

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/337923159>

Self-attentive Pyramid Network for Single Image De-raining

Chapter · December 2019

DOI: 10.1007/978-3-030-36708-4_32

CITATIONS

0

READS

55

4 authors, including:



Tao Dai

Tsinghua University

54 PUBLICATIONS 221 CITATIONS

[SEE PROFILE](#)



Jiawei Li

Tsinghua University

16 PUBLICATIONS 27 CITATIONS

[SEE PROFILE](#)

Some of the authors of this publication are also working on these related projects:



Deep Learning [View project](#)



Ensemble Deep Learning [View project](#)

Self-Attentive Pyramid Network for Single Image De-raining ^{*}

Taian Guo^{1,*}, Tao Dai^{1,*}, Jiawei Li¹, and Shu-Tao Xia²

¹ Tsinghua University, Haidian, Beijing 100084, P. R. China
{gta17, dait14, li-jw15}@mails.tsinghua.edu.cn

² Graduate School at Shenzhen, Tsinghua University, Shenzhen, Guangdong 518055,
P. R. China
xiast@sz.tsinghua.edu.cn

Abstract. Rain Streaks in a single image can severely damage the visual quality, and thus degrade the performance of current computer vision algorithms. To remove the rain streaks effectively, plenty of CNN-based methods have recently been developed, and obtained impressive performance. However, most existing CNN-based methods focus on network design, while rarely exploits spatial correlations of feature. In this paper, we propose a deep self-attentive pyramid network (SAPN) for more powerful feature expression for single image de-raining. Specifically, we propose a self-attentive pyramid module (SAM), which consists of convolutional layers enhanced by self-attention calculation units (SACUs) to capture the abstraction of image contents, and deconvolutional layers to upsample the feature maps and recover image details. Besides, we propose self-attention based skip connections to symmetrically link convolutional and deconvolutional layers to exploit spatial contextual information better. To model rain streaks with various scales and shapes, a multi-scale pooling (MSP) module is also introduced to efficiently leverage features from different scales. Extensive experiments on both synthetic and real-world datasets demonstrate the effectiveness of our proposed method in terms of both quantitative and visual quality.

Keywords: Rain streak removal · Encoder-decoder network · Self-attention.

1 Introduction

Images captured in rain weather are common in real life, thus resulting in images with rain streaks. Such rain streaks would not only affect the visual quality of images, but degrade performance of existing computer vision systems, such as self-driving, video surveillance, and object detection. Therefore, it is of crucial importance to remove rain streaks while recovering image details. Image de-raining has received much attention in recent years, and can be generally divided into video-based [1–4] and single image based methods [5–10]. Most video based methods focus on utilizing the temporal correlations in successive frames,

* * Equal contribution.

2 T. Guo et al.

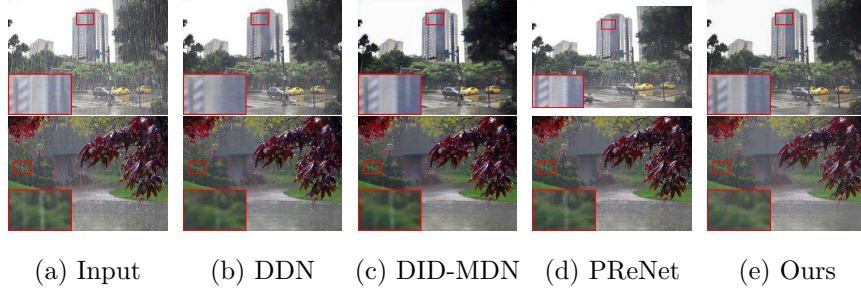


Fig. 1: Sample de-raining results on real-world rainy scenes with long heavy rain streaks. The details in enlarged regions shows that our SAPN removes long heavy rain streaks more cleanly, while keeps the sharp details of the background objects. The two rows demonstrate that our self-attentive network produces better de-raining results on image regions with long heavy rain streaks.

which provide extra temporal information of the rainy scene. In contrast, it is more challenging to perform single image de-raining due to the very limited information from a single image.

In recent years, many single image de-raining methods [5–8, 10] have been proposed. Most traditional image de-raining methods focus on exploiting powerful image prior of rainy images, including sparse prior [7], low rank prior [11] and Gaussian mixture model (GMM) prior [6]. Among them, Luo et al. [5] proposed a dictionary learning based method, which sparsely approximates the patches of the rain layer and the de-rained layer by discriminative sparse codes with a learned dictionary. Li et al. [6] further introduced patch-based Gaussian mixture model (GMM) priors for both the background layer and the rain layer. Zhu et al. [7] introduced three types of priors, and proposed a joint optimization process to alternately remove rain-streak details. However, since such methods rely heavily on handcrafted feature and fixed priors, they are limited in practice due to the diversity of rain streaks (e.g., various shapes, scales and density levels).

Due to the powerful feature representation capability, convolutional neural networks have been widely used in image de-raining, and obtained remarkable performance. For example, Fu et al. [8] proposed a deep detail network to learn the high frequency details during training, since most rain streaks belong to high-frequency information. To consider various shapes and density of rain drops, Zhang et al. [10] proposed a densely connected network with learned rain streak density information to assist rain streak removal process. Since spatial contextual information is important for rain streaks removal, some methods [12, 13] have been developed. Specifically, Li et al. [12] proposed a multi-stage dilated CNN network to obtain a large receptive field size. Recently, Ren et al. [14] proposed a progressive recurrent network (PReNet) to better take advantage of recursive computation and exploit the dependencies of deep features across stages.

Although significant progress has been achieved for single image de-raining, most of existing CNN-based methods focus on the network design, while rarely considering the inherent spatial correlations in feature maps. Meanwhile, self-

Self-Attentive Pyramid Network for Single Image De-raining

3

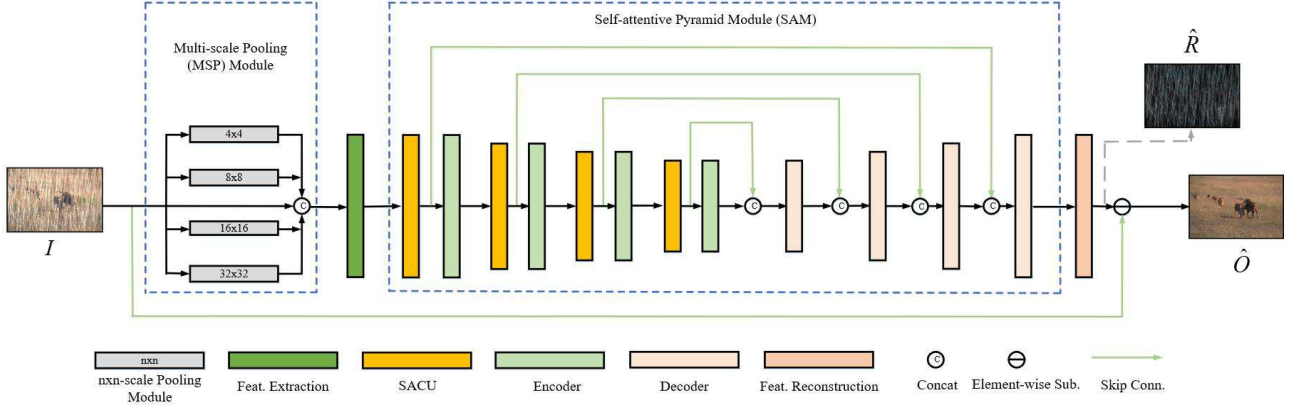


Fig. 2: Framework of Our Self-Attentive Pyramid Network (SAPN)

attention [15] exploits spatial correlations of features by using attention scores to weight all features to obtain salient features. To make full use of the spatial correlations of features, we propose a deep self-attentive pyramid network (SAPN) for single image de-raining, which mainly consists of self-attentive pyramid module (SAM) and multi-scale pooling (MSP) module. Specifically, to efficiently exploit spatial contextual information, we propose the self-attention calculation units (SACUs) based encoding layers to enhance the encoding process, and SACUs based skip connections to enhance the symmetrical decoding process. With the assistance of SACUs, the encoding layers can better utilize the spatial correlations from input features. Besides, our SACUs based skip connections can not only contribute to the propagation of gradient flows, but pass the enhanced original feature signal from convolutional layers to symmetrical deconvolutional layers directly, which is helpful for recovering image details. Furthermore, since feature pyramid is helpful for multi-resolution feature representation, we apply multi-scale pooling in the shallow layers of our network. Extensive experiments on synthetic and real-world datasets demonstrate the superiority of our proposed method in terms of both quantitative and visual quality.

2 Self-attentive pyramid network (SAPN)

2.1 Network Architecture

As shown in Fig. 2, our SAPN consists of four main parts: multi-scale pooling module, shallow feature extraction, self-attentive pyramid module (SAM) and feature reconstruction.

Multi-scale pooling Module. Given I as input rainy image and \hat{R} as estimated rain streak image, then the output of SAPN is represented as follows:

$$\hat{O} = I - \hat{R}, \quad (1)$$

4 T. Guo et al.

where \hat{O} denotes the estimated de-rained output image. To model rain streaks with various scales and shapes from the input image I , we firstly introduce a multi-scale pooling operation $H_{msp}(\cdot)$ to get the multi-scale feature concatenation F_{msp} from I :

$$\begin{aligned} F_{msp} &= H_{msp}(I) \\ &= [S_1^{32 \times 32}(I), S_2^{16 \times 16}(I), S_3^{8 \times 8}(I), S_4^{4 \times 4}(I), I], \end{aligned} \quad (2)$$

where $[\cdot, \cdot, \dots, \cdot]$ denotes channel-wise concatenation and $S_i^{k \times k}(\cdot)$ represents the i -th $k \times k$ -scale pooling operation which is defined as:

$$S_i^{k \times k}(I) = U^{k \times k}(ReLU(Conv^{1 \times 1}(D^{k \times k}(I)))), \quad (3)$$

where $D^{k \times k}(\cdot)$ and $U^{k \times k}(\cdot)$ denote $k \times k$ -scale downsampling and upsampling respectively. $Conv^{1 \times 1}(\cdot)$ denotes a 1×1 convolutional layer.

Shallow feature extraction. After we get multi-scale feature concatenation F_{msp} from Equ. (3), the shallow feature representation F_{fr} can be obtained by

$$F_{fr} = H_{ex}(F_{msp}), \quad (4)$$

where $H_{ex}(\cdot)$ represents two consecutive 3×3 convolutional layers with 64 filters respectively, which are designed to extract the shallow feature representation F_{fr} from F_{msp} .

Self-attentive pyramid module (SAM). Given the shallow feature representation F_{fr} obtained from the above step, the self-attentive pyramid module (SAM), denoted as $H_{sam}(\cdot)$, adopts a pyramid encoder-decoder structure with Self-attention Calculation Units (SACU) embedded in it, and produce a rain streak layer feature representation F_{rs} :

$$F_{rs} = H_{sam}(F_{fr}). \quad (5)$$

The detailed description of SAM is given in Section 2.2.

Feature reconstruction part. After obtaining the rain streak layer feature representation F_{rs} , we can reconstruct the estimated rain streak \hat{R} using the feature reconstruction part $H_{rc}(\cdot)$, which is actually a 3×3 convolutional layer:

$$\hat{R} = H_{rc}(F_{rs}) = H_{sapn}(I), \quad (6)$$

where $H_{sapn}(\cdot)$ represents the function of our proposed SAPN.

Loss function. During the training process, our SAPN is optimized with loss function. To improve not only the pixel-wise reconstruction but the high-level semantic representation, we add perceptual loss to pixel-level L1 loss to get the combined loss L_C :

$$L_C = L_{L1} + \lambda L_P, \quad (7)$$

where λ denotes the trade-off coefficient between the two losses, and the L1 loss L_{L1} and the perceptual loss L_P are defined as:

$$L_{L1} = \frac{1}{CWH} \sum_{c=1}^C \sum_{w=1}^W \sum_{h=1}^H \|\hat{O}^{c,w,h} - O^{c,w,h}\|_1, \quad (8)$$

$$L_P = \frac{1}{CWH} \sum_{c=1}^C \sum_{w=1}^W \sum_{h=1}^H \|(V(\hat{O}))^{c,w,h} - (V(O))^{c,w,h}\|_2^2, \quad (9)$$

where C , W and H denote the channel, width and height dimension of the estimated de-rained image \hat{O} and the ground truth clean image O . $V(\cdot)$ represents the front layers of a pretrained VGG model which is regarded as the high-level feature extractor. The loss function is optimized by Adam optimizer.

After a full glance at the framework of the proposed SAPN, we can conclude that the deep feature representation in our SAPN heavily relies on the self-attentive pyramid module (SAM), which will be shown in the next section.

2.2 Self-attentive Pyramid Module (SAM)

Our SAM is based on the conventional encoder-decoder networks [16], which are widely used in image-to-image tasks. However, most existing encoder-decoder based networks focus on the network design, while rarely exploits the spatial correlations of features and thus limits representation capability of the network. To exploit such correlations inherent in features, we propose a novel self-attentive pyramid module (SAM).

As shown in Fig. 2, the core component of SAPN is self-attentive Pyramid Module (SAM), which is further composed of four Self-attention Calculation Units (SACUs), four encoders and four decoders. The detailed description of SACU will be given in the next section.

Given the feature representation F_{fr} obtained from shallow feature extraction step $H_{ex}(\cdot)$, the original U-net [16] simply encodes the features iteratively and feeds the encoded features to symmetrical decoder. However, the single encoding layer, which consists of several convolutional layers, can not fully utilize the spatial correlations of the features, thus leading to poor ability of modeling the long-range dependency inherent in the features. Given this, we embed self-attention calculation unit (SACU) $H_{sa,i}(\cdot)$ in each encoder $H_{en,i}(\cdot)$ to model the long-range spatial correlations, and thus the encoded features are enhanced before passing through the decoding layer. i -th encoder $H_{en,i}(\cdot)$ is composed of a 3×3 convolutional layer with stride 2 and doubled channels from input, and two 3×3 layers with ReLU activation, which keeps input channels. We can formulate the self-attentive encoding part of the SAM component $H_{sam}(\cdot)$ as:

$$\begin{aligned} F_{sa,i} &= H_{sa,i-1}(F_{en,i-1}), \\ F_{en,i} &= H_{en,i}(F_{sa,i}), \quad i = 1, 2, 3, 4, \end{aligned} \quad (10)$$

where $F_{sa,i}$ denotes the output of i -th SACU and $F_{en,i}$ denotes the output of i -th encoder. $F_{en,0}$ denotes F_{fr} for convenience. With the help of self-attention information obtained from the SACU, the encoding process can be enhanced to get more spatial correlation into consideration.

After the self-attentive encoding part, we obtain $F_{sa,i}$ plus the final output (i.e. $F_{en,4}$) of encoders as the input of following decoding part. Unlike the pyramid network and U-net, which directly utilize the symmetrical encoded features

6 T. Guo et al.

$F_{en,4-i}$ from skip connection as the extra information:

$$F_{de,i} = H_{de,i}([F_{de,i-1}, F_{en,4-i}]), \quad (11)$$

we adopt the extra self-attention information besides the original encoded features $F_{en,4-i}$, which is integrated in features $F_{sa,4-(i-1)}$, to decoder $H_{de,i}(\cdot)$ to get output features $F_{de,i}$. Similar with the encoder design, i -th decoder $H_{de,i}(\cdot)$ starts with a 3×3 deconvolutional layer with stride 2 and keeps the channels, followed by two consecutive 3×3 layers with ReLU activation which halves the channels in the former layer. Specifically, we utilize the obtained self-attention $F_{sa,i}$ to enhance the decoding process, which can be formulated as:

$$F_{de,i} = H_{de,i}([F_{de,i-1}, F_{sa,4-(i-1)}]), \quad (12)$$

where $F_{de,4}$, the output features of the last decoder, is also the final output of SAM and input of feature reconstruction layer, which is also denoted as F_{rs} . $F_{de,0}$, or $F_{en,4}$, is the output features of the last encoding layer and also the input of the first decoding layer. Through the skip connection which delivers the output self-attention $F_{sa,4-(i-1)}$ of $(4-i)$ -th SACU (i.e. $H_{sa,4-(i-1)}(\cdot)$), the i -th decoder $H_{de,i}(\cdot)$ can get not only the symmetrical features but their self-attention information directly since the skip connection structure in SACU. With the help of the extra self-attention information of features, the decoding process can be further enhanced with the spatial correlation provided by the self-attention information, which makes the representation of long-range dependency possible. The experiments in Section 3 demonstrate the performance gain of the utilization of the extra self-attention information. We will give a further explanation to the self-attention calculation unit (SACU) in the next section.

2.3 Self-attention Calculation Unit (SACU)

Self-attention focuses on the attention of feature maps towards themselves, which has been widely researched by previous works [15, 17]. The information provided by self-attention properly handles the problem that long-range feature dependency can not be efficiently convolved by the convolutional layers.

Self-Attentive Pyramid Network for Single Image De-raining

7

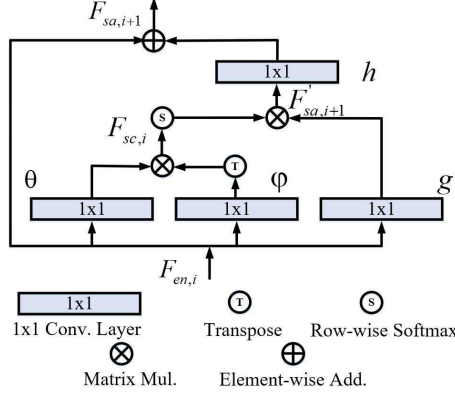


Fig. 3: Self-attention Calculation Unit (SACU)

As shown in Fig. 3, given the output features $F_{en,i}$ of i -th encoder $H_{en,i}(\cdot)$, we firstly obtain three embedding representation $\theta(F_{en,i})$, $\phi(F_{en,i})$, $g(F_{en,i})$ from three different 1×1 convolutional layers $\theta(\cdot)$, $\phi(\cdot)$ and $g(\cdot)$. Then the self-correlation $F_{sc,i}$ of feature map $F_{en,i}$ can be obtained via

$$F_{sc,i} = \theta(F_{en,i})\phi(F_{en,i})^T. \quad (13)$$

Then, we can get the self-weights $F'_{sc,i}$ by softmaxing each row of $F_{sc,i}$, and use it to weight the embedding representation $\theta(F_{en,i})$ by:

$$F'_{sa,i+1} = F'_{sc,i}g(F_{en,i}). \quad (14)$$

After that, the self-attention map $F''_{sa,i+1}$ is obtained via a further 1×1 convolutional layer $h(\cdot)$.

Furthermore, to boost the gradient transmission and avoid the gradient vanishing problem, we add skip connection from the input $F_{en,i}$ to the calculated self-attention map $F''_{sa,i+1}$. To better calibrate the influence between them, unlike the original non-local implementation [17], which regards the balance between the two terms as a hyper-parameter, we bring in a learnable parameter α as a trade-off weight. The final output $F_{sa,i+1}$ of SACU can be formulated as:

$$F_{sa,i+1} = F_{en,i} + \alpha F''_{sa,i+1}. \quad (15)$$

With the learnable parameter α , the weighting of the self-attention information becomes more flexible and thus leads to better utilization of self-attention.

3 Experiments

To validate the advantage of our method, we conduct tremendous experiments on various synthetic datasets and natural rainy images. Since the ground truth images are available in synthetic datasets, PSNR and SSIM are adopted as the evaluation criterion of the de-raining results. We calculate PSNR/SSIM in luminance

8 T. Guo et al.

channel of YCbCr space. Additionally, we compare our proposed SAPN with state-of-the-art de-raining methods, including Deep Detailed Network (DDN) [8], Joing Rain Detection and Removal (JORDER) [9], Density-aware Single Image De-raining using a Multi-stream Dense Network (DID-MDN) [10] and Progressive Recurrent Network (PReNet) [14].

3.1 Datasets

Synthetic Datasets To make a comparison with previous state-of-the-art de-raining approaches, we adopt three public benchmark synthetic datasets to train and evaluate our SAPN, including DDN-Dataset [8], DID-MDN-Dataset [10] and Rain100H [9]. Specifically, DDN-Dataset contains 14,000 rainy-clean image pairs which is synthesized by 1000 clean images. We randomly select 9100 image pairs as training dataset and use the left 4900 image pairs as the testing dataset. DID-MDN-Dataset is composed of 12000 training rainy-clean image pairs and 1201 testing image pairs. Rain100H contains 100 testing images and there are 1800 training image pairs in the corresponding training dataset (i.e. RainTrainH).

Real-world images To validate the effectiveness of the proposed network in real world rainy scenes, we randomly select some images from the previous de-raining works [8, 9, 18, 19] and the internet.

3.2 Training Details

Table 1: Quantitative results of average PSNR(dB)/SSIM compared with state-of-the-art de-raining works. The two best-performing methods are marked in **bold** and underlined respectively.

Dataset	Matric	Input	DDN [8] (CVPR'17)	JORDER [9] (CVPR'17)	DID-MDN [10] (CVPR'18)	Our SAPN
DID-MDN-Dataset	PSNR(dB)	23.63	<u>30.08</u>	26.80	29.36	30.86
	SSIM	.7313	.8788	.8361	<u>.9002</u>	.9230
DDN-Dataset	PSNR(dB)	23.74	<u>30.00</u>	26.47	28.00	30.26
	SSIM	.7499	<u>.8932</u>	.8276	.8776	.9110
Rain100H	PSNR(dB)	13.56	22.26	26.10	<u>26.35</u>	27.06
	SSIM	.3800	.6928	.7971	<u>.8287</u>	.8474

For each of the three datasets, we train our SAPN on a 1080 ti GPU on the training dataset, and evaluate the model on corresponding testing dataset. We train our model for 300, 350 and 50 epochs for DID-MDN-Train, DDN-Train, and RainTrainH respectively. The initial learning rate is set to $2 \cdot 10^{-4}$ and decreased linearly at the end of every epoch. To avoid the problem of over-fitting, we use a



Fig. 4: De-raining results on sample image from DDN-Testset.

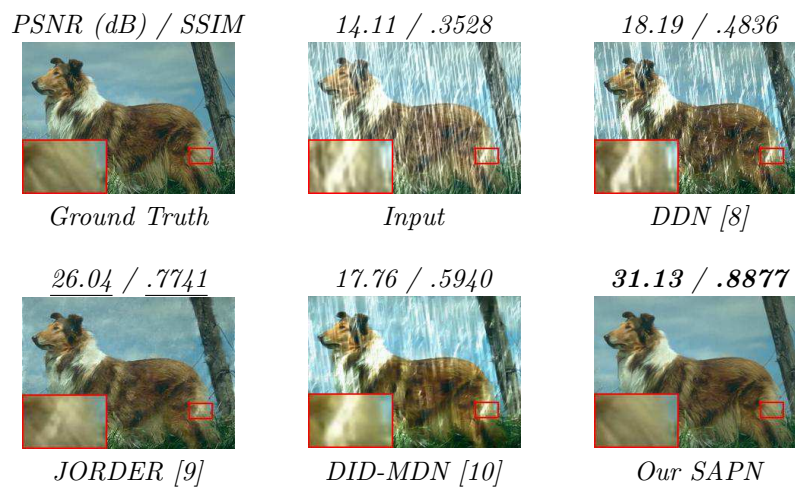


Fig. 5: De-raining results on sample image from Rain100H.



Fig. 6: De-raining results on sample real-world image.

weight decay of 10^{-5} and Adam optimizer with betas 0.5 and 0.999. The model is trained on the Pytorch framework.

3.3 Results on Synthetic Datasets

The details of evaluation results on synthetic datasets are shown in Table 1. Note that the the pretrained model of PResNet [14] is trained from other datasets, thus we did not report the quantitative results for fair comparison. Results show that our method outperforms other state-of-the-arts consistently. This is mainly because our SAPN utilizes both multi-scale information and spatial correlations of features, which enhances the feature representation capability of the network. In contrast, DDN [8] learns the mapping from high-frequency details of rainy images to clean ones with ResNet [20], while DID-MDN [10] utilize multi-scale features and multi-stream DenseNet [18] architecture, both of which do not take spatial correlations inherent in features into consideration.

Besides the quantitative evaluation on synthetic datasets, we also randomly select several images from the testing datasets to validate the visual effect. As shown in Fig. 4 and Fig. 5, our method obtains better visual results. For example, the alphabets in enlarged region in Fig. 4 are clearly recovered by our SAPN, while other methods fail to remove long heavy rain streaks or bring in unpleasant artifacts. Another sample in Fig. 5 also show that our SAPN keeps the background scenes better and removes the rain streaks more cleanly, especially when the image has some long rain streaks or other objects with long shapes, since the adoption of self-attention mechanism enhances the capability of the network to capture long-range dependency and non-local similarity.

3.4 Qualitative Evaluation on Real-world Images

To verify the performance gain of SAPN over previous methods on rainy scenes in real world, we also test our SAPN and other methods on real-world images.

The de-raining results on a randomly selected real world image sample is shown in Fig. 6. Noticeably, our method achieves extremely better results when the rain streaks in rainy image are longer than average, just because we adopt self-attention mechanism in our network design, which can better leverage non-local similarity of input rainy image and attain long-range dependency more effectively and more efficiently. This specialty of our SAPN helps locate rainy areas in input rainy images, leading to better final de-raining results. It is clearly shown in Fig. 6 that our SAPN produces preferable results compared with other methods, which tend to either under de-rain or over de-rain the natural rainy images. Specifically, all other four methods fail to remove all long rain streaks, while JORDER even brings in severe artifacts. In contrast, our method not only removes more rain streaks, but preserves background details better.

3.5 Ablation Study

To verify benefits of each individual component, including multi-scale pooling (MSP) module and SACUs, we train some variants of our SAPN on RainTrainH and evaluate trained models on Rain100H. The results are shown in Table 2.

Table 2: Ablation study of our proposed SAPN on SACUs and multi-scale pooling module on Rain100H.

Methods	U_a	U_b	U_c	U_d
SACUs?		✓		✓
MSP?			✓	✓
PSNR(dB)	26.78	26.95	26.89	27.06

We can conclude that the adoption of SACUs effectively promotes the de-raining results of the basic pyramid encoder-decoder network (U_a), while MSP also improves the performance of the network effectively. The combination of SACUs and MSP leads to our final SAPN architecture (U_d).

4 Conclusion

In this paper, we propose a pyramid encoder-decoder network with self-attention calculation units for single image de-raining. Compared with previous methods which does not exploit spatial correlations of features, our method explicitly learns the self-correlation inherent in output features of each encoder layer, making the encoding and symmetrical decoding process more self-attentive and better resolve the long-range dependency problem in images, leading to better eventual de-raining results, especially in long rain streaks conditions. In order to further improve the de-raining results, we add a multi-scale pooling module before feature extraction, which leads to even higher quantitative performance

12 T. Guo et al.

and much better visual experience. Tremendous experiments on various datasets validate that our network outperforms the state-of-the-art methods.

References

1. Jérémie Bossu, Nicolas Hautière, and Jean-Philippe Tarel, "Rain or snow detection in image sequences through use of a histogram of orientation of streaks," *IJCV*, 2011.
2. Weihong Ren, Jiandong Tian, Zhi Han, Antoni Chan, and Yandong Tang, "Video desnowing and deraining based on matrix decomposition," in *CVPR*, 2017.
3. Wei Wei, Lixuan Yi, Qi Xie, Qian Zhao, Deyu Meng, and Zongben Xu, "Should we encode rain streaks in video as deterministic or stochastic," in *CVPR*, 2017.
4. Kshitiz Garg and Shree K Nayar, "Detection and removal of rain from videos," in *CVPR*, 2004.
5. Yu Luo, Yong Xu, and Hui Ji, "Removing rain from a single image via discriminative sparse coding," in *ICCV*, 2015.
6. Yu Li, Robby T Tan, Xiaojie Guo, Jiangbo Lu, and Michael S Brown, "Rain streak removal using layer priors," in *CVPR*, 2016.
7. Lei Zhu, Chi-Wing Fu, Dani Lischinski, and Pheng-Ann Heng, "Joint bilayer optimization for single-image rain streak removal," in *ICCV*, 2017.
8. Xueyang Fu, Jiabin Huang, Delu Zeng, Yue Huang, Kinghao Ding, and John Paisley, "Removing rain from single images via a deep detail network," in *CVPR*, 2017.
9. Wenhan Yang, Robby T Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan, "Deep joint rain detection and removal from a single image," in *CVPR*, 2017.
10. He Zhang and Vishal M Patel, "Density-aware single image de-raining using a multi-stream dense network," *CVPR*, 2018.
11. Yi Chang, Luxin Yan, and Sheng Zhong, "Transformed low-rank model for line pattern noise removal," in *ICCV*, 2017.
12. Xia Li, Jianlong Wu, Zhouchen Lin, Hong Liu, and Hongbin Zha, "Recurrent squeeze-and-excitation context aggregation net for single image deraining," in *ECCV*, 2018.
13. Guanbin Li, Xiang He, Wei Zhang, Huiyou Chang, Le Dong, and Liang Lin, "Non-locally enhanced encoder-decoder network for single image de-raining," in *MM*, 2018.
14. Dongwei Ren, Wangmeng Zuo, Qinghua Hu, Pengfei Zhu, and Deyu Meng, "Progressive image deraining networks: A better and simpler baseline," *CVPR*, 2019.
15. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin, "Attention is all you need," in *NIPS*, 2017.
16. Olaf Ronneberger, Philipp Fischer, and Thomas Brox, "U-net: Convolutional networks for biomedical image segmentation," in *MICCAI*, 2015.
17. Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He, "Non-local neural networks," in *CVPR*, 2018.
18. Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger, "Densely connected convolutional networks," in *CVPR*, 2017.
19. He Zhang, Vishwanath Sindagi, and Vishal M Patel, "Image de-raining using a conditional generative adversarial network," *arXiv*, 2017.
20. Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *CVPR*, 2016.