a single image blindly for the first time. Nevertheless, as discussed in the previous work (Liu et al. 2018; Cun and Pun 2020; Ulyanov, Vedaldi, and Lempitsky 2018), the watermarked regions should be learnt explicitly using the carefully designed blocks. Thus, we propose a novel method

---

*Corresponding Author

to tackle these issues and follow it for further study.

The basic idea behind the proposed two-stage framework is: *Considering the framework of multi-task learning, watermark removal is more complicated than detection due to the texture non-harmony. Thus, further refinement is necessary.* As shown in Fig. 1, if the background (b), predicted mask (e) and watermark (g) are generated from a single network, the background (b) performs worse compared with other tasks. It might be because the watermark removal should restore the exactly pixel values from the degraded region while the watermark detection only need to gain the binary mask. Consequently, although a single multi-task network can be used for watermark removal as in Hertz et al. (2019), it is still necessary to smooth the watermarked region, and we refine it with another network using the predicted mask (e).

To model the analysis above, in this paper, we propose SplitNet and RefineNet, a novel attention-guided two-stage framework using Residual Block-based UNet (ResUNet) (Hertz et al. 2019) as the backbone of each stage. In SplitNet, we jointly learn three tasks (watermark removal, detection and recovery) using a single encoder and multiple decoders as in Hertz et al. (2019). Differently, we improve the performance of the original method with several modifications. Firstly, inspired by the recent works in multi-domain learning (Wang et al. 2019), we share the parameters in all three decoders while learning the bias for each task separately using channel attentions (Hu et al. 2018). In the second stage, we propose RefineNet to further refine the masked region pixels with the predicted mask and coarser results in SplitNet. Moreover, spatial-separated attention module (Cun and Pun 2020) is involved into RefineNet and learns the masked region specifically. Besides two stage framework, the simple pixel-wise $L_1$ loss in (Hertz et al. 2019) might also produce the blur results in the restored region. Thus, similar to the image super-resolution (Justin, Alexandre, and Li 2016) and reflection removal task (Zhang et al. 2018), we use the deep perceptual loss (Justin, Alexandre, and Li 2016) and SSIM loss (Wang et al. 2004) for better visual and numerical quality. The main contributions of the paper are as follows:

- We consider the task differences in multi-task watermark removal for the first time and formulate it as a two-stage framework by prediction and harmonization.

- In SplitNet, we regard the joint learning framework as a multi-domain learning problem and propose an accurate and compact model by domain (task)-specific attention.

- In RefineNet, we use the predicted mask and involve the mask-guided spatial attention modules for further harmonize the predicted regions.

- The results show that our approach outperforms various state-of-the-art methods by a large margin.

## Related Work

**Visible Watermark Removal** Digital watermarks play an important role in commercial digital copyright protection. Most previous watermark removal methods need to indicate the location of the watermark with the user before removing it by hand-crafted features, such as independent component analysis (Pei and Zeng 2006) and color space transformation (Park, Tai, and Kweon 2012). Besides, multiple images based visible watermark removal has also been widely investigated (Dekel et al. 2017; Gandelsman, Shocher, and Irani 2019). However, these methods require more prior knowledge, and removing the watermark by specific features only works on limited samples.

More recently, deep learning-based methods show great power in many computer vision tasks (He et al. 2016), which has also been used to remove visible watermarks (Li et al. 2019; Cheng et al. 2018; Hertz et al. 2019). However, a pre-trained watermark detector is necessary beforehand (Cheng et al. 2018), or they only consider watermark removal as an image-to-image translation problem (Li et al. 2019) as pix2pix (Isola et al. 2017). Besides, their dataset only contains gray-scale watermarks. Nevertheless, colorful images are more general nowadays. More recently, Hertz et al. (2019) design a novel deep learning-based method to detect, remove, and recover the motif in a single forward. However, their method still pays little attention to learning the degraded region specifically.

**Image In-painting** After identifying the specific watermark location by users, image in-painting techniques can also be used to remove the visible watermarks (Huang and Wu 2004). Thus, attention-guided in-painting methods, such as Partial Convolution (Liu et al. 2018), Contextual Attention (Yu et al. 2018) and Gated Convolution (Yu et al. 2019) may also serve to watermark removal. However, in-painting based watermark removal methods are supposed to point out the watermark location in the image by users. Moreover, in-painting based methods completely erase the masked region, and the background context will be re-generated by the neural network other than re-using the texture from the input. Differently, in the task of watermark removal, the transparent watermark contains both background and foreground information. Ignoring background cures will destroy the original structure, gaining undesired results. More importantly, due to the highly ill-pose of image in-painting, current methods only work on limited types of images, such as faces. Differently, our method can work in the wild with the help of transparent information.

**Image Restoration** Our task can also be considered as image restoration which only needs to restore certain regions. Thus, Blind shadow removal (Wang, Li, and Yang 2018; Ding et al. 2019; Cun, Pun, and Shi 2020), blind image harmonization (Cun and Pun 2020) and reflection separation (Zhang, Ng, and Chen 2018; Yang et al. 2018) share a similar goal with our task. However, the related tasks have limited resources compared with ours. For example, in blind shadow removal, the network only gets supervisions from the shadow-free image and the corresponding shadow mask, and it is impossible to get the ground truth shadow for learning jointly. Reflection removal/separation can generate paired reflection and transmission images from the dataset while it is hard to get the reflected region (mask) as supervision. Luckily, more information can be easily collected and used in our task, such as watermark location and content.
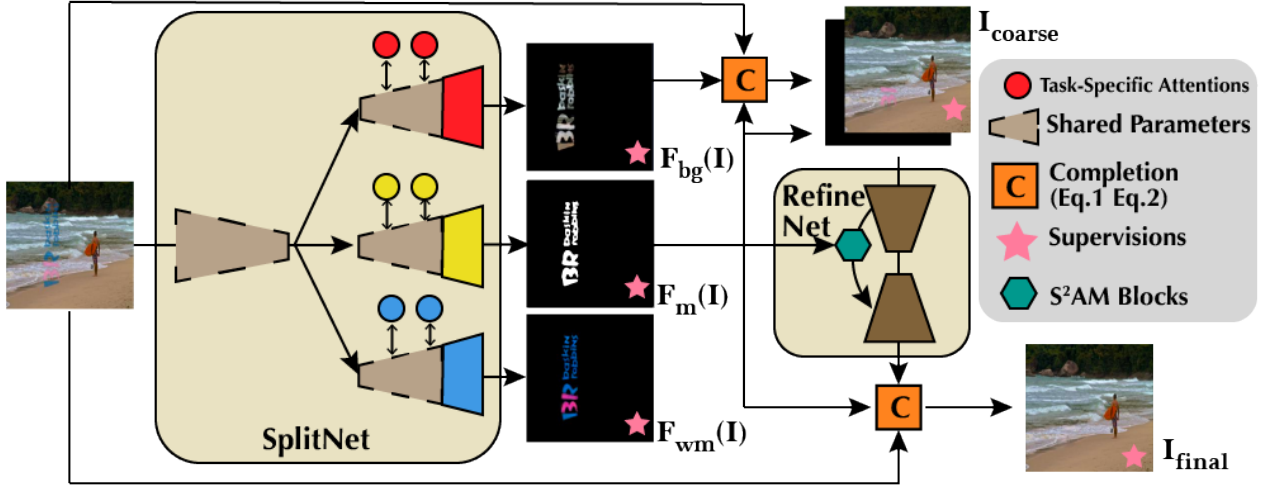
Figure 2: Proposed two-stage framework. We propose SplitNet to gain the coarser results by learning the watermark detection, removal and recovery jointly. Then, we propose RefineNet to smooth the learned region with the predicted mask and the recovered background from the previous stage. Thus, our network can be trained in an end-to-end fashion without any manual intervention. Note that, for clearness, we do not show any skip-connections between all the encoders and decoders.

## Proposed Method

We regard the blind single image-based visible watermark removal as a two-stage task. As shown in Fig. 2, in the first stage, given a single watermarked image $I$, we propose SplitNet $F$, a multi-task ResUNet inspired by multi-domain learning, to generate the coarser intermediate results: the recovered background image $F_{bg}(I)$, the location (mask) of watermark $F_m(I)$ and the recovered watermark $F_{wm}(I)$. Thus, the coarser restored image $I_{coarse}$ can be written as:

$$I_{coarse} = F_{bg}(I) \times F_m(I) + I \times (1 - F_m(I)) \quad (1)$$

As discussed previously, due to the difficulty of the tasks are different, further refinement is necessary for watermark removal. Thus, we propose RefineNet $R$ as the second stage, which uses $I_{coarse}$ and $F_m(I)$ to generate the final result $I_{final}$ and this network smooths the predicted watermarked region using a spatial attention mechanism. Finally, the refined result $I_{final}$ can be formulated by the predicted mask $F_m(I)$ and the original input $I$:

$$I_{final} = R(I_{coarse}, F_m(I)) \times F_m(I) + I \times (1 - F_m(I)) \quad (2)$$

Note that, although the proposed method is a cascaded framework, the inputs of the second network are totally generated by the output of the first stage. Thus, our network can be trained and evaluated in an end-to-end fashion without any manual intervention. Below, we give the details of the proposed SplitNet, RefineNet and the loss functions.

### SplitNet

It has been widely testified in various computer vision tasks (Liu, Johns, and Davison 2019; Ruder 2017; Wang and Forsyth 2009) that learning multiple related targets jointly can boost the performance of a single task. Thus, we propose SplitNet, a multi-task network jointly learn the watermark,

background and mask (as shown in Fig. 2) in the first stage. Similar to (Hertz et al. 2019), our network uses the residual block based UNet (Isola et al. 2017; Hertz et al. 2019; He et al. 2016) (ResUNet) as the basic encoder-decoder structure. Specifically, we build a shared encoder and multiple decoders as the main structure and each encoder/decoder contains 5 different scales of stacked (3 in each scale) ResBlocks to capture the multi-scale features. Getting benefits from the locally skip-connections in ResBlocks and globally skip-connections between the encoder and decoder, this backbone shows superior performance.

Different from previous work, we consider the joint learning framework as a multi-domain learning problem (Xiao, Gu, and Zhang 2020; Wang et al. 2019; Liu, Johns, and Davison 2019) in this task for the first time. In multi-domain learning, to learn an efficient model, almost all the parameters are shared in the training process while each domain needs to be emphasized using different parameters. Similarly, the three tasks in our framework focus on learning *one spatial region* and each task has to learn its specific features for individual reconstruction. In detail, we share the learn-able parameters in all three high-level decoders *altogether* and learn the specific feature for each task by domain attentions *individually*. This strategy helps us build a more efficient and effective model. To model this observation, inspired by SE-Net (Hu et al. 2018), we design the task-specific attentions to re-weight the importance of the channels for each decoder (task). Fig. 3 gives a close look at the proposed framework in the decoder over multiple shared branches. In each level of the decoders, we learn the basis features for all three tasks using the ResBlocks. After that, these basis features are re-weighted by the task-specific attentions. The detailed structure of the task-specific attention is also illustrated in the Fig. 3.

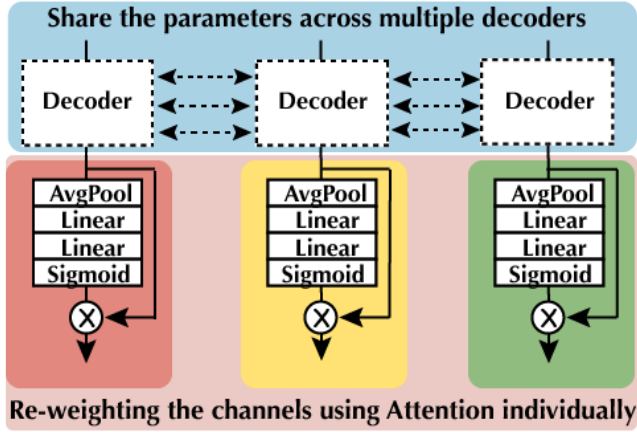Moreover, we optimize the structure of the original Res-

Figure 3: The detailed task-specific decoders in SplitNet. Our method shares the parameters among different tasks in the main stream and learns the features for each task individually using task-specific attentions.

Block in Hertz et al. (2019) for stable training. Since it is not the main contribution of our paper, we give a more detailed explanation in the supplementary material.

## RefineNet

The quality of the restored region remains poor if we predict the watermark-free image and mask from a single multi-task network (as shown in Fig. 1). However, the predicted mask achieves satisfying result, which might be because the texture recovery is far more difficult than detection. Therefore, we propose RefineNet for further refinement using the predicted coarse background and the location of watermark (mask) from SplitNet.

As discussed in previous works (Liu et al. 2018; Ulyanov, Vedaldi, and Lempitsky 2018), if we directly feed the mask and coarser results to UNet, the naive network might not focus on learning the masked region. Thus, inspired by the recently proposed Spatial-Separated Attention Module (S²AM (Cun and Pun 2020)) for image harmonization, we refine the predicted masked regions through the attention-guided network. Nevertheless, as a part of input, taking the mask from user is necessary in the original S²AM. Hence, we regard the predicted mask from SplitNet as a trustworthy label, feeding it into the RefineNet as both additional input channel and the mask of S²AM block (as shown in Fig. 2). Moreover, since the ResBlock based network has been proven to capture the low-level feature well (Hertz et al. 2019), we replace the original eight layers UNet in S²AM with the proposed five layers ResBlock-based backbone as SplitNet. Then, the S²AM is inserted in the two coarser levels of the decoder. Below, we give a short review of S²AM and details of the overall network structure.

S²AM is initially proposed in image harmonization to match the low-level features between the masked foreground and non-masked background. As shown in the Fig. 4, the S²AM block is inserted after the concatenation between the features from encoder and decoder. Then, with the
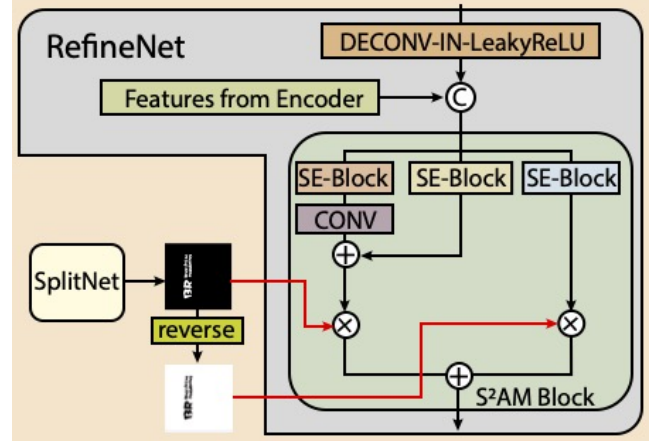


Figure 4: The predicted mask from the SplitNet is used in the S²AM Block to learn the differences in the masked and non-masked regions respectively.

help of the hard-coded mask, three channel attentions (SE-Block (Hu et al. 2018)) are applied to re-weight the original features automatically. More details of S²AM can be found in (Cun and Pun 2020).

## Loss Functions

The proposed framework is supposed to optimize multiple targets at the same time. The predicted results include three outputs from SplitNet: the predicted background $I_{coarse}$, the mask $M'$ of watermark location, and the reconstructed watermark $I_{wm}$. Also, the final prediction $I_{final}$ comes from RefineNet. Below, we introduce the details in each part of the loss functions.

**Losses for Watermark Removal**  Optimizing the losses of watermark removal is the main target of our framework, which contains multiple parts. Firstly, following the previous work (Hertz et al. 2019), we use the relative L1 loss over the ground truth mask $M$ by $\ell_r^{gt} = \frac{1}{sum(M)}||M \times F_{bg}(I) - M \times I_{gt}||_1$ to learn the masked region particularly, where $F_{bg}(I)$ is the direct output for background prediction. As suggested by Hertz et al. (2019), calculating the loss on $M$ only helps to prevent the early over-fitting. Moreover, we add an additional loss over the predicted mask $\ell_r^{pred} = \frac{1}{sum(M)}||M' \times F_{bg}(I) - M \times I_{gt}||_1$ to squeeze the gap between the prediction and ground truth. Because when testing, we need to recover the image using the predicted mask (as in Eq. 1, Eq. 2).

Besides, we also measure the quality of the recovered images $I_{coarse}$ and $I_{final}$. Computing the losses over the full image enables us to interpolate the frequently used perception losses into our system for better visual quality, such as deep perception loss $\Phi_{vgg}$ (Zhang et al. 2018; Justin, Alexandre, and Li 2016) and SSIM loss $\Psi_{ssim}$ (Zhao et al. 2016; Wang et al. 2004). In deep perceptual loss (Zhang et al. 2018), we extract the features between target and the predicted image from the pre-trained VGG16 (Simonyan and Zisserman 2014) in the layers of CONV1_2,

CONV2_2,CONV3_3, and then measure the $L_1$ distance in the feature domain. Finally, the total loss in the coarser stage $\ell_{coarse}$ and the final stage $\ell_{final}$ can be written as:

$$\ell_x = \alpha \sum_{k \in 1,2,3} ||\Phi_{vgg}^{k\_k}(I_x) - \Phi_{vgg}^{k\_k}(I_{gt})||_1 + \ell_r^{gt} \qquad (3)$$
$$+ \beta\Phi_{ssim}(I_x, I_{gt})) + ||I_x - I_{gt}||_1 + \ell_r^{pred},$$

where the recovered images $I_x(x \in \{coarse, final\})$ are from SplitNet and RefineNet by Eq. 1 and Eq. 2. We set all the $\alpha = 0.025$ and $\beta = 0.15$ in $\ell_{coarse}$ and $\ell_{refine}$.

**Losses for Watermark Detection**  Similar to salient object and shadow detection, detecting the watermarked region is a binary pixel-level segmentation. Thus, we choose Binary Cross Entropy as the loss function, which can be written as:

$$\ell_m = M \log(M') + (1 - M) \log(1 - M'), \qquad (4)$$

where $M$ and $M'$ are the ground truth mask and the predicted mask, respectively.

**Losses for Watermark Recovery**  Because the target watermark is also available in the proposed dataset, we measure the restored watermark with the original one over the masked region as $\ell_{wm}$. Following the relative L1 losses $\ell_r^{pred}$ and $\ell_r^{gt}$ for watermark removal, $\ell_{wm}$ is defined as: $\ell_{wm} = \ell_r^{pred''} + \ell_r^{gt''}$, where $\ell_r^{pred''}$ and $\ell_r^{gt''}$ represent the relative L1 loss on the restored watermark $F_{wm}(I)$.

Overall, the total loss $\ell_{all}$ in our algorithm is a combination of the losses above:

$$\ell_{all} = \ell_{coarse} + \ell_{refine} + \ell_{wm} + \ell_m. \qquad (5)$$

## Dataset Acquisition

Because there is no public available dataset containing all necessary information in our framework, we synthesize multiple datasets for different purposes to evaluate the proposed method[1] and we argue that these datasets may be sufficient enough for real-world cases since these watermarks are also man-made by software.

In detail, we choose the background (host) images from VAL2014 subset of the MSCOCO (Lin et al. 2014) dataset, which is similar to the previous work on visual motif removal (Hertz et al. 2019) and image harmonization (Cun and Pun 2020). Differently, to simulate the real-world watermarks, we collect over 1k different and famous logos from the Internet. Then, the watermarked samples are generated by natural images and the watermark (logo) in different locations, semi-transparency and sizes randomly. Each training/testing sample contains the synthesized watermarked image, the original background, the watermark and the mask of the watermark for supervisions. All the watermarks and the background images are non-overlapping in training and validation partition, showing the advantage of the algorithms in the unseen samples. We also analyze the distribution of the size, transparency, and class in the watermark to avoid the dataset bias. More details can be found in the supplementary material. Below, we give the detail settings of the four synthesized datasets:

---

[1]We will open source all the synthesized datasets for further research and comparison.

**LOGO-L**  We synthesize over 12K training and 2K test samples using 40% of the images and logos, respectively. In the dataset, the transparency of the watermarks ranges from 35% to 60%. Also, the watermarks are resized to 35% to 60% of the width (or height) of the host images. Thus, This dataset is an easier one because the watermark size is relatively small and the background is easily recognized through the watermarked area.

**LOGO-H**  We create a harder sub-dataset in LOGO-H containing the same quantity of samples as LOGO-L. The watermark size in this dataset accounts for 60% to 85% of the host image. Besides, we also randomly set the transparency from 60% to 85%. Thus, the watermark in this dataset is difficult to be removed due to the missing texture and the larger degraded regions.

**LOGO-Gray**  In real cases, the embedded watermark is usually a gray-scale image. Therefore, we create a LOGO-Gray sub-dataset to evaluate the performance which only contains gray-scale watermarks. This dataset also includes 12K images for training and 2K images for testing, as in LOGO-L and LOGO-H. The transparency and the size of the watermark are randomly chosen from 35% to 85%.

**LOGO30K**  We create a larger dataset by synthesizing the watermarks with various size, location, and transparency. In LOGO30K, up to 28k and 4k images will be trained and tested respectively. The watermark size and the transparency range from 35% to 85%.

## Experiments

**Implementation details**  We use PyTorch (Paszke et al. 2019) with CUDA v10.0 to implement our algorithm. To simplify the task and compare fairly, we conduct all the experiments on the image with a resolution of $256 \times 256$ and run 100 epochs for converging. The training batch size equals to 4, and all the optimizers are Adam (Kingma and Ba 2014) with the learning rate of $1e^{-3}$. Following previous works (Hertz et al. 2019; Cun and Pun 2020), we evaluate our approach with other state-of-the-art methods on the proposed four datasets using several popular numerical criteria, such as Structural Similarity (SSIM (Wang et al. 2004)), Peak Signal-to-Noise Radio (PSNR) and the deep perceptual similarity (LPIPS(Zhang et al. 2018)).

**Comparisons with state-of-the-art methods**  Nowadays, few visible watermark removal methods are based on deep learning, expect for the most related blind visual motif removal method BVMR (Hertz et al. 2019) and the deep watermark removal methods (Cheng et al. 2018; Li et al. 2019) using the UNet-like (Isola et al. 2017) structures in image-to-image translation. Thus, we also compare our algorithm with some learning-based methods on related tasks. For example, blind image harmonization on a soft-masked version of the spatial-separated attention model (BS$^2$AM(Cun and Pun 2020)); the learning-based reflection separation method (SIRF (Zhang et al. 2018)), which uses context aggregation network and perceptual loss; attention-guided dual hierarchical aggregation network for blind shadow removal (DHAN (Cun, Pun, and Shi 2020)). As shown in

Table 1: Comparisons between the proposed method and other state-of-the-art methods on all synthesized datasets.

| Metrics | LOGO-H | | | LOGO-L | | | LOGO-Gray | | | LOGO-30K | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS | PSNR | SSIM | LPIPS |
| Original | 28.33 | 0.9579 | 0.0748 | 35.12 | 0.9830 | 0.0312 | 31.82 | 0.9730 | 0.0411 | 31.01 | 0.9698 | 0.0555 |
| UNet | 30.51 | 0.9612 | 0.0544 | 34.87 | 0.9814 | 0.0297 | 32.15 | 0.9728 | 0.0353 | 33.27 | 0.9744 | 0.0365 |
| SIRF | 32.35 | 0.9673 | 0.0801 | 36.25 | 0.9825 | 0.0655 | 34.33 | 0.9782 | 0.0672 | 34.63 | 0.9783 | 0.0718 |
| BS$^2$AM | 31.93 | 0.9677 | 0.0445 | 36.11 | 0.9839 | 0.0223 | 32.91 | 0.9754 | 0.0305 | 34.88 | 0.9793 | 0.0283 |
| DHAN | 35.68 | 0.9809 | 0.0661 | 38.54 | 0.9887 | 0.0591 | 36.39 | 0.9836 | 0.0594 | 37.67 | 0.9860 | 0.0624 |
| BVMR | 36.51 | 0.9799 | 0.0237 | 40.24 | 0.9895 | 0.0126 | 38.90 | 0.9873 | 0.0115 | 38.28 | 0.9852 | 0.0181 |
| Ours | **40.05** | **0.9897** | **0.0115** | **42.53** | **0.9924** | **0.0087** | **42.01** | **0.9928** | **0.0073** | **41.27** | **0.9910** | **0.0114** |



(a) Input    (b) UNet    (c) BS$^2$AM    (d) SIRF    (e) GFSR    (f) BVMR    (g) Ours    (h) Target

Figure 5: Comparison with state-of-the-art methods on the proposed four synthesized datasets. These four images are taken from LOGO-30K, LOGO-Gray, LOGO-L and LOGO-H from top to bottom.

Tab. 1, our algorithm achieves much better results in both shallow perceptual metrics (PSNR,SSIM) and deep perceptual criterion (LPIPS). The significant improvement on the harder sub-dataset LOGO-H verifies the assumption that a single model cannot work perfectly in the masked region removal. Besides numeric metrics, our method also shows better visual quality than others. As shown in Fig. 5, naive UNet-based methods (UNet, BS$^2$AM) fail when the watermark is large and the transparency is low. Both SIRF and GFSR can only remove certain parts of the watermark because their methods use limited supervisions. For instance, SIRF separates the layers from the watermark and background, and GFSR learns to remove and detect jointly. Although BVMR outperforms other previous methods, it also shows noticeable artifacts due to the detection errors and texture misunderstanding. Differently, the proposed method uses all the information as supervisions and shows much better results. Interestingly, in some samples (such as the top

one in Fig. 5), our method cannot completely recover the color and texture of the background, while the images look more natural than other methods. This phenomenon also indicates the advantage of the proposed two-stage network with texture harmony in the masked region. We show more comparisons on the synthesized datasets and on "real world" samples (Dekel et al. 2017) in the supplementary material.

**Ablation Studies** As shown in Tab. 2, we evaluate the necessity of each component in our framework by removing or replacing it with other alternatives. We start the ablation study using the basic structure of ResUNet (Hertz et al. 2019) as our SplitNet (ResB). Then, we build the proposed framework by adding the naive ResUNet as RefineNet (ResB). However, the naive RefineNet cannot improve the performance, especially on SSIM. It might be because the refinement is essential in the masked region. Thus, when we add the S$^2$AM module into the ResUNet(ResS$^2$AM), the RefineNet gains the knowledge

Table 2: Ablation study on the LOGO-Gray dataset.

| SplitNet | RefineNet | additional loss | PSNR | SSIM |
|----------|-----------|-----------------|------|------|
| ResB | | | 38.90 | 0.9873 |
| ResB | ResB | | 39.34 | 0.9872 |
| ResB | ResS$^2$AM | | 40.05 | 0.9895 |
| iResB | iResS$^2$AM | | 41.77 | 0.9924 |
| iResBat | iResS$^2$AM | | 41.87 | 0.9925 |
| iResBat | iResS$^2$AM | $\ell_{ssim}$ | 41.89 | 0.9927 |
| iResBat | iResS$^2$AM | $\ell_{vgg} + \ell_{ssim}$ | **42.01** | **0.9928** |



(a) Input  (b) L$_1$ Loss  (c) w/ L$_{SSIM}$  (d) w/ L$_{vgg}$  (e) Full

Figure 6: The visual quality influence of loss functions in the final results. (b) means there is only relatively $\ell_{l1}$ losses for watermark removal, (c) and (d) refer to the additional $\ell_{ssim}$ and $\ell_{vgg}$, respectively. (e) is our full losses in Eq 3.

from the predicted mask and learns to recover the masked region specifically. Besides, the proposed improved Res-Block (iResB) and task-specific attentions (iResBat) in the SplitNet also achieve significant progress. Apart from network structure, the additional perceptual losses also play a critical role in recovering the watermarked region details as shown in Fig. 6, and we also give some numerical comparisons in Table 2 to support our claims on $\ell_{vgg}$ and $\ell_{ssim}$.

**Model Size**  The proposed framework contains two sub-networks, which naturally has more learning-able parameters than our baseline BVMR (Hertz et al. 2019). Consequently, the performance comparison on the fair model sizes is also necessary. As shown in Table 3, we modify the original BVMR (channel equals to 32 and depth equals to 5) to a larger network. For a fair comparison, all the networks are trained under the same loss function. Interestingly, the results of the deeper or heavier network structures in BVMR performs worse than the original one. It is not surprising that fewer parameters in BVMR lead to better results. On the one hand, it might be because the hyper-parameters of the network are the relatively best choice by network architecture searching. On the other hand, it might due to the unstable training in their network (see more discussions in supplementary material). Differently, our network adds relatively few parameters and gains much better performance.

**Intermediate Results**  In Fig. 1, we have shown all the intermediate results in our framework for a better understanding of the proposed pipeline. It is clear that the coarse output from the SplitNet contains noticeable artifacts because SplitNet only includes a naive traditional convolution net-

Table 3: Comparison our method with BVMR.

| Methods | Parameters | PNSR | SSIM |
|---------|-----------|------|------|
| BVMR (original) | **20.51M** | 38.90 | 0.9873 |
| BVMR (channel=44) | 38.77M | 37.18 | 0.9832 |
| BVMR (channel=48) | 46.14M | 37.78 | 0.9847 |
| BVMR (depth=6) | 82.22M | 36.01 | 0.9812 |
| Ours | 32.62M | **42.01** | **0.9928** |

work. However, the final results in RefineNet perform better. It is also noticeable that although the predicted mask is not always completely correct in our network, it has minimal influence for the final prediction. More visualization results can be found in the supplementary material.

**Limitation**  Our method also suffer some limitations. As shown in the Fig. 7, when the detection fails (human in the logo cannot be detected) or the textures in the watermark and background are similar, our network cannot recover the watermark perfectly. However, these issues can often be ameliorated by a larger dataset or a stronger network. Another important limitation is the running speed, the proposed method runs 31fps on 256x256 images which is slower than baseline BVMR (67fps), however, the accuracy is more meaningful than speed since our task is often used as a post-processing tool.
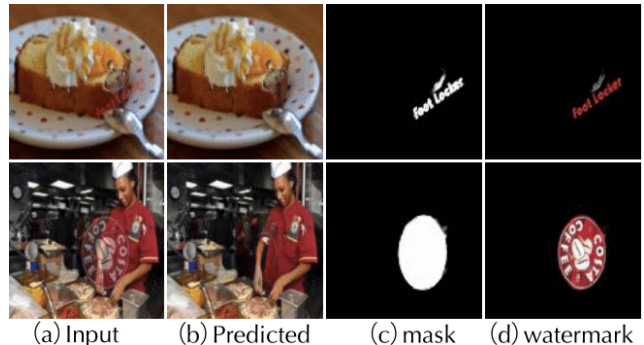


(a) Input  (b) Predicted  (c) mask  (d) watermark

Figure 7: The limitation of the proposed network. Here, (b), (c) and (d) are the predicted results from our framework.

## Conclusion

Following the observation that detection is much easier than removal, in this paper, we propose a novel two-stage framework, SplitNet and RefineNet, for blind single image-based visible watermark removal. The SplitNet obtains the benefits from multi-task learning to generate the coarser outputs (watermark, mask and background). Also, in SplitNet, inspired by multi-domain learning, we build a compact network by sharing the parameters in the main stream decoders, while learning the task-specific attention individually. Then, the RefineNet utilizes the outputs from the previous stage and learns to refine the predicted region with spatial attention mechanism. Besides blind visual motif/watermark removal, our method could also be applied to other related tasks, such

as blind image harmonization, shadow removal, and reflection removal in future work.

## References

Cheng, D.; Li, X.; Li, W.-H.; Lu, C.; Li, F.; Zhao, H.; and Zheng, W.-S. 2018. Large-scale visible watermark detection and removal with deep convolutional networks. In *Chinese Conference on Pattern Recognition and Computer Vision (PRCV)*, 27–40. Springer.

Cun, X.; and Pun, C.-M. 2020. Improving the Harmony of the Composite Image by Spatial-Separated Attention Module. *IEEE Transactions on Image Processing* 29: 4759–4771.

Cun, X.; Pun, C.-M.; and Shi, C. 2020. Towards Ghost-Free Shadow Removal via Dual Hierarchical Aggregation Network and Shadow Matting GAN. *Proceedings of the AAAI Conference on Artificial Intelligence* 34(07): 10680–10687.

Dekel, T.; Rubinstein, M.; Liu, C.; and Freeman, W. T. 2017. On the effectiveness of visible watermarks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2146–2154.

Ding, B.; Long, C.; Zhang, L.; and Xiao, C. 2019. AR-GAN: Attentive Recurrent Generative Adversarial Network for Shadow Detection and Removal. In *The IEEE International Conference on Computer Vision (ICCV)*.

Gandelsman, Y.; Shocher, A.; and Irani, M. 2019. doubledip": Unsupervised image decomposition via coupled deep-image-priors. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 6, 2.

He, K.; Zhang, X.; Ren, S.; and Sun, J. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 770–778.

Hertz, A.; Fogel, S.; Hanocka, R.; Giryes, R.; and Cohen-Or, D. 2019. Blind visual motif removal from a single image. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 6858–6867.

Hu, J.; Shen, L.; Albanie, S.; Sun, G.; and Wu, E. 2018. Squeeze-and-Excitation Networks. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* .

Huang, C.-H.; and Wu, J.-L. 2004. Attacking visible watermarking schemes. *IEEE transactions on multimedia* 6(1): 16–30.

Isola, P.; Zhu, J. Y.; Zhou, T.; and Efros, A. A. 2017. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 5967–5976.

Justin, J.; Alexandre, A.; and Li, F.-F. 2016. Perceptual Losses for Real-Time Style Transfer and Super-Resolution. *Proceedings of the European Conference on Computer Vision (ECCV)* .

Kingma, D. P.; and Ba, J. 2014. Adam: A Method for Stochastic Optimization. *ICLR* .

Li, X.; Lu, C.; Cheng, D.; Li, W.-H.; Cao, M.; Liu, B.; Ma, J.; and Zheng, W.-S. 2019. Towards Photo-Realistic Visible Watermark Removal with Conditional Generative Adversarial Networks. In *International Conference on Image and Graphics*, 345–356. Springer.

Lin, T.-Y.; Maire, M.; Belongie, S. J.; Hays, J.; Perona, P.; Ramanan, D.; Dollár, P.; and Zitnick, C. L. 2014. Microsoft COCO: Common Objects in Context. In *European Conference on Computer Vision (ECCV)*, 740–755.

Liu, G.; Reda, F. A.; Shih, K. J.; Wang, T.-C.; Tao, A.; and Catanzaro, B. 2018. Image inpainting for irregular holes using partial convolutions. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 85–100.

Liu, S.; Johns, E.; and Davison, A. J. 2019. End-to-end multi-task learning with attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 1871–1880.

Park, J.; Tai, Y.-W.; and Kweon, I. S. 2012. Identigram/watermark removal using cross-channel correlation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 446–453. IEEE.

Paszke, A.; Gross, S.; Massa, F.; Lerer, A.; Bradbury, J.; Chanan, G.; Killeen, T.; Lin, Z.; Gimelshein, N.; Antiga, L.; Desmaison, A.; Kopf, A.; Yang, E.; DeVito, Z.; Raison, M.; Tejani, A.; Chilamkurthy, S.; Steiner, B.; Fang, L.; Bai, J.; and Chintala, S. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems 32*, 8024–8035. Curran Associates, Inc. URL http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf.

Pei, S.-C.; and Zeng, Y.-C. 2006. A novel image recovery algorithm for visible watermarked images. *IEEE Transactions on Information Forensics and Security* 1(4): 543–550.

Ruder, S. 2017. An overview of multi-task learning in deep neural networks. *arXiv preprint arXiv:1706.05098* .

Simonyan, K.; and Zisserman, A. 2014. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556* .

Ulyanov, D.; Vedaldi, A.; and Lempitsky, V. 2018. Deep image prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 9446–9454.

Ulyanov, D.; Vedaldi, A.; and Lempitsky, V. S. 2016. Instance Normalization: The Missing Ingredient for Fast Stylization. *arXiv preprint arXiv:1607.08022* .

Wang, G.; and Forsyth, D. 2009. Joint learning of visual attributes, object classes and visual saliency. In *IEEE 12th International Conference on Computer Vision (ICCV)*, 537–544. IEEE.

Wang, J.; Li, X.; and Yang, J. 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.

Wang, X.; Cai, Z.; Gao, D.; and Vasconcelos, N. 2019. Towards universal object detection by domain attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 7289–7298.

Wang, Z.; Bovik, A. C.; Sheikh, H. R.; and Simoncelli, E. P. 2004. Image quality assessment: from error visibility to structural similarity. *IEEE transactions on image processing* 13(4): 600–612.

Xiao, J.; Gu, S.; and Zhang, L. 2020. Multi-Domain Learning for Accurate and Few-Shot Color Constancy. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 3258–3267.

Xu, B.; Wang, N.; Chen, T.; and Li, M. 2015. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853* .

Yang, J.; Gong, D.; Liu, L.; and Shi, Q. 2018. Seeing deeply and bidirectionally: A deep learning approach for single image reflection removal. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 654–669.

Yu, J.; Lin, Z.; Yang, J.; Shen, X.; Lu, X.; and Huang, T. S. 2018. Generative image inpainting with contextual attention. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 5505–5514.

Yu, J.; Lin, Z.; Yang, J.; Shen, X.; Lu, X.; and Huang, T. S. 2019. Free-form image inpainting with gated convolution. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 4471–4480.

Zhang, R.; Isola, P.; Efros, A. A.; Shechtman, E.; and Wang, O. 2018. The Unreasonable Effectiveness of Deep Features as a Perceptual Metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Zhang, X.; Ng, R.; and Chen, Q. 2018. Single Image Reflection Separation with Perceptual Losses. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* 4786–4794.

Zhao, H.; Gallo, O.; Frosio, I.; and Kautz, J. 2016. Loss functions for image restoration with neural networks. *IEEE Transactions on computational imaging* 3(1): 47–57.
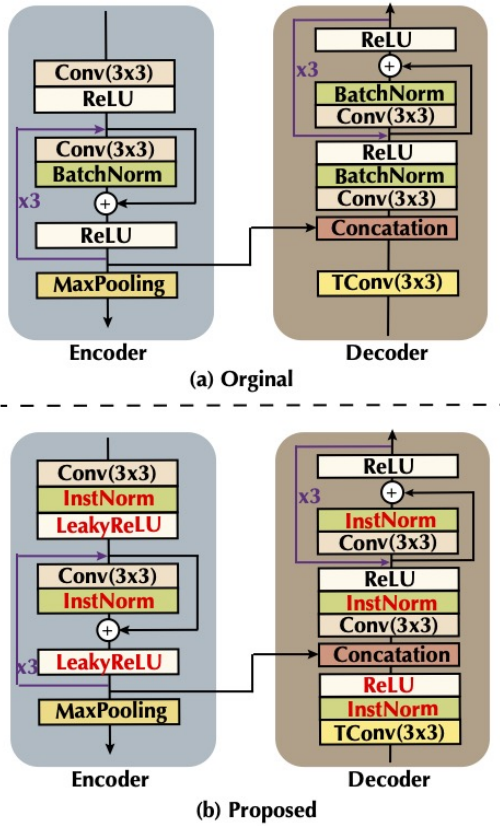
# Supplemental Materials



Figure 8: The differences between the proposed improved ResBlock and the original ResBlock. We mark the modification in red for easier recognization. Although we merely change few layers, our method achieves better and robuster results.

**Differences in ResBlocks**   We find the training process is unstable if we directly deploy the network structure (Hertz et al. 2019). Thus, we make several improvements for better stability. As shown in Figure 8, we illustrate a single layer of the proposed improved ResBlock in the ResBlock based UNet comparing with the original. In detail, we replace all the batch normalization with instance normalization (Ulyanov, Vedaldi, and Lempitsky 2016) because our task is more similar to the style transformation task in the specific region, and the normalization should be applied for each sample. Then, we concatenate all the features after the non-linear activation other than the mixture of convolutional feature and non-linear activation in the previous work (Hertz et al. 2019). In the decoder, we also replace the original ReLU with LeakyReLU (Xu et al. 2015), which is similar to UNet in (Isola et al. 2017).

The proposed structure can hugely improve network stability and performance. As shown in Figure 9, we plot the PSNR between the original ResBlock (ResB) and the proposed block (iResB) on the validation set of LOGO-L dataset. With the help of the carefully designed block, the
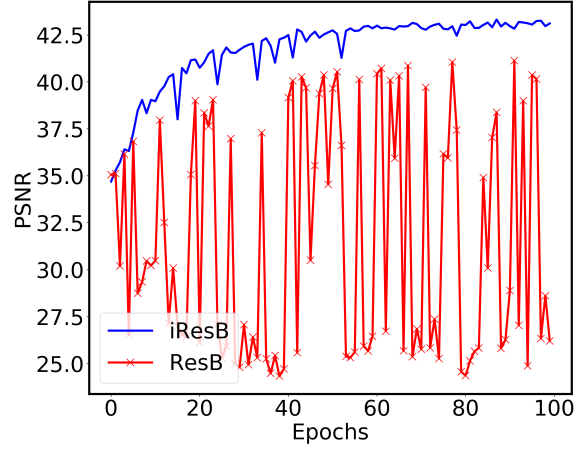


Figure 9: The PSNR of the ResBlock and the proposed iResBlock under the *same* network structure on the validation set of LOGO-L. It is obvious that the proposed block is more stable in training.

proposed network can achieve better performance, and the volatility of the curve is more stable.

**Dataset Analysis**   Using the training set of LOGO-Gray dataset as an example, we analyze the distribution of the type, transparency, size of the watermark in the proposed dataset. In Figure 12, the proposed dataset contains around 300 different classes (logos) for training. Almost each class has similar samples in the overall dataset, which avoids the training bias. We also plot the distribution of the transparency over all samples (as shown in Figure 10). It shows that the transparency is also equally distributed. Moreover, we analyze the watermark size among this dataset as shown in Figure 11, our dataset contains more small watermarks which is similar to the real-world distribution.
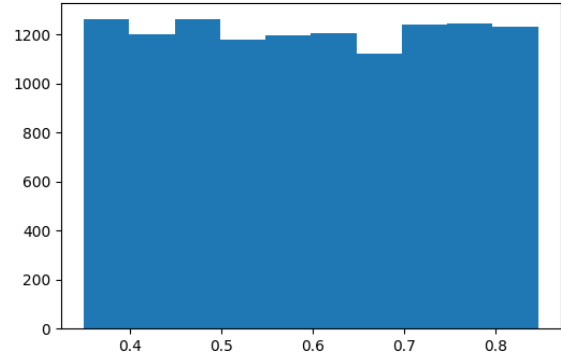


Figure 10: The histogram of the distribution on watermark transparency. The percentage of the transparency are distributed equally among the dataset.

**More results for comparison**   We plot more comparisons on the synthesized dataset to show the benefits of the proposed method as Figure 14 and Figure 15. Besides, we evaluate the baseline network (BVMR) with ours on the "real"
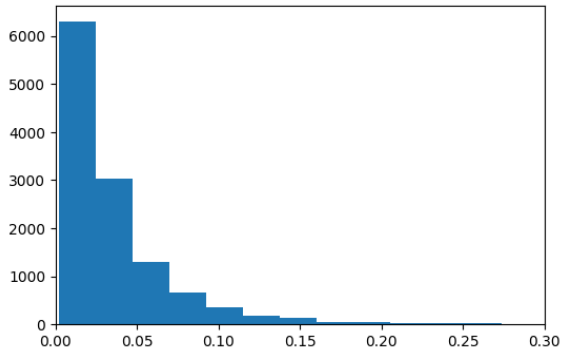
Figure 11: We plot the distribution of the watermark area in the host image. Most watermarks are much smaller than the host image. This is in accord with the real world samples.

dataset in (Dekel et al. 2017). This dataset only contains the input samples and we only present some comparisons visually. From Figure 13, our network can remove and recover the watermark successfully. Notice that, these results are generated using the pre-trained model from LOGO-Gray dataset. This experiment shows that the proposed method generalizes well on the novel scenes.
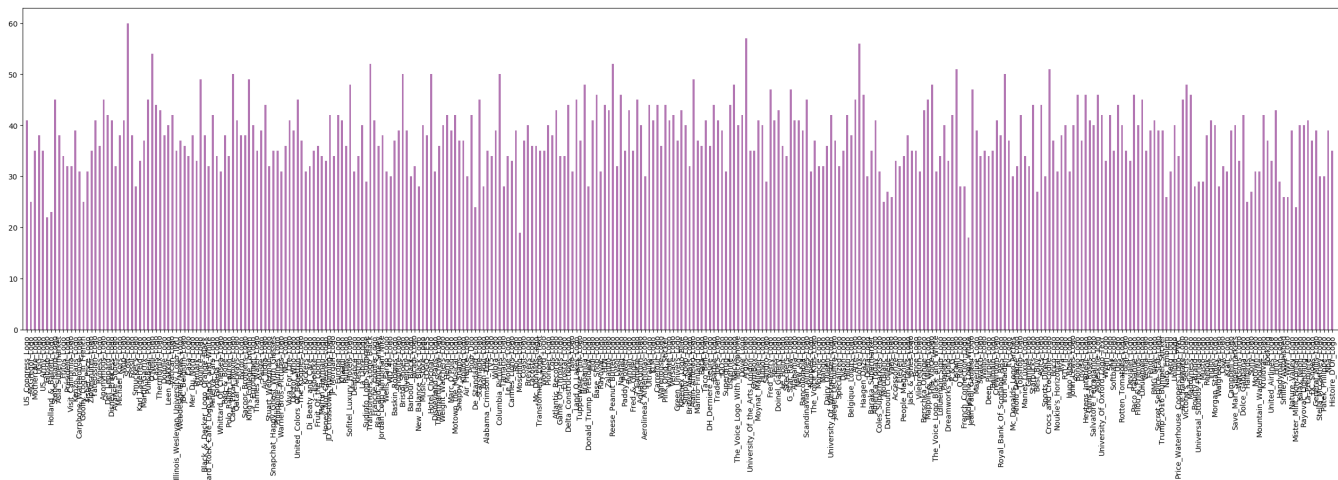
Figure 12: The logo classes analysis in the training set of LOGO-Gray Dataset. The training subset of LOGO-Gray dataset contains 300+ different logos. Almost all the classes are equally distributed in the dataset.
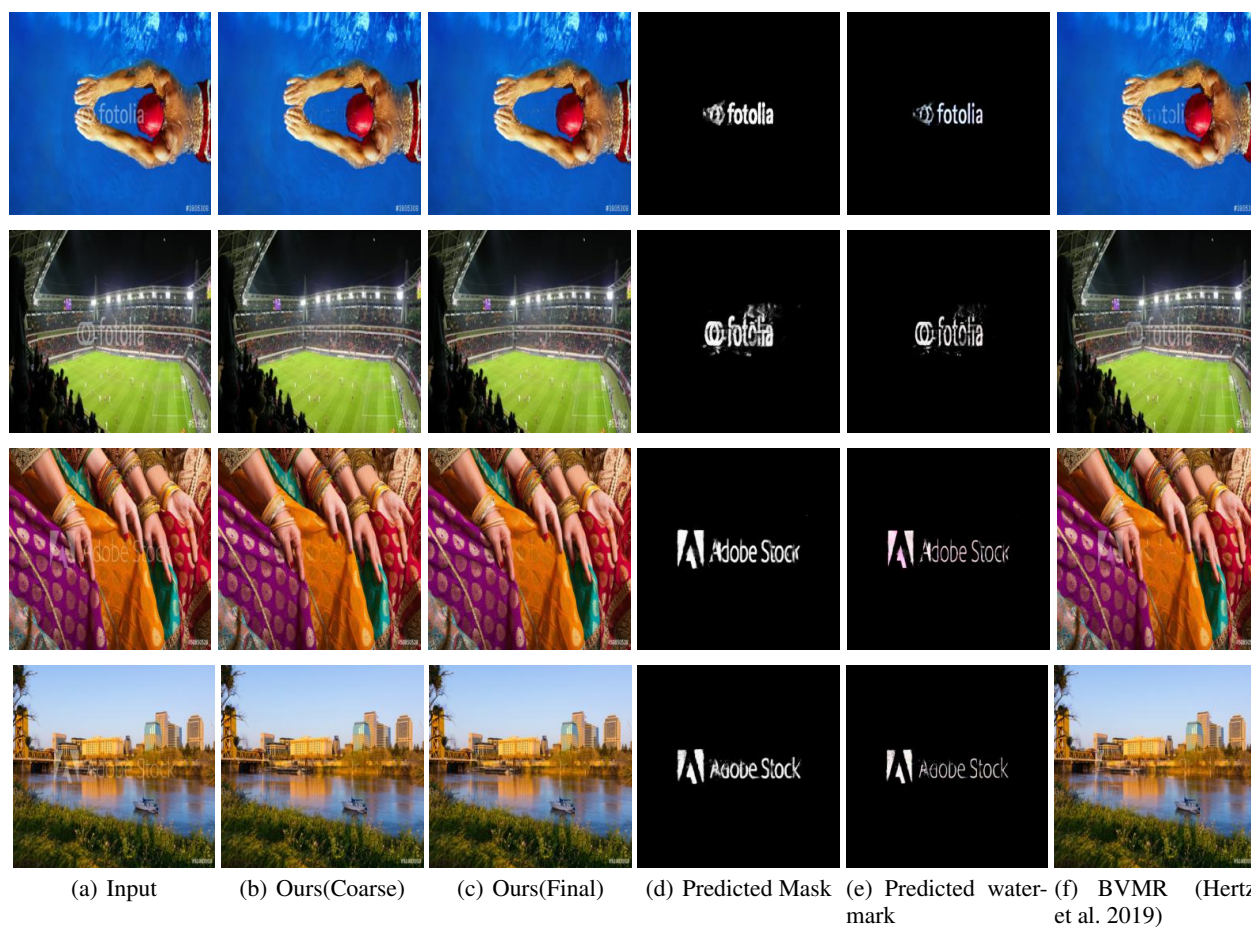


(a) Input　　(b) Ours(Coarse)　　(c) Ours(Final)　　(d) Predicted Mask　(e) Predicted watermark　(f) BVMR (Hertz et al. 2019)

Figure 13: More comparisons between our method and BVMR (Hertz et al. 2019) over the dataset in (Dekel et al. 2017).

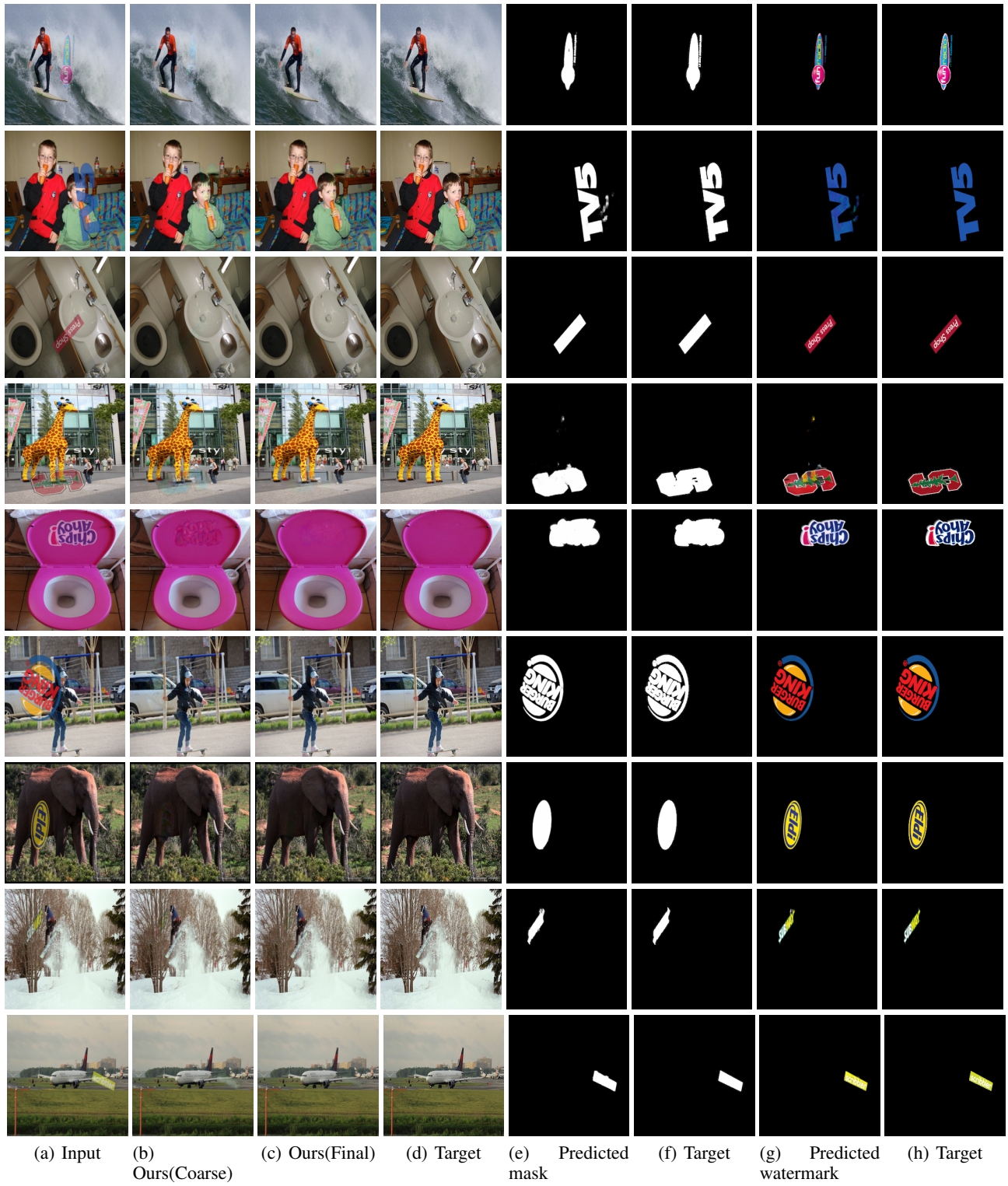(a) Input  (b) Ours(Coarse)  (c) Ours(Final)  (d) Target  (e) Predicted mask  (f) Target  (g) Predicted watermark  (h) Target

Figure 14: More intermedia results of the proposed two-stage framework.

(a) Input    (b) UNet    (c) BS$^2$AM    (d) SIRF    (e) GFSR    (f) BVMR    (g) Ours    (h) Target
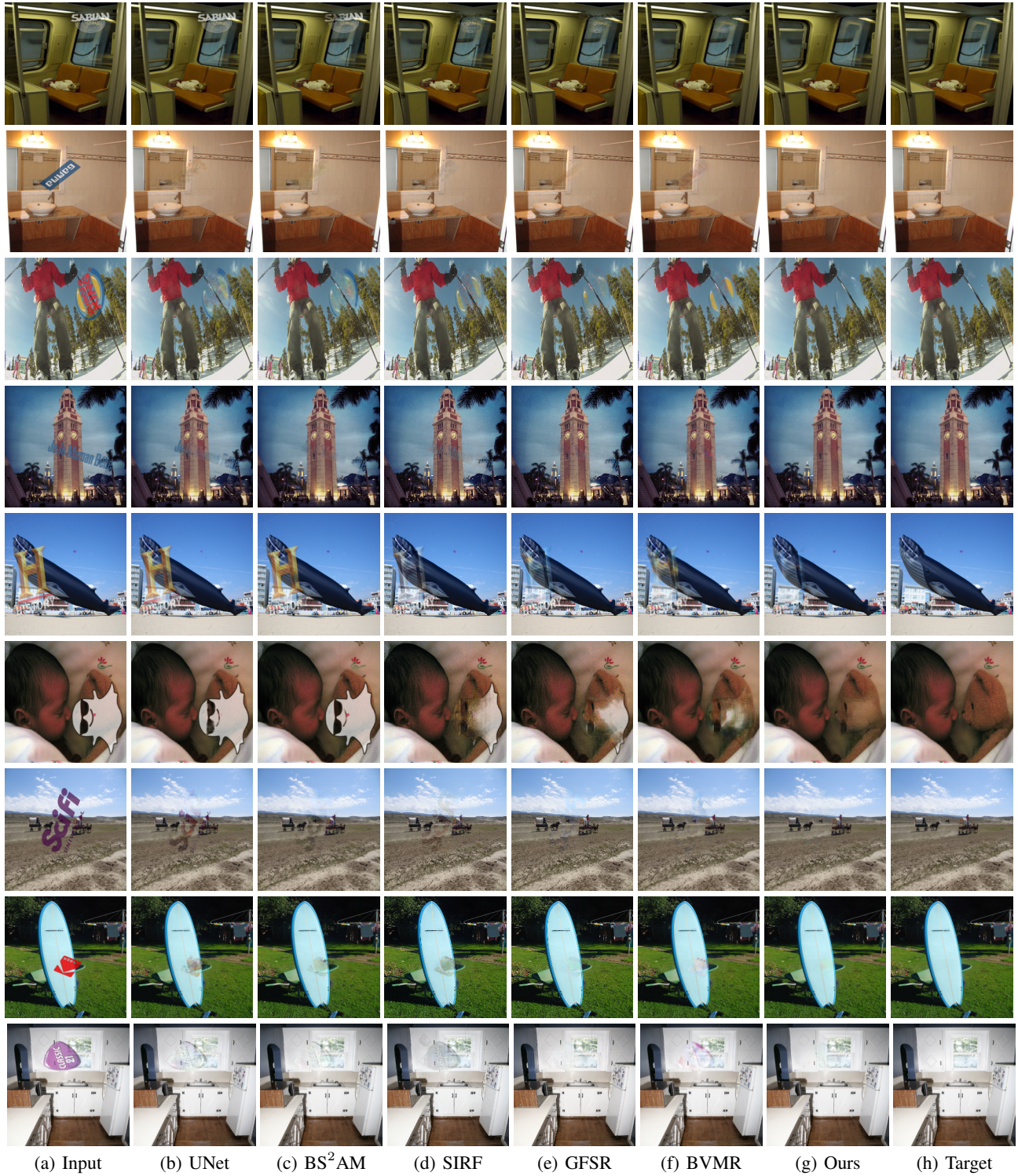
Figure 15: More comparisons between our method and other state-of-the-art methods over different proposed datasets.