

# SDWNet: A Straight Dilated Network with Wavelet Transformation for image Deblurring

Wenbin Zou<sup>1,\*</sup>, Mingchao Jiang<sup>2,\*</sup>, Yunchen Zhang<sup>3,\*</sup>, Liang Chen<sup>1,†</sup>, Zhiyong Lu<sup>2</sup>, Yi Wu<sup>1</sup>  
 Fujian Provincial Key Laboratory of Photonics Technology, Fujian Normal University, Fuzhou, China.<sup>1</sup>  
 JOYY AI GROUP, Guangzhou, China.<sup>2</sup>  
 China Design Group Co., Ltd., Nanjing, China.<sup>3</sup>

alexzou14@foxmail.com, jiangshaoyu1993@gmail.com, cydiachen@cydiachen.tech,  
 cl.0827@126.com, yong1514@gmail.com, wuyi@fjnu.edu.cn

## Abstract

Image deblurring is a classical computer vision problem that aims to recover a sharp image from a blurred image. To solve this problem, existing methods apply the Encode-Decode architecture to design the complex networks to make a good performance. However, most of these methods use repeated up-sampling and down-sampling structures to expand the receptive field, which results in texture information loss during the sampling process and some of them design the multiple stages that lead to difficulties with convergence. Therefore, our model uses dilated convolution to enable the obtainment of the large receptive field with high spatial resolution. Through making full use of the different receptive fields, our method can achieve better performance. On this basis, we reduce the number of up-sampling and down-sampling and design a simple network structure. Besides, we propose a novel module using the wavelet transform, which effectively helps the network to recover clear high-frequency texture details. Qualitative and quantitative evaluations of real and synthetic datasets show that our deblurring method is comparable to existing algorithms in terms of performance with much lower training requirements. The source code and pre-trained models are available at <https://github.com/FlyEgle/SDWNet>.

## 1. Introduction

With the increasing ease of access to images, it is inevitable that blurred images will be obtained in different ways. It is increasingly important to eliminate the blur and restore a clear image. Since the process of image blurring is a one-to-many process, image deblurring is a notoriously difficult ill-posed problem in the field of im-

\*Equal contribution

†Corresponding author

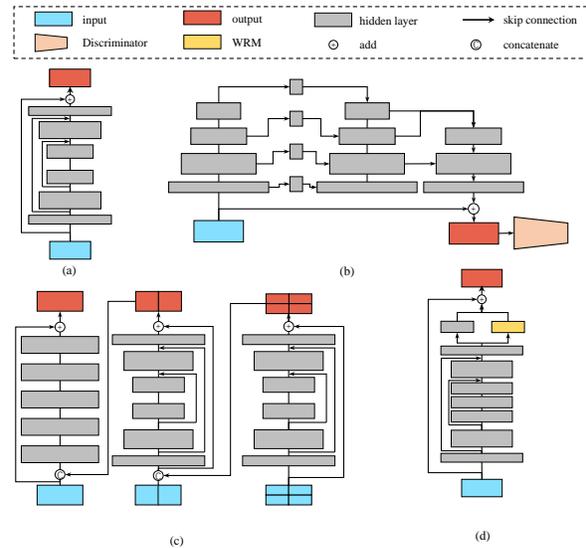


Figure 1. Schematic diagram of current mainstream network architecture. (a). Encode-Decode structure. (b). Generating adversarial network (GAN) structure (c). Coarse-to-fine structure. (d). Ours

age processing [1]. To address this problem, a number of optimization-based [2, 3, 4, 5, 6, 7, 8] and learning-based methods [9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19] have been proposed to learn the mapping function between the clear image and blurry image pairs. Most traditional deblurring methods [5, 6, 7, 8] tackle this problem via estimating blur kernel. Due to blur kernels in natural images are very complex, estimating the best blur kernel is a very tricky task. Therefore, inaccurate estimation of blur kernels results in poorly recovered images.

Recently, the convolution neural network-based (CNN-based) algorithms achieve remarkable progress in image deblurring. Gong *et al.* [11] employ estimated dense motion flow maps to help the model learn the mapping between

clear and blurred images. Then, Nah *et al.* [12] propose a multiscale loss function to implement a coarse-to-fine processing method and achieve good performance. However, this network is complex and very difficult to train. To address the difficulty of training, Tao *et al.* [13] and Gao *et al.* [14] improve the work by using shared network weights between different scales to achieve excellent performance. On this basis, Zhang *et al.* [20] propose an end-to-end CNN multilayer model similar to spatial pyramid matching. Kupyn *et al.* [16, 17] propose DeblurGAN and DeblurGAN-v2 based on adversarial learning to recover more realistic texture details from the blurry image. Shen *et al.* [18] propose a human-aware attentive deblurring network to remove the motion blur between foreground humans and background. Suin *et al.* [19] propose an efficient deblurring design built on new convolutional modules that learn the transformation of features using global attention and adaptive local filters to achieve superior performance. However, significant challenges remain in single image deblurring, as follows:

1. Most of the above methods employ an Encode-Decode structure to learn the features of different receptive fields, as in Figure 1 (a). However, the repeated up-sampling and down-sampling contained in the Encode-Decode structure results in the loss of texture details, which affects the recovery of the image seriously.
2. Some current image deblurring methods use GAN structures to obtain realistic texture details, as in Figure 1 (b). Since the GAN structure requires a joint generator and discriminator for training, it leads to unstable network performance.
3. Most current image deblurring methods tend to design a coarse-to-fine structure to achieve superior PSNR performance, as shown in Figure 1 (c). However, coarse-to-fine structures are often very complex and computationally intensive resulting in a slow convergence process.

In this paper, we address the above challenges using the method of dilated convolution and wavelet transform. We propose a novel image deblurring method that exploits the deblurring cues at different receptive field via a dilated convolution model. Specifically, we propose a simple yet efficient end-to-end CNN model in the wavelet domain called straight dilated network with wavelet transformation (SDWNet), as in Figure 1 (d). It consists of the dilated convolution module and the wavelet reconstruction module. The dilated convolution module uses dilated convolution to obtain a larger field of perception for this network. This helps the model to capture similar features at a distance and thus facilitates image recovery. The wavelet reconstruction module provides additional information for spatial domain re-

construction by exploiting the frequency domain properties of the wavelet transform. Extensive experiments and ablation analysis demonstrate that with the assistance of the dilated convolution module and the wavelet reconstruction module, our SDWNet can achieve state-of-the-art performance.

The contributions of this work are summarized as follows:

- We propose a dilated convolution module. Unlike previous deblurring networks that use repeated up-sampling and down-sampling to obtain large receptive fields, we use the dilated convolution with different dilated rates to obtain features with different receptive fields. This module facilitates the network to capture non-local similar features and recovers a clear image.
- We propose a wavelet reconstruction module. Instead of performing deblurring in a single spatial/frequency domain, we use the information recovered in the frequency domain to complement the spatial domain, so that the recovered image contains more high-frequency details.
- We propose a novel CNN-based image deblurring method. Different from previous deblurring methods that use a coarse to fine structure, we use a simple and streamlined structure to achieve results that are competitive with state-of-the-art methods. This structure effectively solves the problem of difficult training and slow convergence.

## 2. Related Work

### 2.1. Deep Image Deblurring

Recently, deep learning methods have achieved remarkable success in low-level computer vision tasks including image denoise [21], image super-resolution [22, 23, 24], and image deblurring [9, 12, 13, 16, 17]. Many researchers tend to use deep learning methods to design an end-to-end model to achieve excellent performance. Sun *et al.* [9] design a CNN-based model to remove non-uniform motion blur by estimating the blur kernel. Due to the complexity of blurring in real scene images, the blur kernel estimation does not remove the fuzziness completely. Many deep learning-based methods tend to predict clear images directly from blurred images. Nah *et al.* [12] propose a multi-scale CNN model using a coarse-to-fine strategy, which can directly recover latent images without assuming any blur kernel. Because this network does not share parameters at different scales, it leads to increased computation and inference time. To address this problem, Tao *et al.* [13] propose an encoder-decoder structure with jump connections

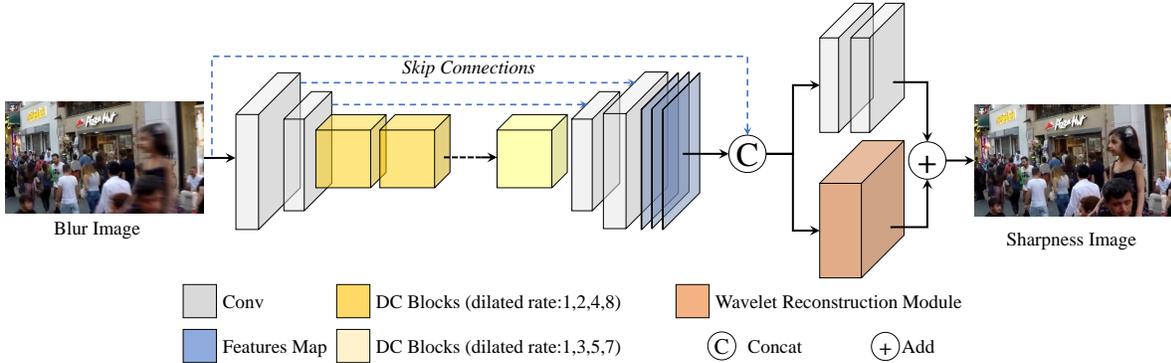


Figure 2. Network architecture of our proposed SDWNet.

and parameter sharing at three scales, which effectively reduces computational effort and achieves better deblurring performance. Kupyn *et al.* [16, 17] propose DeblurGAN and DeblurGAN-v2 using adversarial learning and pyramidal structures to effectively recover clear images. Most of these networks perform alternate down-sampling and up-sampling of deep features to obtain large fields of perception. However, alternate up-sampling and down-sampling can cause a lot of information to be lost in the image recovery process, resulting in poor image recovery results. To address this problem, we use dilated convolution with different dilated rates to obtain information about the different receptive fields, thus making the recovered image clearer.

## 2.2. Dilated Convolution

Dilated convolution can obtain data features of different receptive fields by the jump step size and keeps the parameters constant. On this basis, dilated convolution has been successfully applied in many advanced vision tasks. Yu *et al.* [25] introduce dilated convolution for use on semantic segmentation, which significantly improves the segmentation performance. Zhou *et al.* [26] propose a cascade dilated module for medical image segmentation using convolutional layers with different dilated rates. Then, Brehm *et al.* [27] introduce dilated convolution to the task of image deblurring and achieved excellent performance. They design an atrous convolution block using different dilated rates to recover sharper images. Due to the dilated rate of the atrous blocks follows the semantic segmentation method, the network does not completely cover all the pixel points, resulting in a still blurry image. We are inspired by their network and carefully adjust the dilated rate to obtain almost complete coverage of the receptive field.

## 2.3. Related Application based on WT

The wavelet transform is widely used in image processing tasks because it separates high-frequency information from low-frequency information in an image and is reversible. Many researchers introduce wavelet transforms

into image restoration tasks [28, 15, 24]. Min *et al.* [28] use the wavelet transform to separate the frequency information from the blurred image and then recover the image, effectively weakening the smoothing characteristics of the image. Zhang *et al.* [15] propose double discrete wavelet transform to enhance the blurred image processing capabilities. Liu *et al.* [24] propose a multilayer wavelet CNN using the U-Net structure, resulting in a clearer recovered image. These methods all use a direct mix of all frequency information, leading to problems with different frequency information interacting with each other and creating wrong textures. Therefore, we propose a wavelet transform reconstruction module that effectively recovers a clear image.

## 3. Proposed Method

### 3.1. Framework

In this section, we describe our proposed straight dilated network with wavelet transformation in detail. Since complex models can bring problems such as unstable training and slow convergence, thus we used a plain network structure, as shown in Figure 2. Our SDWNet mainly consists of three parts: the shallow feature extraction layer, the dilated convolution (DC) module, and the reconstruction module. To obtain a larger perceptual field, we first utilize a kernel size of  $7 \times 7$  convolution to extract shallow features. Inspired by the [26], we propose the dilated convolution blocks for fusing multi-receptive field information by using different dilated rates. Then, our network uses cascading multiple DC blocks to learn the broad contextual information. Due to the wavelet transform is an effective tool for recovering high-frequency information, we propose a wavelet reconstruction module as a parallel reconstruction branch, thereby preserving the desired fine texture in the final output image.

Unlike the cascade of multiple dilated convolution blocks with the same dilated rate in [26], we designed two dilated convolution blocks with different dilated rates to obtain richer receptive field information. Besides, we add

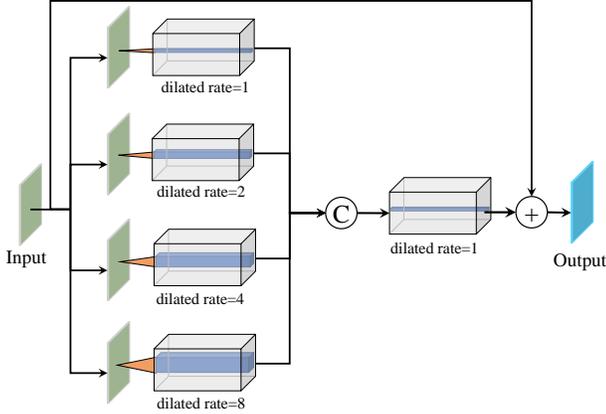


Figure 3. The architecture of our proposed dilated convolution block (DCB). We control the output of the different receptive fields by modifying the dilated rate of the intermediate parallel dilated convolution. We set the dilated rate to  $\{1, 3, 5, 7\}$  in the last layer of DCB. All other DCB dilated rates are set to  $\{1, 2, 4, 8\}$ .

jump connections to make full use of the information from the shallow features. Instead of other wavelet transform methods that predict the four frequency subbands directly, our method uses a shared network to recover the four frequency subbands separately, thus effectively avoiding artifacts caused by the interaction of different frequency subbands.

Given an input blurred image  $\mathbf{I}_{blur}$ , the proposed model predicts a residual image  $\mathbf{R}$  to which the degraded input image  $\mathbf{I}_{blur}$  is added to obtain:  $\mathbf{X} = \mathbf{I}_{blur} + \mathbf{R}$ . We optimize our SDWNet with the following loss function:

$$\mathcal{L}_{total} = \mathcal{L}_{char}(\mathbf{X}, \mathbf{Y}) + \lambda * \mathcal{L}_{ssim}(\mathbf{X}, \mathbf{Y}), \quad (1)$$

where  $\mathbf{Y}$  represents the ground-truth image, and  $\mathcal{L}_{char}$  is the Charbonnier loss [29]:

$$\mathcal{L}_{char} = \frac{1}{N} \sum_{i=1}^N \sqrt{\|\mathbf{X}^i - \mathbf{Y}^i\|^2 + \epsilon^2}, \quad (2)$$

with constant  $\epsilon$  empirically set to  $10^{-3}$  for all the experiments. In addition,  $\mathcal{L}_{ssim}$  is the ssim loss, defined as:

$$\mathcal{L}_{ssim} = \frac{1}{N} \sum_{i=1}^N SSIM(\mathbf{X}^i, \mathbf{Y}^i), \quad (3)$$

where  $SSIM(\cdot)$  denotes the SSIM [30] operator. The parameter  $\lambda$  in Eq. (1) is a hyper-parameter used to control the composition of the SSIM loss function. The following experiments will verify it. Next, we describe each key element of our method.

### 3.2. Dilated Convolution Module

We now give more details about our proposed dilated convolution module, which contains  $n$  dilated convolution

blocks (DCB). The DC Module is formulated as:

$$F_n = H_{DCB}^n(H_{DCB}^{n-1}(\cdots H_{DCB}^1(F_0)\cdots)), \quad (4)$$

where  $H_{DCB}^n$  denotes the function of  $n$ -th DCB.  $F_n$  and  $F_1$  represent the input and output of the DC Module. DCB is composed of multiple dilated convolutions with different dilated rates in parallel, as shown in Figure 3. It can be expressed as follows:

$$F_{dr-1} = H_{dr-1}(F_{input}), \quad (5)$$

$$F_{dr-2} = H_{dr-2}(F_{input}), \quad (6)$$

$$F_{dr-4} = H_{dr-4}(F_{input}), \quad (7)$$

$$F_{dr-8} = H_{dr-8}(F_{input}), \quad (8)$$

$$F_{dr-cat} = Concat(F_{dr-1}, F_{dr-2}, F_{dr-4}, F_{dr-8}), \quad (9)$$

where  $H_{dr-1}$ ,  $H_{dr-2}$ ,  $H_{dr-4}$ , and  $H_{dr-8}$  denote dilated convolution operations with dilated rates of 1, 2, 4 and 8, respectively.  $F_{dr-1}$ ,  $F_{dr-2}$ ,  $F_{dr-4}$ , and  $F_{dr-8}$  denote the output of dilated convolutions with different dilated rates. Inspired by [31], we attach fine-grained control on receptive fields. On shallow layers, we adopt regular dilated rates of 1, 2, 4, and 8. On the last layer, we adopt a non-overlapped dilated rate of 1, 3, 5, and 7 to avoid gridding effects for the image deblurring tasks. Then, we use a dilated convolution with a dilated rate of 1 to fuse features from different receptive fields. Finally, we superimpose the fused features onto the input features to get the output. The output features can be written as:

$$F_{fuse} = H_{fuse}(F_{dr-cat}), \quad (10)$$

$$F_{out} = F_{input} + F_{fuse}, \quad (11)$$

where  $H_{fuse}$  denotes the dilated convolution used to fuse the features.  $F_{fuse}$  and  $F_{out}$  denote the fused features and output features.

### 3.3. Wavelet Reconstruction Module

The wavelet reconstruction module (WRM) mainly uses the wavelet transform to convert spatial domain information to the wavelet domain for recovery. As shown in Figure 4, the input feature  $F_{input}$  can be divided into four different frequency sub-bands by the discrete wavelet transform. These frequency sub-bands can be defined as follows:

$$\{F_{LL}, F_{LH}, F_{HL}, F_{HH}\} = \mathbf{DWT}(F_{input}), \quad (12)$$

where  $\mathbf{DWT}(\cdot)$  denotes the operation of the discrete wavelet transform.  $F_{LL}$ ,  $F_{LH}$ ,  $F_{HL}$ , and  $F_{HH}$  denote the feature of four frequency sub-bands, respectively. To avoid interference between the different frequency subbands, each of the four subbands is fed into a 3-layer convolutional network

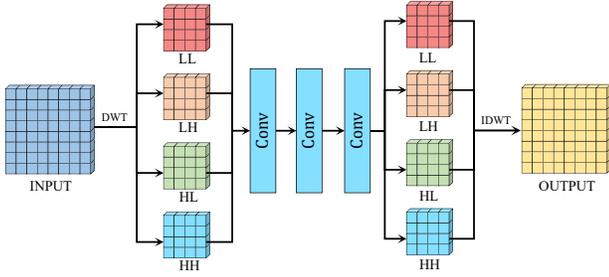


Figure 4. The architecture of our proposed wavelet reconstruction module (WRM). We decompose the input features into four frequency subbands by the wavelet transforms: LL, LH, HL and HH. Then, the corresponding frequency sub-bands are recovered by three-layer convolution. The final output is obtained using the wavelet inverse transform.

for recovery. We can express it as:

$$F_{LL,r} = H_{conv \times 3}(F_{LL}), \quad (13)$$

$$F_{LH,r} = H_{conv \times 3}(F_{LH}), \quad (14)$$

$$F_{HL,r} = H_{conv \times 3}(F_{HL}), \quad (15)$$

$$F_{HH,r} = H_{conv \times 3}(F_{HH}), \quad (16)$$

where  $H_{conv \times 3}(\cdot)$  denotes the 3-layer convolution network.  $F_{LL,r}$ ,  $F_{LH,r}$ ,  $F_{HL,r}$ , and  $F_{HH,r}$  represent the four frequency sub-band features recovered by the 3-layer convolution network. We finally use the discrete wavelet inverse transform to reconstruct the recovered frequency sub-bands into output features  $F_{out}$ . It can be formulated as:

$$F_{out} = \text{IDWT}(F_{LL,r}, F_{LH,r}, F_{HL,r}, F_{HH,r}), \quad (17)$$

where  $\text{IDWT}(\cdot)$  denotes the discrete wavelet inverse transform operation.

## 4. Experiments with Analysis

### 4.1. Data Set

The following are the training and test sets that we use:

The **GoPro** Dataset [12] uses the GoPro Hero 4 camera to capture 240 frames per second (fps) video sequences, and generates blurred images through averaging consecutive short-exposure frames. It is a common benchmark for image motion blurring, containing 3,214 blurry/clear image pairs. We follow the same split [12], to use 2,103 pairs for training and the remaining 1,111 pairs for evaluation.

The **HIDE** Dataset [18] is specifically collected for human-aware motion deblurring and its test set contains 2,025 images. While the GoPro and HIDE datasets are synthetically generated, the image pairs of the RealBlur dataset are captured in real-world conditions.

The **RealBlur** dataset [32] has two subsets: (1). RealBlur-J is formed with the camera JPEG outputs, and

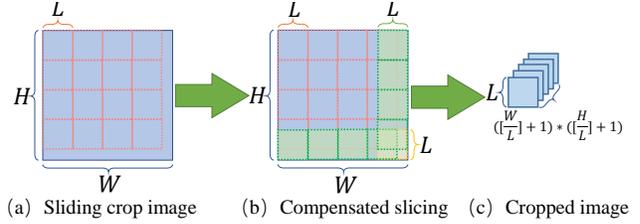


Figure 5. Dataset sliding crop, where  $\lfloor \cdot \rfloor$  denotes represents the rounding down operation. (a) represents the cropping of a block of size  $L \times L$  from an image of species  $H \times W$ . (b) indicates compensatory cropping at the edges of images not covered in (a), marked as green and yellow blocks. The result is (c) a series of the cropped image.

(2). RealBlur-R is generated offline by applying white balance, demosaicking, and denoising operations to the RAW images.

### 4.2. Dataset Sliding Crop

To further improve the robustness of the network, we perform appropriate sliding window slicing on GoPro, as shown in Figure 5. The GoPro dataset images are all  $1280 \times 720$  resolution, so we use a step of 240 to perform  $480 \times 480$  size sliding window slicing in the order of left, right, top and bottom, and compensated slicing on the edge part. Finally, we can crop out 24 patches from each image. Thus, we can crop up to 50472 patches from the original dataset.

### 4.3. Implementation Details

We implement all of the models using PyTorch [33]. Our SDWNet is an end-to-end trainable network and requires no pretraining. Compare with others methods, our network needs fewer training epochs. In the training stage, we use the AdamW [34] optimizer to train our model. We set the input size to  $416 \times 416$  and the batch size to 8. The initial learning rate of  $4 \times 10^{-4}$ , and we use the Cosine Annealing strategy [35] to steadily decrease the learning rate. Weights decay is setting as  $1 \times 10^{-4}$  for the regularization model. We set the hyperparameter  $\lambda$  in the loss function to 1. And, we use a data augmentation strategy of random rotation, random flip, and RGB channel shuffle. We first use the GoPro datasets to train 1500 epochs with the above configuration. Then, we train 50 epochs on the GoPro crop datasets with the best model to get the best results. Besides, all experiments are conducted on the desktop computer with two NVIDIA Tesla V100 GPUs.

### 4.4. Image Deblurring Results

**Quantitative results.** Quantitative analyses are performed to evaluate the performance of the SDWNet for image deblurring. More precisely, we quantitatively assess the average performance of PSNR and SSIM over GoPro and



Figure 6. Qualitative comparison with the leading deblurring algorithms: SRN [13], DeblurGAN-v2 [17], Gao *et al.* [14], MTRNN [10], and DBGAN [38]. From the figure, we can see that our method can generate the right and clear details of the image.

HIDE datasets. We compare our SDWNet with the excellent deblurring methods [5, 7, 36, 11, 16, 12, 37, 17, 10, 39] of the past and the experimental results are shown in Table 1. From Table 1 it can be seen that our method can achieve better performance compared with other deblurring

methods. Compared with the previous DMPHN method, our method achieves 0.16dB improvement in PSNR and 0.027 improvements in SSIM. It is worth noting that not only does our method achieve the best performance on the GoPro dataset, but it also achieves the best results on the

Table 1. Quantitative comparisons of our models with the state-of-the-art deblurring methods on GoPro [12] and HIDE [18] datasets (PSNR(dB)/SSIM). Best and second-best results are **highlighted** and underlined. \* represents the training results of our method on a cropped GoPro dataset.

Method	GoPro		HIDE	
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$
Xu <i>et al.</i> [5]	21.00	0.741	-	-
Hyun <i>et al.</i> [7]	23.64	0.824	-	-
Whyte <i>et al.</i> [36]	24.60	0.846	-	-
Gong <i>et al.</i> [11]	26.40	0.863	-	-
DeblurGAN [16]	28.70	0.858	24.51	0.871
Nah <i>et al.</i> [12]	29.08	0.914	25.73	0.874
Zhang <i>et al.</i> [37]	29.19	0.931	-	-
DeblurGAN-v2 [17]	29.55	0.934	26.61	0.875
SRN [13]	30.26	0.934	28.36	0.915
Shen <i>et al.</i> [18]	-	-	28.89	0.930
Gao <i>et al.</i> [14]	30.90	0.935	29.07	0.913
DBGAN [38]	31.10	0.942	28.94	0.915
MT-RNN [10]	31.15	0.945	<u>29.15</u>	0.918
DMPHN [39]	31.20	0.940	29.09	0.924
<b>SDWNet(Ours)</b>	<u>31.26</u>	<u>0.966</u>	28.99	<u>0.957</u>
<b>SDWNet*(Ours)</b>	<b>31.36</b>	<b>0.967</b>	<b>29.23</b>	<b>0.963</b>

Table 2. Deblurring comparisons on the RealBlur dataset [32] under two different settings: 1). applying our GoPro trained model directly on the RealBlur set (to evaluate generalization to real images), 2). Training and testing on RealBlur data where methods are denoted with symbol \*. The PSNR/SSIM scores for other evaluated approaches are taken from the RealBlur benchmark [32].

Method	RealBlur-R		RealBlur-J	
	PSNR $\uparrow$	SSIM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$
Hu <i>et al.</i> [40]	33.67	0.916	26.41	0.803
Nah <i>et al.</i> [12]	32.51	0.841	27.87	0.827
DeblurGAN [16]	33.79	0.903	27.97	0.834
Pan <i>et al.</i> [41]	34.01	0.916	27.22	0.790
Xu <i>et al.</i> [5]	34.46	0.937	27.14	0.830
DeblurGAN-v2 [17]	35.26	0.944	<b>28.70</b>	0.866
Zhang <i>et al.</i> [37]	35.48	0.947	27.80	0.847
SRN [13]	35.66	0.947	28.56	<u>0.867</u>
DMPHN [39]	<u>35.70</u>	<u>0.948</u>	28.42	0.860
<b>SDWNet(Ours)</b>	<b>35.85</b>	<b>0.948</b>	<u>28.61</u>	<b>0.867</b>
DeblurGAN-v2* [17]	36.44	0.935	29.69	0.870
<b>SDWNet*(Ours)</b>	<b>38.21</b>	<b>0.963</b>	<b>30.73</b>	<b>0.896</b>

HIDE at the same time.

To demonstrate the generalization performance of our method in real scenarios, we also perform experimental validation on the RealBlur dataset, as shown in Table 2. Compared to previous best deblurring methods, our SDWNet achieves the best performance on the RealBlur-R dataset. Our method achieves a 2.51dB PSNR performance gain on the RealBlur-R dataset. On the RealBlur-J dataset, we obtain similar PSNR and SSIM performance to the previous best methods.

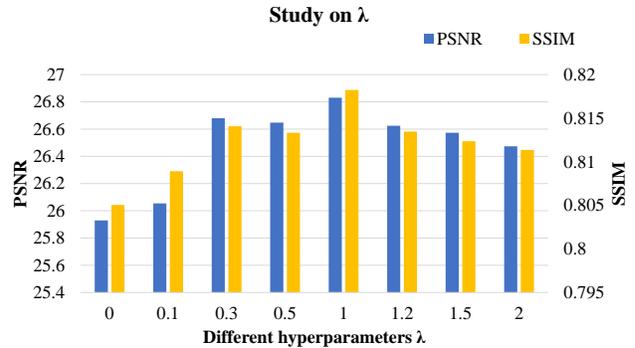


Figure 7. Study on the influence of  $\lambda$ . When  $\lambda = 1$ , the model achieves best results.

**Visual results.** A qualitative analysis of the effect of our SDWNet on image deblurring compared with other methods is shown in Figure 6. We compare the visual deblurring result of our method with the previous methods. To fully demonstrate the superiority of our method, we have zoomed in on the details in the image shown. It is worth noting that many of the detailed textures in the blurred image are difficult to determine. As a result, the repeated up-sampling and downsampling process can cause the texture orientation to change, which can affect image performance. We solve this problem effectively by the dilated convolution to enable the recovered details to be correct. Our method also uses the wavelet transform to convert the features to the frequency domain for recovery, ensuring full recovery of high-frequency details. Therefore, our approach visual results in superior performance. Using the second image in Figure 6 as an example, the images recovered by the older methods still have some blurring. The more recent methods of the last few years have produced images with some error texture. However, our proposed SDWNet can accurately recover a clear image. It demonstrates that our method outperforms other methods in qualitative analysis.

**Performance and efficiency comparison.** In addition to the superior PSNR and SSIM of our model, we also compare the parameters and running times of our method with the previous methods. The results of the experiment are shown in Table 3. Our method has competitive PSNR and SSIM performance to other superior methods, but the parameters and FLOPs of our method are much smaller than other methods, and our method is the fastest in Table 3. Notably, our method achieves better PSNR and SSIM performance than DMPHN [39] using only one-third of the parameters and FLOPs of DMPHN [39]. It efficiently demonstrates that the efficient deblurring performance of our network structure.

#### 4.5. Ablation Studies

In this subsection, we design a series of ablation experiments to analyze the effectiveness of each of the modules

Table 3. Performance and efficiency comparison on the GoPro [12] test dataset. Runtimes are computed with the Nvidia Titan Xp GPU. FLOPs are computed with the input size of  $256 \times 256$ .

Method	DeblurGAN-v2 [17]	DBGAN [38]	DMPHN [39]	Suin <i>et al.</i> [19]	MPRNet [42]	SDWNet (Ours)
Params (M)	60.9	11.6	21.7	23.0	20.1	<b>7.2</b>
Flops (G)	411.34	660.2	678.56	536.74	660.2	<b>181.31</b>
Time (s)	0.21	0.83	1.07	0.34	0.18	<b>0.14</b>
PSNR	29.55	31.10	31.20	31.85	<b>32.66</b>	31.36
SSIM	0.934	0.942	0.940	0.948	0.959	<b>0.967</b>

Table 4. Model Policy with depths and widths on network performance with an input of  $96 \times 96$

Model Setting	PSNR	SSIM
$d = 10, w = 16$	27.18	0.833
$d = 10, w = 32$	27.70	0.848
$d = 16, w = 32$	<b>28.24</b>	<b>0.863</b>
$d = 20, w = 16$	27.21	0.852
$d = 20, w = 32$	27.08	0.828

Table 5. Ablation studies on ELU, Bilinear, Wide, Dilated rate, and WRM. The PSNR Performance on Gopro test dataset.

Operation	ELU	Bilinear	Dilated rate	WRM	PSNR
Baseline	✗	✗	✗	✗	27.36
	✓	✗	✗	✗	27.64
	✓	✓	✗	✗	27.87
	✓	✓	✓	✗	28.04
	✓	✓	✓	✓	<b>28.36</b>

we propose. We use the GoPro test set for evaluation and performed 200 epochs of training on an image patch of size  $96 \times 96$ .

**Model Design Policy.** We explore the impact of different depths and widths on network performance, as shown in Table 4. Where the depth represents the number of DCBs we set and width represents the number of channels in our intermediate features. As can be seen from the experimental results, the width has a greater effect on our model than the depth. Our model works best at  $d = 16$  and  $w = 32$ . Therefore, our final model is set to  $d = 16$  and  $w = 32$ .

**Effectiveness of Each Operation.** We set up a baseline of DCBs, where the activation function is ReLU and the upsampling is deconvolution. We demonstrate the effectiveness of our proposed module by modifying the corresponding activation function and upsampling method, as shown in Table 5. From the experiment, it is known that the ELU activation function and bilinear upsampling obtain better performance than ReLU and deconvolution. We adjust the dilated rate of the last layer of the DC module to help improve the performance of the network. And the WRM can help the network recover high-frequency details in the frequency domain, enabling network performance to be improved. These comparisons show that our proposed methods are useful for image deblurring.

**Effectiveness of  $\lambda$  in Loss Function.** We conduct trade-

off experiments for the Charbonnier loss and SSIM loss, as shown in Figure 7. We find the optimal hyperparameter  $\lambda$  by adjusting the value of the hyperparameter  $\lambda$  so that the network performance can be optimized. We set the hyperparameters  $\lambda$  to 0, 0.1, 0.3, 0.5, 1 and 2 respectively. We can notice from the graph that the best performance is obtained when the hyperparameter  $\lambda = 1$ .

## 5. Conclusion

In this work, we propose a novel dilated convolution network structure for image deblurring. For this structure, we propose two modules: the dilated convolution module and the wavelet reconstruction module. Specifically, the dilated convolution module use dilated convolution with different dilated rates, which effectively helps the network to obtain different receptive field information. The wavelet reconstruction module exploits the properties of the wavelet transform to provide high-frequency information for spatial domain reconstruction, resulting in clearer images. The quantitative and qualitative results show that our algorithm can effectively restore a clearer image than other methods. Our approach is simple in structure and easily transferable to other high-level tasks. In the future, this approach will be explored to facilitate other image restoration tasks such as image denoising and super-resolution.

## Acknowledgments

This work was supported in part by the National Nature Science Foundation of China under Grant No. 61901117, U1805262, 61971165, in part by the Natural Science Foundation of Fujian Province under Grant No. 2019J05060, 201-9J01271, in part by the Fujian Provincial Education Department Project under Grant No. JT180094, JT180095, in part by the Special Fund for Marine Economic Development of Fujian Province under Grant No. ZHHY-2020-3, in part by the research program of Fujian Province under Grant No. 2018H6007, the Special Funds of the Central Government Guiding Local Science and Technology Development under Grant No. 2017L3009, and the National Key Research and Development Program of China under Grant No. 2016YFB-1001001.

## References

- [1] G. Rioux, C. Scarvelis, R. Choksi, T. Hoheisel, and P. Maréchal. Blind deblurring of barcodes via kullback-leibler divergence. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 43(1):77–88, 2021.
- [2] Qi Shan, Jiaya Jia, and Aseem Agarwala. High-quality motion deblurring from a single image. *Acm transactions on graphics (tog)*, 27(3):1–10, 2008.
- [3] D. Zoran and Y. Weiss. From learning models of natural image patches to whole image restoration. In *2011 International Conference on Computer Vision*, pages 479–486, 2011.
- [4] Z. Dou, B. Zhang, and X. Yu. A new alternating minimization algorithm for image segmentation. In *6th International Conference on Wireless, Mobile and Multi-Media (ICWMMN 2015)*, pages 181–184, 2015.
- [5] L. Xu, S. Zheng, and J. Jia. Unnatural l0 sparse representation for natural image deblurring. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1107–1114, 2013.
- [6] R. Fattal and A. Goldstein. Blur-kernel estimation from spectral irregularities. In *IEEE International Conference on Computer Vision (ICCV)*, 2012.
- [7] Tae Hyun Kim, Byeongjoo Ahn, and Kyoung Mu Lee. Dynamic scene deblurring. In *IEEE International Conference on Computer Vision (ICCV)*, 2013.
- [8] Li Xu and Jiaya Jia. Two-phase kernel estimation for robust motion deblurring. In *European conference on computer vision*, pages 157–170. Springer, 2010.
- [9] J. Sun, Wenfei Cao, Zongben Xu, and J. Ponce. Learning a convolutional neural network for non-uniform motion blur removal. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 769–777, 2015.
- [10] Dongwon Park, Dong Un Kang, Jisoo Kim, and Se Young Chun. Multi-temporal recurrent neural networks for progressive non-uniform single image deblurring with incremental temporal training. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*. Springer International Publishing, 2020.
- [11] Dong Gong, Jie Yang, Lingqiao Liu, Yanning Zhang, Ian Reid, Chunhua Shen, Anton Van Den Hengel, and Qinfeng Shi. From motion blur to motion flow: A deep learning solution for removing heterogeneous motion blur. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2319–2328, 2017.
- [12] S. Nah, T. H. Kim, and K. M. Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 257–265, 2017.
- [13] X. Tao, H. Gao, X. Shen, J. Wang, and J. Jia. Scale-recurrent network for deep image deblurring. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8174–8182, 2018.
- [14] H. Gao, X. Tao, X. Shen, and J. Jia. Dynamic scene deblurring with parameter selective sharing and nested skip connections. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3843–3851, 2019.
- [15] Y. Zhang and K. Hirakawa. Blur processing using double discrete wavelet transform. In *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1091–1098, 2013.
- [16] O. Kupyn, V. Budzan, M. Mykhailych, D. Mishkin, and J. Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8183–8192, 2018.
- [17] O. Kupyn, T. Martyniuk, J. Wu, and Z. Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8877–8886, 2019.
- [18] Z. Shen, W. Wang, X. Lu, J. Shen, H. Ling, T. Xu, and L. Shao. Human-aware motion deblurring. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5571–5580, 2019.
- [19] M. Suin, K. Purohit, and A. N. Rajagopalan. Spatially-attentive patch-hierarchical network for adaptive motion deblurring. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [20] K. Zhang, W. Luo, Y. Zhong, L. Ma, B. Stenger, W. Liu, and H. Li. Deblurring by realistic blurring. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2734–2743, 2020.
- [21] Xiao-Jiao Mao, Chunhua Shen, and Yu-Bin Yang. Image denoising using very deep fully convolutional encoder-decoder networks with symmetric skip connections. *arXiv preprint arXiv:1603.09056*, 2, 2016.
- [22] B. Lim, S. Son, H. Kim, S. Nah, and K. M. Lee. Enhanced deep residual networks for single image super-resolution. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1132–1140, 2017.
- [23] T. Guo, H. S. Mousavi, T. H. Vu, and V. Monga. Deep wavelet prediction for image super-resolution. In *2017 IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1100–1109, 2017.
- [24] P. Liu, H. Zhang, K. Zhang, L. Lin, and W. Zuo. Multi-level wavelet-cnn for image restoration. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 886–88609, 2018.
- [25] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015.
- [26] Xiao-Yun Zhou, Jian-Qing Zheng, and Guang-Zhong Yang. Atrous convolutional neural network (acnn) for biomedical semantic segmentation with dimensionally lossless feature maps. *arXiv preprint arXiv:1901.09203*, page 68, 2019.

- [27] S. Brehm, S. Scherer, and R. Lienhart. High-resolution dual-stage multi-level feature aggregation for single image and video deblurring. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1872–1881, 2020.
- [28] C. Min, G. Wen, B. Li, and F. Fan. Blind deblurring via a novel recursive deep cnn improved by wavelet transform. *IEEE Access*, 6:69242–69252, 2018.
- [29] P. Charbonnier, L. Blanc-Feraud, G. Aubert, and M. Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *Proceedings of 1st International Conference on Image Processing*, volume 2, pages 168–172 vol.2, 1994.
- [30] H. Zhao, O. Gallo, I. Frosio, and J. Kautz. Loss functions for image restoration with neural networks. *IEEE Transactions on Computational Imaging*, 3(1):47–57, 2017.
- [31] Jie Liu, Chuming Li, Feng Liang, Chen Lin, Ming Sun, Junjie Yan, Wanli Ouyang, and Dong Xu. Inception convolution with efficient dilation search. *arXiv preprint arXiv:2012.13587*, 2020.
- [32] Jaesung Rim, Haeyun Lee, Jucheol Won, and Sunghyun Cho. Real-world blur dataset for learning and benchmarking deblurring algorithms. In *European Conference on Computer Vision*, pages 184–201. Springer, 2020.
- [33] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.
- [34] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [35] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [36] Andrew Zisserman Oliver Whyte, Josef Sivic and Jean Ponce. Non-uniform deblurring for shaken images. In *International Journal of Computer Vision*, page 168–186, 2012.
- [37] Jiawei Zhang, Jinshan Pan, Jimmy Ren, Yibing Song, Linchao Bao, Rynson WH Lau, and Ming-Hsuan Yang. Dynamic scene deblurring using spatially variant recurrent neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2521–2529, 2018.
- [38] Kaihao Zhang, Wenhan Luo, Yiran Zhong, Lin Ma, Bjorn Stenger, Wei Liu, and Hongdong Li. Deblurring by realistic blurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2737–2746, 2020.
- [39] Hongguang Zhang, Yuchao Dai, Hongdong Li, and Piotr Koniusz. Deep stacked hierarchical multi-patch network for image deblurring. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5978–5986, 2019.
- [40] Zhe Hu, Sunghyun Cho, Jue Wang, and Ming-Hsuan Yang. Deblurring low-light images with light streaks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3382–3389, 2014.
- [41] Jinshan Pan, Deqing Sun, Hanspeter Pfister, and Ming-Hsuan Yang. Blind image deblurring using dark channel prior. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1628–1636, 2016.
- [42] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shabbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *CVPR*, 2021.