

Two-stage Progressive Residual Dense Attention Network for Image Denoising

Wencong Wu^a, An Ge^a, Guannan Lv^a, Yuelong Xia^a, Yungang Zhang^{a,*}, Wen Xiong^{a,*}

^a*School of Information Science and Technology, Yunnan Normal University, Kunming 650500, Yunnan Province, China*

Abstract

Deep convolutional neural networks (CNNs) for image denoising can effectively exploit rich hierarchical features and have achieved great success. However, many deep CNN-based denoising models equally utilize the hierarchical features of noisy images without paying attention to the more important and useful features, leading to relatively low performance. To address the issue, we design a new Two-stage Progressive Residual Dense Attention Network (TSP-RDANet) for image denoising, which divides the whole process of denoising into two sub-tasks to remove noise progressively. Two different attention mechanism-based denoising networks are designed for the two sequential sub-tasks: the residual dense attention module (RDAM) is designed for the first stage, and the hybrid dilated residual dense attention module (HDRDAM) is proposed for the second stage. The proposed attention modules are able to learn appropriate local features through dense connection between different convolutional layers, and the irrelevant features can also be suppressed. The two sub-networks are then connected by a long skip connection to retain the shallow feature to enhance the denoising performance. The experiments on seven benchmark datasets have verified that compared with many state-of-the-art methods, the proposed TSP-RDANet can obtain favorable results both on synthetic and real noisy image denoising. The code of our TSP-RDANet is available at <https://github.com/WenCongWu/TSP-RDANet>.

Keywords: Image denoising, CNN, residual dense attention, hybrid dilated residual dense attention.

1. Introduction

Image denoising is a hot topic in computer vision tasks, which is a key preparatory work for the subsequent high-level tasks so that the performance of these tasks can be improved. Image denoising seeks to solve an inverse problem to obtain the ‘clear’ image from an image contaminated by noise, and it can be described as $x = y - N$ for the additive white Gaussian noise (AWGN), where x , y , and N denote the denoised ‘clear’ image, the noisy image, and the noise, respectively. Many image denoising methods have been designed in recent decades. For example, Dabov et al. [1] presented the block-matching and 3-D filtering (BM3D) for image denoising, which applies the non-local self-similarity (NSS) and an enhanced sparse learning scheme in the transform domain to enhance model performance. The trilateral weighted sparse coding (TWSC) [2] was proposed for practical denoising scenes. Later, Ou et al. [3] designed a weighted group sparse coding model using multi-scale NSS to improve the denoising performance. Some researchers have also tried to tackle the real-world noise removal task. For example, the multi-channel weighted nuclear norm minimization (MCWNNM) [4] was developed for noise removal in real scenes. Although these methods can obtain impressive denoising performance, they need to manually set parameters for different denoising tasks, and they are generally equipped with complex optimization algorithms, which result in high computational cost.

Various deep neural networks (DNNs) based denoising models have been proposed in recent years. Compared with the above mentioned traditional denoising techniques, the DNN-based denoising models need fewer model hyper-parameters, and generally have faster denoising speeds. For example, Zhang et al. [5] presented the denoising CNN (DnCNN), which surpasses most of the traditional methods. Peng et al. [6] designed the dilated residual network, named as the DSNet, where the symmetric skip connection is applied in different convolutional layers to extract hierarchical features, the denoising results of the DSNet therefore are improved. Zhang et al. [7] developed the

*Corresponding author

Email addresses: yungang.zhang@ynnu.edu.cn (Yungang Zhang), wen.xiong@ynnu.edu.cn (Wen Xiong)

residual dense network (RDNet) to restore the degraded image, which uses the densely connected layer to capture rich hierarchical features. Jia et al. [8] proposed the dense dense U-Net (DDUNet) for removing the noise from noisy images, which adopts the multi-scale dense connection to obtain more image features. Although these DNN-based models are able to extract rich features to achieve favorable denoising performance, they treat all the extracted features equally, which may be an obstacle to improving model performance.

To better utilize the more useful and informative image features, several attention-based denoising methods have been developed. Tian et al. [9] presented an attention-guided denoising CNN (ADNet) to eliminate the Gaussian and real noise, and the model contains an attention block to capture the noise information in the noisy images with complex backgrounds. Anwar et al. [10] designed a real image denoising network (RIDNet), which exploits the feature attention block to capture the channel dependencies to further improve denoising quality.

Recently, progressive models have also been developed to promote image denoising effect. Zamir et al. [11] presented a multi-stage architecture MPRNet for image restoration, which utilizes encoder-decoders to capture multi-scale features, and the supervised attention module is developed to refine the filtered features and the degraded images in each stage. Bai et al. [12] developed a progressive denoising network (MSPNet), which decomposes the overall denoising process into multiple sub-steps to improve noise reduction gradually. Although the progressive denoising networks such as the MPRNet and MSPNet can improve denoising performance, the complex network structures bring large numbers of parameters, and their application in real denoising scenarios is therefore limited.

It also can be seen that in the current image restoration area, the Transformers-based models can generally outperform most CNN-based models, nevertheless some researchers have pointed out that the CNN-based models can also reach or even surpass transformers as appropriate network structures and components can be provided [13]. Moreover, the transformer-based models generally suffer from the problems of large model scales and difficulty in training. Instead, the CNN-based models are relatively easier to be deployed on various terminal devices with limited computing resources. Therefore, the CNN-based denoising models still enrich great research potential.

Motivated by the promising performance of the progressive denoising scheme and the attention-based feature learning, we propose a two-stage progressive residual dense attention network for image denoising, named as the TSP-RDANet, which uses two heterogeneous networks to enhance feature interaction and feature representation ability. The residual dense attention module (RDAM) and the hybrid dilated residual dense attention module (HDRDAM) are developed to capture rich hierarchical features, respectively. The residual learning is embedded into these modules to promote denoising performance, and the network training can be accelerated as well. The long skip connections are employed between these attention modules to infuse the relevant and useful features. The major contributions of this work are listed as follows:

(1) We novelly design the residual dense attention module (RDAM) and the hybrid dilated residual dense attention module (HDRDAM). These two different attention modules are utilized in our two-stage progressive denoising models. The RDAM and HDRDAM use densely connected layers to extract rich local features, in which the irrelevant features are filtered by the attention mechanism, and the residual learning is applied both in the RDAM and HDRDAM to enhance the denoising performance of the network.

(2) A novel two-stage progressive residual dense attention network (TSP-RDANet) is proposed for image denoising, which decomposes the entire denoising process into two sub-tasks to progressively restore a noisy image.

(3) Experiments on multiple synthetic and real-world datasets have verified that our TSP-RDANet obtains promising denoising performance compared with many other state-of-the-art models.

The rest of this paper has the following organization. Section 2 shows the related techniques for image denoising. In Section 3, we introduce the proposed TSP-RDANet model and the proposed attention modules. Section 4 offers our experimental details and results. The conclusion is reported in Section 5.

2. Related work

2.1. CNN based image denoising techniques

During the past years, a number of image denoising methods based on convolutional neural networks have been developed [5]. For instance, a denoising CNN (DnCNN) model was proposed by Zhang et al. [5], which uses rectified linear units (ReLU) [14], batch normalization (BN) [15], and residual learning [16] to speed up its training procedure and promote model performance. The image restoration CNN (IRCNN) [17] utilizes the dilated convolution [18] to

enlarge its receptive field to obtain more contextual information. Peng et al. [6] applied the symmetric skip connection and dilated convolution to develop the dilated residual networks (DSNet) for image denoising. The FFDNet [19] introduces an adjustable noise level map to increase the flexibility of Gaussian noise removal, and the FFDNet employs the downsampled sub-images to achieve fast image denoising.

Generally, the image denoising methods based on CNN can obtain better performance than many traditional denoising methods. However, the early CNN-based denoising models are mainly developed for synthetic noise removal, and it is inconvenient to use them in practical denoising scenes. Therefore, many blind denoising methods have been proposed, including the VDN [20], AINDNet [21], BUIFD [22], CBDNet [23], VDIR [24], DCBDNet [25]. Many of these blind denoising models contain a noise estimation sub-network to estimate noise distribution to achieve flexibility in practical image denoising, but the noise estimation networks increase the complexity of the models. Other blind denoising models without the noise estimation sub-networks have also been developed. For instance, the SADNet [26] is a spatial-adaptive model, where the deformable convolution [27, 28] is applied to promote its performance. Quan et al. [29] presented a complex-valued denoising network (CDNet), which investigates the advantages of complex-valued operations in CNN. Li et al. [30] developed the AirNet, which adopts a single AirNet model to achieve the recovery of multiple types of degraded images, and excellent restoration performance was produced.

Some researchers have tried to expand the width of CNN denoising models instead of increasing network depth. For instance, Pan et al. [31] presented a dual CNN (DualCNN) for low-level vision tasks, which uses two different branches to restore the main parts and details of the degraded images, respectively. The restored structures and details from the two branches are then fused to predict the recovery images. A batch-renormalization [32] denoising network (BRDNet) was presented by Tian et al. [33], where two sub-networks were combined to extract the complementary features. Later, a dual denoising network (DudeNet) [34] was developed by the same author. A dual adversarial network (DANet) was designed in [35] for practical noise removal. It can be found that the dual network structure can be an effective way to improve denoising performance.

2.2. Attention

Different attention mechanisms are also widely used in image denoising models, which can obtain more effective features and further promote the denoising performance. The RIDNet proposed by Anwar et al. [10] employs the feature attention to capture the channel relationships in order to suppress irrelevant features. The ADNet [9] utilizes attention-guided feature learning to improve its denoising result. Jiang et al. [36] developed an enhanced frequency fusion network (EFF-Net) for noise removal, in which a dynamic hash attention was designed to enhance its denoising performance. Ren et al. [37] presented the DeamNet, where the dual element-wise attention mechanism module is designed to promote the denoising effect of the network. Zamir et al. [11] presented the CycleISP framework for practical image denoising, which employs spatial attention [38] and channel attention [39] mechanisms to exploit the inter-channel and inter-spatial dependencies. Wu et al. [40] proposed a dual residual attention network for obtaining the denoised image, where the dual-branch structure was used to capture complementary features, and the attention mechanism in [39, 38] was adopted to filter unimportant features. Huang et al. [41] designed a prior-guided dynamic tunable network to achieve real-world noise removal, where the global spatial and channel attention was proposed to capture non-local features. Zhuge et al. [42] presented an enhanced feature denoising network, where the cross-channel attention was embedded into the network to enhance the interaction of channel features. Thakur et al. [43] designed a blind Gaussian denoising network, which applied a multi-scale pixel attention to capture salient features from multiple scales.

2.3. Progressive denoising models

Much previous work has demonstrated that compared with the single-stage counterparts, the multi-stage/progressive models can achieve more effective performance in image restoration, such as image denoising [11, 12, 44], image de-raining [45, 46], and image deblurring [47, 48, 49]. For example, A multi-stage progressive image restoration framework (MPRNet) was developed in [11], which adopts the supervised attention module between different stages to achieve progressive learning. Bai et al. [12] designed a multi-stage progressive denoising network (MSPNet), which divides the overall denoising process into three sub-steps, and excellent noise reduction results were reported. Tian et al. [44] developed a multi-stage CNN with the wavelet transform (MWDCNN), which utilizes the enhanced residual dense architectures to capture enough features for image denoising. A progressive recurrent network (PReNet) was

presented in [46] for image deraining, where a recurrent layer was applied to capture the feature dependencies between different stages. Fu et al. [50] developed a lightweight pyramid network (LPNet) for rain removal, where the laplacian pyramid was used to generate the derained images through multiple different sub-networks. Nah et al. [47] developed a multi-scale progressive network architecture for image deblurring, where the Gaussian pyramid was utilized to generate downsampled blurry images, and these images were fed into progressive sub-networks to produce the deblurred images. Zhang et al. [48] designed a deep multi-patch hierarchical network for blur removal, which decomposes the blurred image into multiple patches gradually, and the image patches are sent into different encoder-decoders to generate the deblurred result.

3. Proposed method

In this section, the architecture of our progressive denoising model TSP-RDANet is illustrated. The proposed residual dense attention module (RDAM) and the hybrid dilated residual dense attention module (HDRDAM) are also introduced in details.

3.1. Architecture of the TSP-RDANet

The framework of our TSP-RDANet is displayed in Fig. 1, which mainly comprises two stages: Stage 1 and Stage 2. The TSP-RDANet can extract rich image features through densely connected modules, including RDAMs and HDRDAMs. To obtain a desirable trade-off between the network performance and complexity, five RDAMs and five HDRDAMs are used for the first and the second stages denoising network, respectively. The strided convolutions and the transpose convolutions between the RDAMs are used for obtaining the multi-scale information and enlarging the receptive field. In HDRDAMs, the dilated convolutions are utilized to expand the receptive field to obtain more contextual information as well. In addition, the TSP-RDANet model adopts long skip connections to fuse the shallow and deep convolutional features, allowing the model to fully utilize the features, thereby improving its repair performance. Furthermore, the long skip connection enables the sub-networks of the two stages to interact on the different salient features to promote the expressive ability of the model.

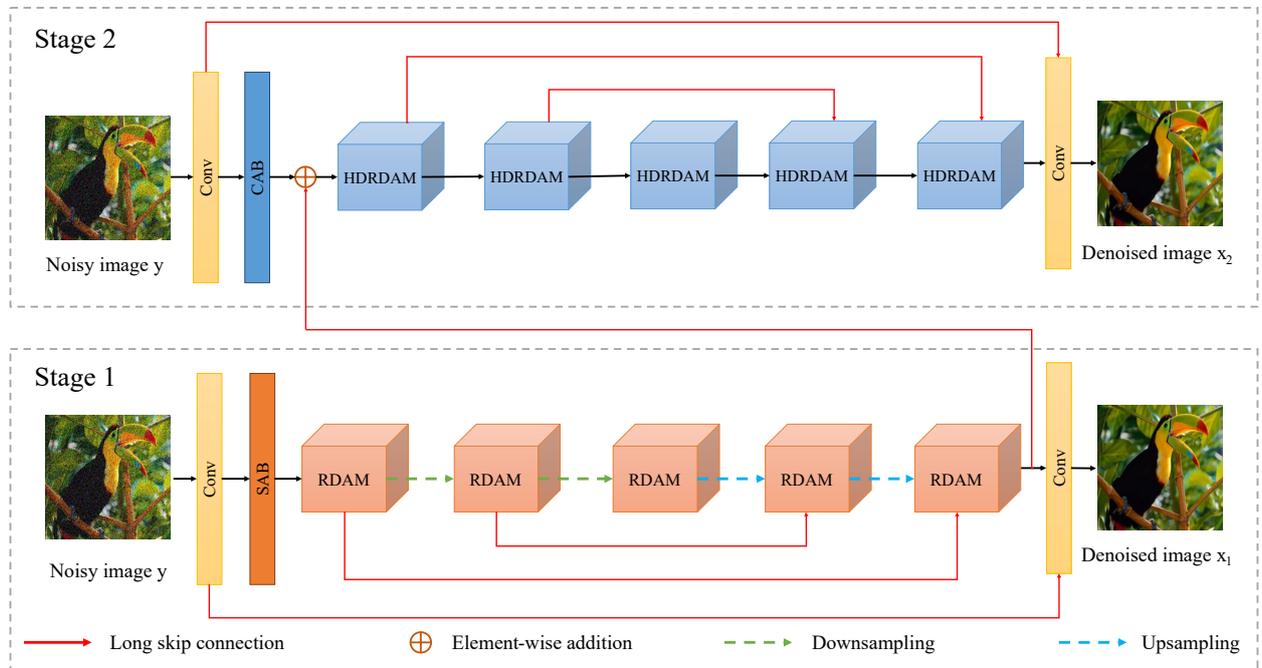


Figure 1: The architecture of the proposed TSP-RDANet.

The whole denoising process of the TSP-RDANet model is represented as follows:

$$x_1, x_2 = \mathcal{F}_{TSP-RDANet}(y), \quad (1)$$

where x_1 and x_2 denote the denoised image from different stages, and y represents the noisy image. Specifically, the first stage denoising network contains two convolutional layers, a spatial attention block (SAB), five RDAMs, long skip connections, the 2×2 strided convolution (SConv), and 2×2 transposed convolution (TConv). The SConv and TConv are respectively applied for image downsampling and upsampling to augment the receptive field of the denoising network, and the multi-scale features can be extracted as well. The denoising procedure of the first stage can be formulated as follows:

$$\begin{aligned} O_{Conv}^1 &= K * y, \\ O_{SAB} &= SAB(O_{Conv}^1), \\ O_{RDAM}^1 &= RDAM(O_{SAB}), \\ O_{RDAM}^2 &= RDAM(SConv(O_{RDAM}^1)), \\ O_{RDAM}^3 &= RDAM(SConv(O_{RDAM}^2)), \\ O_{RDAM}^4 &= RDAM(TConv(O_{RDAM}^3) + O_{RDAM}^2), \\ O_{RDAM}^5 &= RDAM(TConv(O_{RDAM}^4) + O_{RDAM}^1), \\ x_1 &= K * (O_{RDAM}^5 + O_{Conv}^1), \end{aligned} \quad (2)$$

where O_{Conv}^1 and $*$ denote the output of a convolutional layer and convolution operation respectively, and K is the standard convolutional kernel to expand the number of the feature maps. Moreover, O_{SAB} and O_{RDAM}^j ($j \in \{1, 2, 3, 4, 5\}$) is the output of the SAB and the RDAM, respectively.

The second stage denoising network comprises two convolutional layers, a channel attention block (CAB), five HDRDAMs, and long skip connections. The denoising procedure of Stage 2 can be defined as follows:

$$\begin{aligned} O_{Conv}^2 &= K * y, \\ O_{CAB} &= CAB(O_{Conv}^2), \\ O_{HDRDAM}^1 &= HDRDAM(O_{CAB} + O_{RDAM}^5), \\ O_{HDRDAM}^2 &= HDRDAM(O_{HDRDAM}^1), \\ O_{HDRDAM}^3 &= HDRDAM(O_{HDRDAM}^2), \\ O_{HDRDAM}^4 &= HDRDAM(O_{HDRDAM}^3 + O_{HDRDAM}^2), \\ O_{HDRDAM}^5 &= HDRDAM(O_{HDRDAM}^4 + O_{HDRDAM}^1), \\ x_2 &= K * (O_{HDRDAM}^5 + O_{Conv}^2), \end{aligned} \quad (3)$$

where O_{CAB} and O_{HDRDAM}^i ($i \in \{1, 2, 3, 4, 5\}$) are the output of the CAB and the HDRDAM respectively, and O_{Conv}^2 denotes the output of a convolutional layer. The different hierarchical features are fused by the long skip connections between the proposed attention modules, and the restoration effect therefore can be improved.

3.2. The residual dense attention module

The structure of the residual dense attention module (RDAM) is shown in Fig. 2, which is applied for the first stage of the proposed TSP-RDANet. The RDAM contains the dense block (DB), spatial attention block (SAB) [38], and residual learning (RL) [16]. The feature maps f_{RDAM} processed by the RDAM can be expressed as follows:

$$\begin{aligned} O_{RDAM} &= RDAM(f_{RDAM}), \\ &= SAB(DB(f_{RDAM})) + f_{RDAM}, \end{aligned} \quad (4)$$

where O_{RDAM} is the output of the RDAM.

The DB (dense block) consists of eight standard convolutions (Conv), eight rectified linear units (ReLU) [14], and eight concatenation operations (Concat). Specifically, these convolutional layer extracts rich local and hierarchical features, and these features are non-linearly transformed by using the ReLU, then fused by the Concat. Since all the hierarchical features extracted by the DB are treated equally, the important image features are not paid extra attention, which may result in the degradation of the image restoration performance.

The spatial attention block (SAB) is therefore used here to solve the issue, which can disclose the inter-relationship of spatial features. More importantly, the SAB can give more attention on the essential image features and filter the irrelevant ones. The SAB is composed of the Conv, ‘ \otimes ’, Concat, GAP (Global Average Pooling), ReLU, GMP (Global Max Pooling), and Sigmoid [51] activation. The symbols ‘ \otimes ’ and ‘ \oplus ’ in Fig. 2 denote the element-wise product and the RL, respectively.

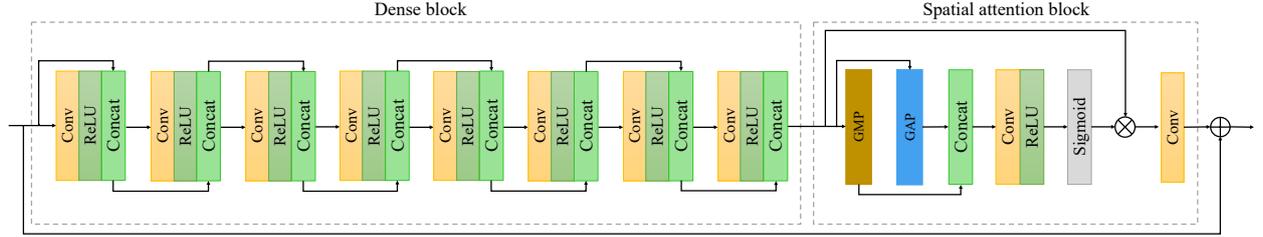


Figure 2: The structure of the designed RDAM.

3.3. The hybrid dilated residual dense attention module

The structure of the designed hybrid dilated residual dense attention module (HDRDAM) is presented in Fig. 3, which is used in the Stage 2 of our denoising network. The HDRDAM includes the hybrid dilated dense block (HDDB), channel attention block (CAB) [39], and residual learning (RL) [16]. The process of the feature maps f_{HDRDAM} passing through the HDRDAM can be formulated as follows:

$$\begin{aligned} O_{HDRDAM} &= HDRDAM(f_{HDRDAM}), \\ &= CAB(HDDB(f_{HDRDAM})) + f_{HDRDAM}, \end{aligned} \quad (5)$$

where O_{HDRDAM} is the output of the HDRDAM.

The HDDB contains eight hybrid dilated convolutions (r -DConv in Fig. 3) [18], eight rectified linear units (ReLU) [14], and eight concatenation operations (Concat), where the ‘ r ’ denotes the dilated rate and the range of its value is [1, 4]. Like the RDAM, the HDRDAM is also able to extract rich local features via dense connection between the dilated convolutional layers. Moreover, the hybrid dilated convolution can not only capture more useful information by enlarging the receptive field, but also can effectively remove the possible gridding phenomenon [52].

The CAB contains the GAP, 1-DConv, ReLU, Sigmoid and ‘ \otimes ’, where the 1-DConv and the symbol ‘ \otimes ’ in Fig. 3 represent the dilated convolution with the dilated rate 1 and the element-wise product, respectively. The CAB is utilized for capturing the inter-dependencies between the spatial features, and it is also used for suppressing the unimportant features. The symbol ‘ \oplus ’ in Fig. 3 denotes the RL.

3.4. Loss function

We utilize different loss functions for synthetic noise removal and real noise elimination. For the Gaussian denoising, we use the loss function L_{mse} that calculates the difference between the recovery image x_i ($i \in \{1, 2\}$) and the ground-truth image x_{gt} , which can be denoted as:

$$\mathcal{L} = \sum_{i=1}^2 L_{mse}(x_i, x_{gt}), \quad (6)$$

where L_{mse} denotes the mean squared error (MSE). In each denoising stage, the L_{MSE} can be designed as:

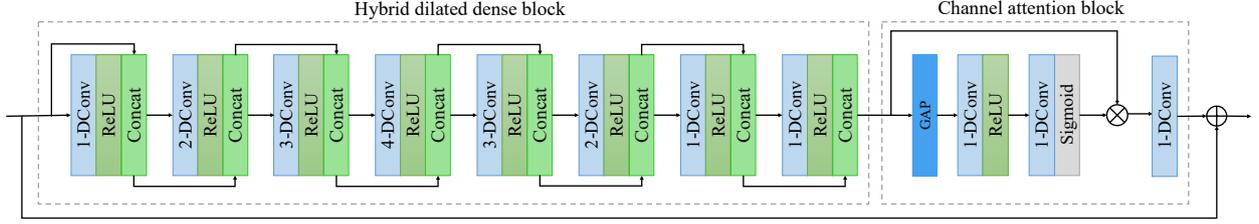


Figure 3: The structure of the designed HDRDAM.

$$\mathcal{L} = \|x_i - x_{gt}\|^2, i \in \{1, 2\} \quad (7)$$

For real noise removal, we refer to [53] to adopt the Charbonnier loss [54] and the edge loss [55] to optimize the TSP-RDANet model, which have been verified can effectively retain image details and edge textures in the recovery image. The whole loss function is expressed as:

$$\mathcal{L} = \sum_{i=1}^2 [L_{char}(x_i, x_{gt}) + \lambda_{edge} L_{edge}(x_i, x_{gt})], \quad (8)$$

where L_{Char} and L_{edge} represent the Charbonnier loss [54] and the edge loss [55], respectively. We also set λ_{edge} to 0.1 according to [40]. For each stage, the Eq. (8) can be further formulated as follows:

$$\mathcal{L} = \sqrt{\|x_i - x_{gt}\|^2 + \epsilon^2} + 0.1 * \sqrt{\|\Delta(x_i) - \Delta(x_{gt})\|^2 + \epsilon^2}, i \in \{1, 2\} \quad (9)$$

where the constant ϵ is equal to 10^{-3} , and the Δ denotes the Laplacian operator [56].

4. Experiments

In this section, we introduce the datasets used in our experiments, the experimental settings, and the experimental results. For results comparison, the qualitative effect of the predicted images generated by different denoising models, the peak signal-to-noise ratio (PSNR), and the structural similarity index measure (SSIM) [57] were utilized. In our comparative experiments, it can be noticed that the methods being compared in each task are not entirely the same, or some results of some methods are missing. This is because some compared methods did not release their code, or they could not be reproduced on our devices.

4.1. Datasets

For the synthetic noise reduction evaluation, the Flick2K dataset [58], which contains 2650 high-resolution color images, was adopted as the training set to train our TSP-RDANet model. To facilitate model training and accelerate the training process, these high-resolution images were randomly cut into image patches with the size of 128×128 . The image patches were grayscaled to train the TSP-RDANet model for single-channel noisy removal evaluation. To produce the clean/noisy image pairs, the additive white Gaussian noise (AWGN) with different noise levels ([0, 50]) is randomly added to the ground-truth image patches. Five public benchmarks were utilized as the test datasets, including the CBS68 [59], Kodak24 [60], McMaster [61], Set12 [59], and BSD68 [59].

For the performance assessment of the real-world image denoising tasks, we adopted the SIDD medium dataset [62] as the training dataset, which comprises 320 pairs of high-resolution noisy images and their near noise-free counterparts. These high-resolution image pairs were also randomly cropped into 128×128 image patches to facilitate the training procedure. The SIDD validation set [62] and DND sRGB dataset [63] were selected as the test sets for real noise removal.

4.2. Experimental settings

All of our experiments were implemented on the Ubuntu 18.04 from a PC equipped with a CPU of Intel(R) Core(TM) i7-11700KF, 32GB RAM, and an Nvidia GeForce RTX 3080Ti GPU. We employed Python 3.8 and Pytorch 1.8 to code the training and testing packages of the TSP-RDANet model. It costed approximately 48 and 50 hours to train the proposed TSP-RDANet for grayscale and color images, respectively. Training the TSP-RDANet model for real noise reduction costed around 180 hours.

The Adam algorithm [64] was utilized to optimize the parameters of the TSP-RDANet. In every training batch, 4 image patches with the size of 128×128 are fed into the TSP-RDANet model for both the synthetic and real image denoising models. For the AWGN noise elimination models, the total number of iterations is 5×10^5 . The learning rate (lr) is firstly set to 1.0×10^{-4} and decreases by half after each 1×10^5 iterations. We run 120 epochs to train the real-world denoising model, and the initial lr is 2.0×10^{-4} . Moreover, we utilized the cosine annealing strategy [65] to gradually reduce the lr to 1.0×10^{-6} .

4.3. Ablation study

In this section, the effectiveness of different stages and different numbers of the RDAM and HDRDAM modules in the proposed TSP-RDANet model are explored. We implemented an ablation experiment on the three denoising models, including: Stage 1 only, Stage 2 only, and the whole TSP-RDANet. Moreover, we also discussed the performance impact of different numbers of RDAM and HDRDAM on the TSP-RDANet model. The average PSNR and SSIM of these denoising models on the BSD68 dataset [59] are used for comparison.

Table 1: The ablation results of different stages on the BSD68 dataset.

Models	Stage 1 only	Stage 2 only	TSP-RDANet (whole)
PSNR	27.78	27.84	27.88
SSIM	0.781	0.783	0.785

In our ablation experiments, the AWGN with noise level 35 was added in the tested images, and the corresponding denoising results of three models are reported in Table 1. It can be seen that compared with Stage 1 only and Stage 2 only, the whole TSP-RDANet model achieves the best performance. In terms of the PSNR values, the denoising result of the TSP-RDANet model surpasses the other two denoising models by 0.10 dB and 0.04 dB. In terms of the SSIM values, the denoising result of the TSP-RDANet model exceeds the rest two denoising models by 0.004 and 0.002. It can be found that by utilizing a progressive strategy, the performance of the TSP-RDANet can be gradually improved.

Table 2: The denoising performances (PSNR/SSIM) using different numbers of the RDAM and HDRDAM on the BSD68 dataset.

Noise levels	TSP-RDANet _{1×1}	TSP-RDANet _{3×3} ¹	TSP-RDANet _{7×7} ³	TSP-RDANet _{5×5} ²
$\sigma=10$	33.36/0.914	33.75/0.925	33.82/0.926	33.80/0.927
$\sigma=20$	29.91/0.840	30.21/0.854	30.29/0.857	30.28/0.858
$\sigma=30$	28.07/0.783	28.37/0.799	28.45/0.803	28.44/0.804
$\sigma=40$	26.84/0.737	27.15/0.756	27.24/0.760	27.23/0.761
$\sigma=50$	25.91/0.698	26.26/0.720	26.35/0.725	26.34/0.726

The denoising results of the TSP-RDANet model equipped with different numbers of the RDAM and HDRDAM modules were also evaluated. Table 2 shows the comparison of average PSNR and SSIM values at different noise levels, where k in TSP-RDANet _{$k \times k$} ^{m} is the numbers of the RDAM and HDRDAM modules in two stages, and m is the number of the downsampling and upsampling operations in Stage 1. In Table 2, one can see that compared with the TSP-RDANet_{1×1} and TSP-RDANet_{3×3}¹, the TSP-RDANet_{5×5}² achieved the best denoising effects on PSNR and SSIM values. Although the TSP-RDANet_{5×5}² is slightly lower than the TSP-RDANet_{7×7}³ on average PSNR value, its performance on average SSIM outperforms the TSP-RDANet_{7×7}³. After full consideration of model performance and complexity, we chose the TSP-RDANet_{5×5}² as our ‘baseline’.

4.4. Synthetic noise removal evaluation

In this subsection, we give the experimental results on synthetic noisy images. The Set12 and BSD68 datasets were applied for the grayscale image denoising test. The CBS68, Kodak24, and McMaster datasets were employed for color image noise reduction experiments. For both the grayscale and color images, the noisy images were produced by adding the AWGN with noise levels of 15, 25, and 50 respectively into the ground-truth images. In Tables 3-5, the red and blue numbers denote the top two denoising results, respectively.

The proposed TSP-RDANet was compared with the classical and the state-of-the-art denoising models, including the BM3D [1], IRCNN [17], DnCNN [5], FFDNet [19], BUIFD [22], DudeNet [34], MWDCNN [44], ADNet [9], CDNet [29], DSNetB [6], AINDNet [21], RIDNet [10], VDN [20], BRDNet [33], and AirNet [30].

4.4.1. Grayscale image denoising evaluation

The PSNR values of different denoising models on the Set12 dataset are displayed in Table 3. One can find that compared with other models, our TSP-RDANet obtains the leading average PSNR results at all the compared noise levels, especially at the higher noise levels. The results reveal that our model is more powerful on discriminating between normal and noisy signals, which benefits from the utilization of the progressive scheme and attention-guided feature filtering.

Table 3: The PSNR comparisons of multiple denoising models on the Set12 dataset.

Noise levels	Models	C.man	House	Peppers	Starfish	Monarch	Airplane	Parrot	Lena	Barbara	Boat	Man	Couple	Average
$\sigma=15$	BM3D [1]	31.91	34.93	32.69	31.14	31.85	31.07	31.37	34.26	33.10	32.13	31.92	31.10	32.37
	DnCNN-S [5]	32.61	34.97	33.30	32.20	33.09	31.70	31.83	34.62	32.64	32.42	32.46	32.47	32.86
	IRCNN [17]	32.55	34.89	33.31	32.02	32.82	31.70	31.84	34.53	32.43	32.34	32.40	32.40	32.77
	FFDNet [19]	32.43	35.07	33.25	31.99	32.66	31.57	31.81	34.62	32.54	32.38	32.41	32.46	32.77
	BUIFD [22]	31.74	34.78	32.80	31.92	32.77	31.34	31.39	34.38	31.68	32.18	32.25	32.22	32.46
	DudeNet [34]	32.71	35.13	33.38	32.29	33.28	31.78	31.93	34.66	32.73	32.46	32.46	32.49	32.94
	MWDCNN [44]	32.53	35.09	33.29	32.28	33.20	31.74	31.97	34.64	32.65	32.49	32.46	32.52	32.91
	TSP-RDANet	32.55	35.35	33.28	32.18	33.30	31.72	31.94	34.80	32.76	32.57	32.49	32.62	32.96
$\sigma=25$	BM3D [1]	29.45	32.85	30.16	28.56	29.25	28.42	28.93	32.07	30.71	29.90	29.61	29.71	29.97
	DnCNN-S [5]	30.18	33.06	30.87	29.41	30.28	29.13	29.43	32.44	30.00	30.21	30.10	30.12	30.43
	IRCNN [17]	30.08	33.06	30.88	29.27	30.09	29.12	29.47	32.43	29.92	30.17	30.04	30.08	30.38
	FFDNet [19]	30.10	33.28	30.93	29.32	30.08	29.04	29.44	32.57	30.01	30.25	30.11	30.20	30.44
	BUIFD [22]	29.42	33.03	30.48	29.21	30.20	28.99	28.94	32.20	29.18	29.97	29.88	29.90	30.12
	DudeNet [34]	30.23	33.24	30.98	29.53	30.44	29.14	29.48	32.52	30.15	30.24	30.08	30.15	30.52
	MWDCNN [44]	30.19	33.33	30.85	29.66	30.55	29.16	29.48	32.67	30.21	30.28	30.10	30.13	30.55
	TSP-RDANet	30.28	33.63	30.95	29.49	30.58	29.16	29.53	32.80	30.36	30.44	30.17	30.37	30.65
$\sigma=50$	BM3D [1]	26.13	29.69	26.68	25.04	25.82	25.10	25.90	29.05	27.22	26.78	26.81	26.46	26.72
	DnCNN-S [5]	27.03	30.00	27.32	25.70	26.78	25.87	26.48	29.39	26.22	27.20	27.24	26.90	27.18
	IRCNN [17]	26.88	29.96	27.33	25.57	26.61	25.89	26.55	29.40	26.24	27.17	27.17	26.88	27.14
	FFDNet [19]	27.05	30.37	27.54	25.75	26.81	25.89	26.57	29.66	26.45	27.33	27.29	27.08	27.32
	BUIFD [22]	25.44	29.76	26.50	24.87	26.49	25.34	25.07	28.81	25.49	26.59	26.87	26.34	26.46
	DudeNet [34]	27.22	30.27	27.51	25.88	26.93	25.88	26.50	29.45	26.49	27.26	27.19	26.97	27.30
	MWDCNN [44]	26.99	30.58	27.34	25.85	27.02	25.93	26.48	29.63	26.60	27.23	27.27	27.11	27.34
	TSP-RDANet	27.53	30.89	27.59	25.96	27.06	25.97	26.67	29.91	27.11	27.52	27.33	27.31	27.57

It is worth noting from Table 3 that the BM3D is superior to our proposed TSP-RDANet on the ‘‘Barbara’’ image (Fig. 4 (v)) at noise levels 15, 25, and 50, which is we think is due to the fact that the ‘Barbara’ image contains rich repetitive structures, that can be effectively learned by the non-local self-similarity (NSS) based method. Furthermore, the denoising effect of the DudeNet is better than that of the TSP-RDANet on the ‘‘C.man’’, ‘‘Peppers’’, ‘‘Starfish’’, and ‘‘Airplane’’ images (Fig. 4 (i)-(iv)) at noise level 15. It can be found from Fig. 4 (i)-(iv) that these images contain many low-frequency regions, and their texture structures are not severely damaged under weak noise, which are beneficial to the DudeNet using the dual-branch structure to predict the noise map and obtain high-quality denoised image. Additionally, the proposed TSP-RDANet model has a larger receptive field, which may make the model insensitive to the weak noise in the low-frequency image areas.

Table 4 shows the averaged SSIM results of the compared methods on the Set12 dataset, and the noisy images at the noise levels of 15, 25, and 50 were used for evaluation. The proposed TSP-RDANet shares the leading position with the ADNet at the noise level 15. It can also be discovered that our model surpasses other compared methods at the more challenging higher noise levels, which again verifies the discrimination ability of our model.

The BSD68 dataset [59] was also used to evaluate grayscale image denoising performance at noise levels 15, 25, and 50. Table 5 reports the averaged PSNR and SSIM values of the compared denoising methods. It can be seen that the proposed TSP-RDANet achieves the best performance at noise levels 25 and 50 than other denoising methods, and obtains competitive results at noise level 15.

The visual comparison was also implemented in our experiments. Fig. 5 shows the comparison between different denoising methods on synthetic noise removal of a grayscale image, where the image used for evaluation is the

Table 4: The SSIM comparisons of multiple denoising models on the Set12 dataset.

Models	BM3D [1]	IRCNN [17]	DnCNN-S [5]	FFDNet [19]	BUIFD [22]	ADNet [9]	CDNet [29]	TSP-RDANet
$\sigma=15$	0.896	0.901	0.903	0.903	0.899	0.905	0.903	0.905
$\sigma=25$	0.851	0.860	0.862	0.864	0.855	0.865	0.865	0.867
$\sigma=50$	0.766	0.780	0.783	0.791	0.755	0.791	0.792	0.798

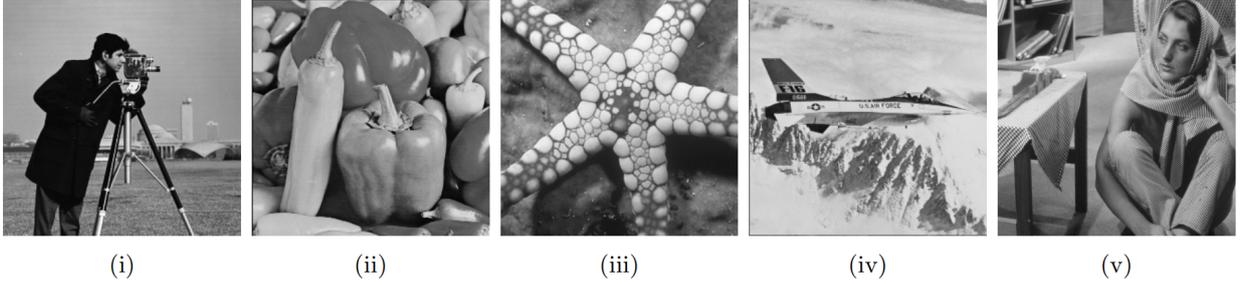


Figure 4: Five grayscale images from the Set12 dataset. (i) “C.man” image, (ii) “Peppers” image, (iii) “Starfish” image, (iv) “Airplane” image, (v) “Barbara” image.

Table 5: The PSNR and SSIM comparisons of image denoising models on the BSD68 dataset.

Metrics	Models	$\sigma=15$	$\sigma=25$	$\sigma=50$
PSNR	BM3D [1]	31.07	28.57	25.62
	DnCNN-S [5]	31.72	29.23	26.23
	IRCNN [17]	31.63	29.15	26.19
	FFDNet [19]	31.63	29.19	26.29
	BUIFD [22]	31.35	28.75	25.11
	DSNetB [6]	31.69	29.22	26.29
	ADNet [9]	31.74	29.25	26.29
	AINDNet [21]	31.69	29.26	26.32
	CDNet [29]	31.74	29.28	26.36
	DudeNet [34]	31.78	29.29	26.31
	MWDCNN [44]	31.77	29.28	26.29
	TSP-RDANet	31.77	29.34	26.45
SSIM	BM3D [1]	0.872	0.802	0.687
	DnCNN-S [5]	0.891	0.828	0.719
	IRCNN [17]	0.888	0.825	0.717
	FFDNet [19]	0.890	0.830	0.726
	BUIFD [22]	0.886	0.819	0.682
	ADNet [9]	0.892	0.829	0.722
	CDNet [29]	0.892	0.831	0.727
	TSP-RDANet	0.892	0.832	0.731

“test004” from the BSD68 dataset. The image was contaminated by the AWGN with a standard deviation of 50. We zoom in a region (green box) for visual detail comparison (red box). One can find that the TSP-RDANet model can remove noise more effectively, meanwhile the model can keep more details of the image. It can be verified that the proposed TSP-RDANet can generate robust denoising performance in grayscale image denoising, both subjectively and objectively.

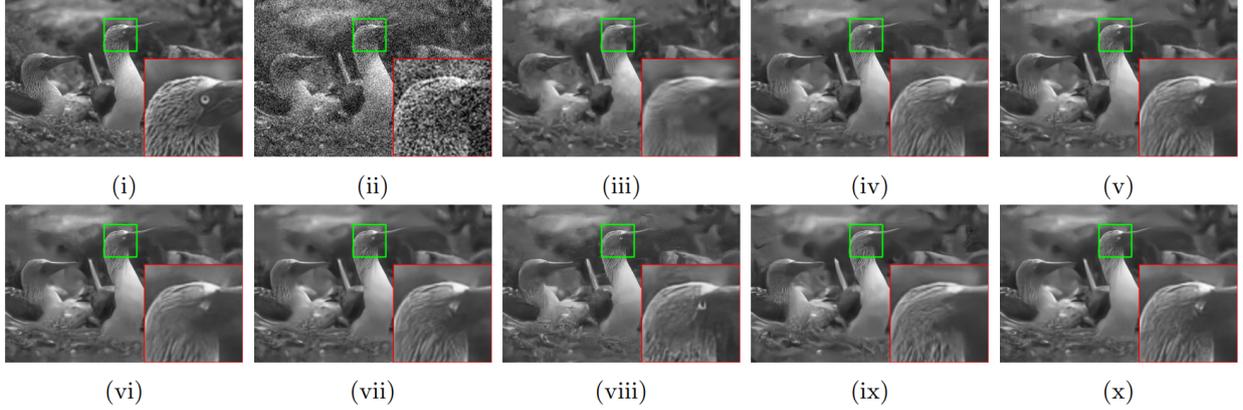


Figure 5: Qualitative comparison between different methods on the image “test004”. (i) Original / PSNR (dB), (ii) Noisy / 14.15, (iii) BM3D / 27.27, (iv) IRCNN / 27.51, (v) DnCNN-S / 27.60, (vi) DnCNN-B / 27.57, (vii) FFDNet / 27.70, (viii) BUIFD / 27.17, (ix) ADNet / 27.64, (x) TSP-RDANet / 27.70.

4.4.2. Color image denoising evaluation

For the synthetic noise removal evaluation on color images, we adopted three public and commonly used datasets, including the CBS68 [59], Kodak24 [60], and McMaster [61] datasets. These datasets were polluted by the AWGN with standard deviations of 15, 25, and 50. Table 6, Table 7, and Table 8 list the averaged PSNR and SSIM results of our TSP-RDANet and the compared denoising methods, and the red and blue numbers are the best and second-best noise reduction results, respectively. Similar with the performance on grayscale images, the proposed model can obtain competitive results at noise level 15. At the noise level of 25 and 50, our model achieves the highest denoising performance.

The qualitative comparison was also implemented for color image denoising. Fig. 6 and Fig. 7 display the visual assessment of our TSP-RDANet and other models on synthetic color image denoising. The results on the “163085” image from the CBS68 dataset and “kodim21” image from the Kodak24 dataset are presented. These images were contaminated by the AWGN with the noise level of 50. We zoom in an image region (red box) for a more detailed comparison (green box). One can find from Fig. 6 and Fig. 7 that the TSP-RDANet obtains a desirable trade-off between noise reduction and image texture retaining.

It should be noted from Tables 3-8 that with the increase of noise level, the denoising result of the TSP-RDANet is better than other compared models. This is due to the fact that the designed RDAM and HDRDAM can greatly increase the receptive field of the TSP-RDANet, which helps it extract more contextual information and cope with stronger noise [5, 19].

4.5. Real-world noise reduction evaluation

For real-world noise reduction assessment, we utilized the SIDD validation set [62] and the DND sRGB images [63] as the test datasets. The classical and the state-of-the-art blind denoising models, such as the TWSC [2], MCWNNM [4], DnCNN-B [5], RIDNet [10], CBDNet [23], VDN [20], AINDNet [21], DANet+ [35], DeamNet [37], VDIR [24], SADNet [26], and CycleISP [11], were applied for comparing. Table 9 lists the averaged PSNR and SSIM values of our TSP-RDANet and other denoising methods, and the best and the second performances are emphasized in red and blue, respectively. One can find that the TSP-RDANet obtains the leading results compared

Table 6: The PSNR and SSIM comparisons of multiple denoising models on the CBSD68 dataset.

Metrics	Models	$\sigma=15$	$\sigma=25$	$\sigma=50$
PSNR	CBM3D [1]	33.52	30.71	27.38
	IRCNN [17]	33.86	31.16	27.86
	CDnCNN-S [5]	33.89	31.23	27.92
	FFDNet [19]	33.87	31.21	27.96
	BUIFD [22]	33.65	30.76	26.61
	DSNetB [6]	33.91	31.28	28.05
	VDN [20]	33.90	31.35	28.19
	RIDNet [10]	34.01	31.37	28.14
	ADNet [9]	33.99	31.31	28.04
	BRDNet [33]	34.10	31.43	28.16
	DudeNet [34]	34.01	31.34	28.09
	AirNet [30]	33.92	31.26	28.01
	MWDCNN [44]	34.18	31.45	28.13
	TSP-RDANet	34.14	31.52	28.33
SSIM	IRCNN [17]	0.929	0.882	0.790
	CDnCNN-B [5]	0.929	0.883	0.790
	FFDNet [19]	0.929	0.882	0.789
	ADNet [9]	0.933	0.889	0.797
	BUIFD [22]	0.930	0.882	0.777
	AirNet [30]	0.933	0.888	0.798
	TSP-RDANet	0.932	0.889	0.804

Table 7: The PSNR and SSIM comparisons of different models on the Kodak24 dataset.

Metrics	Models	$\sigma=15$	$\sigma=25$	$\sigma=50$
PSNR	CBM3D [1]	34.28	31.68	28.46
	IRCNN [17]	34.56	32.03	28.81
	CDnCNN-S [5]	34.48	32.03	28.85
	FFDNet [19]	34.63	32.13	28.98
	BUIFD [22]	34.41	31.77	27.74
	DSNetB [6]	34.63	32.16	29.05
	ADNet [9]	34.76	32.26	29.10
	BRDNet [33]	34.88	32.41	29.22
	DudeNet [34]	34.81	32.26	29.10
	AirNet [30]	34.68	32.21	29.06
	MWDCNN [44]	34.91	32.40	29.26
TSP-RDANet	34.98	32.54	29.46	
SSIM	IRCNN [17]	0.920	0.877	0.793
	CDnCNN-B [5]	0.920	0.876	0.791
	FFDNet [19]	0.922	0.878	0.794
	ADNet [9]	0.924	0.882	0.798
	BUIFD [22]	0.923	0.879	0.786
	AirNet [30]	0.924	0.882	0.799
MWDCNN [44]	0.927	0.886	0.806	
TSP-RDANet	0.926	0.887	0.812	

Table 8: The PSNR and SSIM comparisons of image denoising models on the McMaster dataset.

Metrics	Models	$\sigma=15$	$\sigma=25$	$\sigma=50$
PSNR	CBM3D [1]	34.06	31.66	28.51
	IRCNN [17]	34.58	32.18	28.91
	CDnCNN-S [5]	33.44	31.51	28.61
	FFDNet [19]	34.66	32.35	29.18
	BUIFD [22]	33.84	31.06	26.20
	DSNetB [6]	34.67	32.40	29.28
	ADNet [9]	34.93	32.56	29.36
	BRDNet [33]	35.08	32.75	29.52
	AirNet [30]	34.70	32.44	29.26
	TSP-RDANet	35.06	32.81	29.74
SSIM	IRCNN [17]	0.920	0.882	0.807
	CDnCNN-B [5]	0.904	0.869	0.799
	FFDNet [19]	0.922	0.886	0.815
	ADNet [9]	0.927	0.894	0.825
	BUIFD [22]	0.901	0.847	0.733
	AirNet [30]	0.925	0.891	0.822
	TSP-RDANet	0.928	0.897	0.836

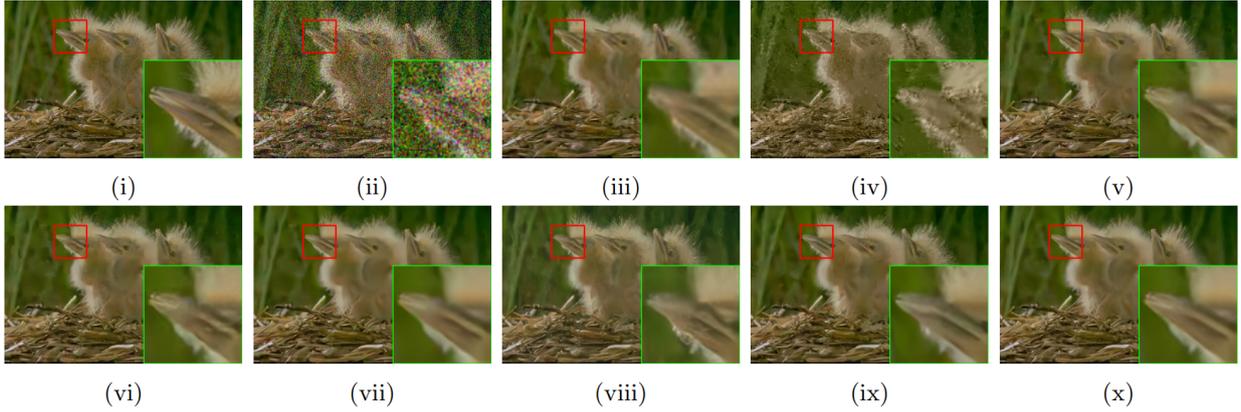


Figure 6: Qualitative comparison on the image “163085”. (i) Original / PSNR (dB), (ii) Noisy / 14.15, (iii) CBM3D / 28.34, (iv) MCWNNM / 23.85, (v) IRCNN / 28.69, (vi) CDnCNN-B / 28.68, (vii) FFDNet / 28.75, (viii) BUIFD / 26.88, (ix) ADNet / 28.80, (x) TSP-RDANet / 29.01.

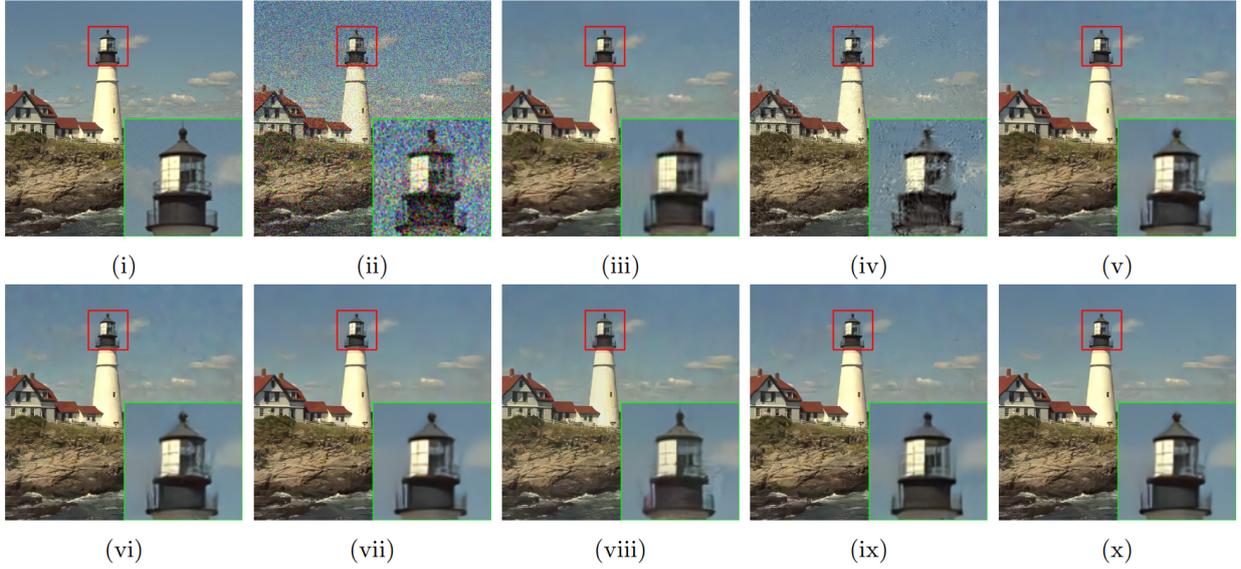


Figure 7: Qualitative comparison on the image “kodim21”. (i) Original / PSNR (dB), (ii) Noisy / 14.15, (iii) CBM3D / 27.44, (iv) MCWNNM / 23.35, (v) IRCNN / 27.97, (vi) CDnCNN-B / 28.03, (vii) FFDNet / 28.07, (viii) BUIFD / 27.61, (ix) ADNet / 28.15, (x) TSP-RDANet / 28.28.

with other methods. On the SIDD validation set, the proposed TSP-RDANet achieves 0.30 dB, 0.63 dB, 0.15 dB, 0.23 dB, 0.32 dB, 0.12 dB, and 0.06 dB performance improvement on the state-of-the-art methods VDN, AINDNet, DANet+, DeamNet, VDIR, SADNet, and CycleISP, respectively.

Table 9: The performance (PSNR/SSIM) comparison between different models on the SIDD and DND datasets.

Models	SIDD	DND
TWSC [2]	35.33/0.933	37.94/0.940
MCWNNM [4]	33.40/0.879	37.38/0.929
DnCNN-B [5]	23.66/0.583	37.90/0.943
RIDNet [10]	38.71/0.954	39.23/0.953
CBDNet [23]	30.78/0.951	38.06/0.942
VDN [20]	39.28/0.956	39.38/0.952
AINDNet [21]	38.95/0.952	39.37/0.951
DANet+ [35]	39.43/0.956	39.58/0.955
DeamNet [37]	39.35/0.955	39.63/0.953
VDIR [24]	39.26/0.955	39.63/0.953
SADNet [26]	39.46/0.957	39.59/0.952
CycleISP [11]	39.52/0.957	39.56/0.956
TSP-RDANet	39.58/0.958	39.70/0.954

Fig. 8 shows the visual denoising results of the proposed TSP-RDANet and other compared models for real image denoising. The image “10_3” from the SIDD was used for visual presentation. An image region (green box) is enlarged for a more detailed visual comparison (red box). As can be seen from Fig. 8, the TSP-RDANet obtains a more visually appealing result. Moreover, the TSP-RDANet also surpasses all other denoising models on PSNR value.

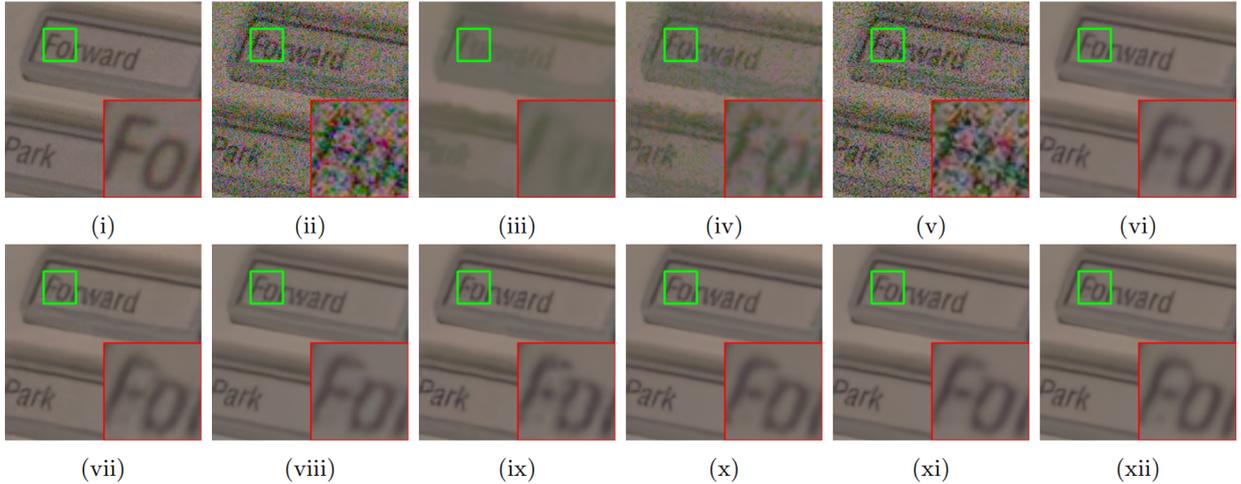


Figure 8: Qualitative comparison of different denoising models on the image “10.3”. (i) Original / PSNR (dB), (ii) Noisy / 18.25, (iii) TWSC / 30.42, (iv) MCWNNM / 28.63, (v) CDnCNN-B / 20.76, (vi) VDN / 36.39, (vii) AINDNet / 36.24, (viii) DANet+ / 36.37, (ix) VDIR / 36.35, (x) SADNet / 36.72, (xi) CycleISP / 36.72, (xii) TSP-RDANet / 36.84.

4.6. Model computational complexity evaluation

The model complexity of all the compared models was also evaluated in our experiments, which was quantified by the running time and network parameters. The BM3D, TWSC, and MCWNNM were performed in the Matlab (R2020a) platform, and the remaining methods were executed in the PyCharm (2021) software. We randomly selected three clean color images with the sizes of 256×256 , 512×512 , and 1024×1024 to obtain the runtime of every model on each image. We performed single-channel processing of these color images to obtain their grayscale counterparts. The clean grayscale and color images were then added to the AWGN at noise level 25 to generate the noisy images for evaluation. The public pytorch-OpCounter package¹ was used to obtain the model parameters of the evaluated denoising models. In order to make a more precise and objective comparison, we ran 20 executions of each denoising model to obtain the average runtime, which was used for comparing the complexity of different denoising models.

Table 10 reports the runtime of all tested models. It can be seen that the running speeds of our TSP-RDANet outperform the BM3D, TWSC, MCWNNM, AINDNet, and VDIR at all image sizes. Although the denoising speed of the proposed TSP-RDANet is slower than the IRCNN, DnCNN-B, FFDNet, BUlFD, DANet+, CycleISP, ADNet, and SADNet, the denoising effect of our model outperforms these models. The numbers of the network parameters of different models were offered in Table 11. One can find that the proposed TSP-RDANet owns fewer network parameters than many denoising methods such as the DANet+, VDN, SADNet, AINDNet, CBDNet, and AirNet. Meanwhile, our model can achieve better quantitative and qualitative denoising results than these models (see Sec. 4.4 and Sec. 4.5). It can be verified from our experiments that our model obtains a favorable balance between denoising result and complexity.

5. Conclusion

In this work, we introduce a new two-stage progressive network (TSP-RDANet) for image noise removal. The proposed model decomposes the entire denoising process into two sub-stage tasks to progressively obtain more effective noise reduction. Especially, we design the novel residual dense attention module (RDAM) for the first stage of the model. The hybrid dilated residual dense attention module (HDRDAM) is developed for the second stage. The RDM and HDRDAM consist of the dense block, attention block, and residual learning, which makes them be capable of

¹<https://github.com/Lyken17/pytorch-OpCounter>

Table 10: Runtime results (in seconds) of the compared models on different synthetic grayscale and color noisy images.

Devices	Models	256 × 256		512 × 512		1024 × 1024	
		Grayscale	Color	Grayscale	Color	Grayscale	Color
CPU	BM3D [1]	0.458	0.593	2.354	3.771	9.782	12.818
	TWSC [2]	12.314	34.41	53.155	140.964	221.507	608.492
	MCWNNM [4]	-	62.777	-	277.623	-	1120.112
GPU	IRCNN [17]	0.030	0.030	0.030	0.030	0.030	0.030
	DnCNN-B [5]	0.032	0.032	0.037	0.037	0.057	0.057
	FFDNet [19]	0.031	0.030	0.031	0.030	0.032	0.030
	BUIFD [22]	0.035	0.037	0.050	0.053	0.112	0.123
	DANet+ [35]	-	0.025	-	0.027	-	0.041
	ADNet [9]	0.031	0.033	0.035	0.045	0.051	0.093
	SADNet [26]	0.030	0.030	0.043	0.044	0.101	0.102
	CycleISP [11]	-	0.055	-	0.156	-	0.545
	AirNet [30]	-	0.143	-	0.498	-	2.501
	VDN [20]	0.144	0.162	0.607	0.597	2.367	2.376
	AINDNet [21]	-	0.531	-	2.329	-	9.573
	VDIR [24]	-	0.385	-	1.622	-	6.690
	TSP-RDANet	0.261	0.263	0.625	0.627	1.979	2.004

Table 11: Number of parameters (in K) of the compared image denoising models.

Models	Number of model parameters	
	Grayscale	Color
IRCNN [17]	186	188
DnCNN-B [5]	668	673
FFDNet [19]	485	852
BUIFD [22]	1075	1085
RIDNet [10]	1497	1499
BRDNet [33]	1113	1117
ADNet [9]	519	521
DudeNet [34]	1077	1079
DeamNet [37]	1873	1876
VDIR [24]	-	2227
CycleISP [11]	-	2837
CBDNet [23]	-	4365
VDN [20]	7810	7817
DANet+ [35]	-	9154
AINDNet [21]	-	13764
SADNet [26]	3450	3451
AirNet [30]	-	8930
TSP-RDANet	2846	2848

extracting rich local features and filtering irrelevant and redundant features to further improve their learning ability. The long skip connection is utilized to link the two sub-networks to achieve progressive image denoising. Downsampling and the dilated convolution are applied to the TSP-RDANet model to enlarge its receptive field size and further capture more contextual information to enhance its restoration capability. Compared with the state-of-the-art models, our TSP-RDANet achieves excellent and competitive denoising results on multiple denoising tasks.

The experimental results yet also reveal that the discriminative ability of the proposed model under weak noise levels still needs to be further investigated. Moreover, although the proposed model can produce competitive blind denoising performance, the model needs clean/noisy image pairs to train. In our following work, we will focus on evolving the model into a self-supervised or unsupervised manner.

6. Acknowledgements

This work is supported by the Natural Science Foundation of China (61863037, 41971392) and the Applied Basic Research Foundation of Yunnan Province (202001AT070077).

References

- [1] K. Dabov, A. Foi, V. Katkovnik, K. O. Egiazarian, Image denoising by sparse 3-d transform-domain collaborative filtering, *IEEE Trans. Image Process.* 16 (8) (2007) 2080–2095.
- [2] J. Xu, L. Zhang, D. Zhang, A trilateral weighted sparse coding scheme for real-world image denoising, *European Conference on Computer Vision* 11212 (2018) 21–38.
- [3] Y. Ou, M. N. S. Swamy, J. Luo, B. Li, Single image denoising via multi-scale weighted group sparse coding, *Signal Process.* 200 (2022) 108650.
- [4] J. Xu, L. Zhang, D. Zhang, X. Feng, Multi-channel weighted nuclear norm minimization for real color image denoising, *IEEE International Conference on Computer Vision* (2017) 1105–1113.
- [5] K. Zhang, W. Zuo, Y. Chen, D. Meng, L. Zhang, Beyond a gaussian denoiser: Residual learning of deep CNN for image denoising, *IEEE Trans. Image Process.* 26 (7) (2017) 3142–3155.
- [6] Y. Peng, L. Zhang, S. Liu, X. Wu, Y. Zhang, X. Wang, Dilated residual networks with symmetric skip connection for image denoising, *Neurocomputing* 345 (2019) 67–76.
- [7] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, Y. Fu, Residual dense network for image restoration, *IEEE Trans. Pattern Anal. Mach. Intell.* 43 (7) (2021) 2480–2495.
- [8] F. Jia, W. H. Wong, T. Zeng, Ddunet: Dense dense u-net with applications in image denoising, in: *International Conference on Computer Vision Workshops*, 2021, pp. 354–364.
- [9] C. Tian, Y. Xu, Z. Li, W. Zuo, L. Fei, H. Liu, Attention-guided CNN for image denoising, *Neural Networks* 124 (2020) 117–129.
- [10] S. Anwar, N. Barnes, Real image denoising with feature attention, *International Conference on Computer Vision* (2019) 3155–3164.
- [11] S. W. Zamir, A. Arora, S. H. Khan, M. Hayat, F. S. Khan, M. Yang, L. Shao, Cycleisp: Real image restoration via improved data synthesis, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2693–2702.
- [12] Y. Bai, M. Liu, C. Yao, C. Lin, Y. Zhao, Mspnet: Multi-stage progressive network for image denoising, *Neurocomputing* 517 (2023) 71–80.
- [13] Z. Wang, Y. Bai, Y. Zhou, C. Xie, Can cnns be more robust than transformers?, in: *International Conference on Learning Representations*, 2023.
- [14] A. Krizhevsky, I. Sutskever, G. E. Hinton, Imagenet classification with deep convolutional neural networks, *Advances in Neural Information Processing Systems* (2012) 1106–1114.
- [15] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, *International Conference on Machine Learning* 37 (2015) 448–456.
- [16] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, *IEEE Conference on Computer Vision and Pattern Recognition* (2016) 770–778.
- [17] K. Zhang, W. Zuo, S. Gu, L. Zhang, Learning deep CNN denoiser prior for image restoration, *IEEE Conference on Computer Vision and Pattern Recognition* (2017) 2808–2817.
- [18] F. Yu, V. Koltun, Multi-scale context aggregation by dilated convolutions, in: *International Conference on Learning Representations*, 2016.
- [19] K. Zhang, W. Zuo, L. Zhang, Ffdnet: Toward a fast and flexible solution for cnn-based image denoising, *IEEE Trans. Image Process.* 27 (9) (2018) 4608–4622.
- [20] Z. Yue, H. Yong, Q. Zhao, D. Meng, L. Zhang, Variational denoising network: Toward blind noise modeling and removal, *Advances in Neural Information Processing Systems* (2019) 1688–1699.
- [21] Y. Kim, J. W. Soh, G. Y. Park, N. I. Cho, Transfer learning from synthetic to real-noise denoising with adaptive instance normalization, *IEEE Conference on Computer Vision and Pattern Recognition* (2020) 3479–3489.
- [22] M. E. Helou, S. Süsstrunk, Blind universal bayesian image denoising with gaussian noise level learning, *IEEE Trans. Image Process.* 29 (2020) 4885–4897.
- [23] S. Guo, Z. Yan, K. Zhang, W. Zuo, L. Zhang, Toward convolutional blind denoising of real photographs, *IEEE Conference on Computer Vision and Pattern Recognition* (2019) 1712–1722.
- [24] J. W. Soh, N. I. Cho, Variational deep image restoration, *IEEE Trans. Image Process.* 31 (2022) 4363–4376.

- [25] W. Wu, S. Liao, G. Lv, P. Liang, Y. Zhang, Image blind denoising using dual convolutional neural network with skip connection (2023). [arXiv:2304.01620](https://arxiv.org/abs/2304.01620).
- [26] M. Chang, Q. Li, H. Feng, Z. Xu, Spatial-adaptive network for single image denoising, in: *European Conference on Computer Vision*, Vol. 12375, 2020, pp. 171–187.
- [27] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, Y. Wei, Deformable convolutional networks, in: *IEEE International Conference on Computer Vision*, 2017, pp. 764–773.
- [28] X. Zhu, H. Hu, S. Lin, J. Dai, Deformable convnets V2: more deformable, better results, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 9308–9316.
- [29] Y. Quan, Y. Chen, Y. Shao, H. Teng, Y. Xu, H. Ji, Image denoising using complex-valued deep CNN, *Pattern Recognit.* 111 (2021) 107639.
- [30] B. Li, X. Liu, P. Hu, Z. Wu, J. Lv, X. Peng, All-in-one image restoration for unknown corruption, *IEEE Conference on Computer Vision and Pattern Recognition* (2022) 17431–17441.
- [31] J. Pan, D. Sun, J. Zhang, J. Tang, J. Yang, Y. Tai, M. Yang, Dual convolutional neural networks for low-level vision, *Int. J. Comput. Vis.* 130 (6) (2022) 1440–1458.
- [32] S. Ioffe, Batch renormalization: Towards reducing minibatch dependence in batch-normalized models, in: *Advances in Neural Information Processing Systems*, 2017, pp. 1945–1953.
- [33] C. Tian, Y. Xu, W. Zuo, Image denoising using deep CNN with batch renormalization, *Neural Networks* 121 (2020) 461–473.
- [34] C. Tian, Y. Xu, W. Zuo, B. Du, C. Lin, D. Zhang, Designing and training of a dual CNN for image denoising, *Knowl. Based Syst.* 226 (2021) 106949.
- [35] Z. Yue, Q. Zhao, L. Zhang, D. Meng, Dual adversarial network: Toward real-world noise removal and noise generation, in: *European Conference on Computer Vision*, Vol. 12355, 2020, pp. 41–58.
- [36] B. Jiang, J. Li, H. Li, R. Li, D. Zhang, G. Lu, Enhanced frequency fusion network with dynamic hash attention for image denoising, *Inf. Fusion* 92 (2023) 420–434.
- [37] C. Ren, X. He, C. Wang, Z. Zhao, Adaptive consistency prior based deep network for image denoising, *IEEE Conference on Computer Vision and Pattern Recognition* (2021) 8596–8606.
- [38] S. Woo, J. Park, J. Lee, I. S. Kweon, CBAM: convolutional block attention module, *European Conference on Computer Vision* 11211 (2018) 3–19.
- [39] J. Hu, L. Shen, G. Sun, Squeeze-and-excitation networks, *IEEE Conference on Computer Vision and Pattern Recognition* (2018) 7132–7141.
- [40] W. Wu, S. Liu, Y. Zhou, Y. Zhang, Y. Xiang, Dual residual attention network for image denoising (2023). [arXiv:2305.04269](https://arxiv.org/abs/2305.04269).
- [41] J. Huang, Z. Zhao, C. Ren, Q. Teng, X. He, A prior-guided deep network for real image denoising and its applications, *Knowl. Based Syst.* 255 (2022) 109776.
- [42] R. Zhuge, J. Wang, Z. Xu, Y. Xu, Single image denoising with a feature-enhanced network, *Neural Networks* 168 (2023) 313–325.
- [43] R. K. Thakur, S. K. Maji, Multi scale pixel attention and feature extraction based neural network for image denoising, *Pattern Recognit.* 141 (2023) 109603.
- [44] C. Tian, M. Zheng, W. Zuo, B. Zhang, Y. Zhang, D. Zhang, Multi-stage image denoising with the wavelet transform, *Pattern Recognit.* 134 (2023) 109050.
- [45] Y. Zheng, X. Yu, M. Liu, S. Zhang, Residual multiscale based single image deraining, in: *British Machine Vision Conference*, 2019, p. 147.
- [46] D. Ren, W. Zuo, Q. Hu, P. Zhu, D. Meng, Progressive image deraining networks: A better and simpler baseline, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3937–3946.
- [47] S. Nah, T. H. Kim, K. M. Lee, Deep multi-scale convolutional neural network for dynamic scene deblurring, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 257–265.
- [48] H. Zhang, Y. Dai, H. Li, P. Koniusz, Deep stacked hierarchical multi-patch network for image deblurring, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 5978–5986.
- [49] M. Suin, K. Purohit, A. N. Rajagopalan, Spatially-attentive patch-hierarchical network for adaptive motion deblurring, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3603–3612.
- [50] X. Fu, B. Liang, Y. Huang, X. Ding, J. W. Paisley, Lightweight pyramid networks for image deraining, *IEEE Trans. Neural Networks Learn. Syst.* 31 (6) (2020) 1794–1807.
- [51] J. Han, C. Moraga, The influence of the sigmoid function parameters on the speed of backpropagation learning, *International Workshop on Artificial Neural Networks* 930 (1995) 195–201.
- [52] P. Wang, P. Chen, Y. Yuan, D. Liu, Z. Huang, X. Hou, G. W. Cottrell, Understanding convolution for semantic segmentation, *IEEE Winter Conference on Applications of Computer Vision* (2018) 1451–1460.
- [53] S. W. Zamir, A. Arora, S. H. Khan, M. Hayat, F. S. Khan, M. Yang, L. Shao, Multi-stage progressive image restoration, in: *IEEE Conference on Computer Vision and Pattern Recognition*, 2021, pp. 14821–14831.
- [54] W. Lai, J. Huang, N. Ahuja, M. Yang, Deep laplacian pyramid networks for fast and accurate super-resolution, *IEEE Conference on Computer Vision and Pattern Recognition* (2017) 5835–5843.
- [55] K. Jiang, Z. Wang, P. Yi, C. Chen, B. Huang, Y. Luo, J. Ma, J. Jiang, Multi-scale progressive fusion network for single image deraining, *IEEE Conference on Computer Vision and Pattern Recognition* (2020) 8343–8352.
- [56] B. Kamgar-Parsi, A. Rosenfeld, Optimally isotropic laplacian operator, *IEEE Trans. Image Process.* 8 (10) (1999) 1467–1472.
- [57] Z. Wang, A. C. Bovik, H. R. Sheikh, E. P. Simoncelli, Image quality assessment: from error visibility to structural similarity, *IEEE Trans. Image Process.* 13 (4) (2004) 600–612.
- [58] B. Lim, S. Son, H. Kim, S. Nah, K. M. Lee, Enhanced deep residual networks for single image super-resolution, in: *IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2017, pp. 1132–1140.
- [59] S. Roth, M. J. Black, Fields of experts: A framework for learning image priors, *IEEE Computer Society Conference on Computer Vision and Pattern* (2005) 860–867.
- [60] R. Franzen, Kodak lossless true color image suite, <http://r0k.us/graphics/kodak/> (1999).
- [61] L. Zhang, X. Wu, A. Buades, X. Li, Color demosaicking by local directional interpolation and nonlocal adaptive thresholding, *Journal of*

Electronic Imaging 20 (2) (2011) 1–17.

- [62] A. Abdelhamed, S. Lin, M. S. Brown, A high-quality denoising dataset for smartphone cameras, IEEE Conference on Computer Vision and Pattern Recognition (2018) 1692–1700.
- [63] T. Plotz, S. Roth, Benchmarking denoising algorithms with real photographs, in: IEEE Conference on Computer Vision and Pattern Recognition, 2017, pp. 1586–1595.
- [64] D. P. Kingma, J. Ba, Adam: A method for stochastic optimization, in: International Conference on Learning Representations, 2015.
- [65] I. Loshchilov, F. Hutter, SGDR: stochastic gradient descent with warm restarts, in: International Conference on Learning Representations, 2017.