

# Efficient and Accurate Multi-scale Topological Network for Single Image Dehazing

Qiaosi Yi<sup>†</sup>, Juncheng Li<sup>†</sup>, Faming Fang<sup>\*</sup>, Aiwen Jiang, Guixu Zhang

**Abstract**—Single image dehazing is a challenging ill-posed problem that has drawn significant attention in the last few years. Recently, convolutional neural networks have achieved great success in image dehazing. However, it is still difficult for these increasingly complex models to recover accurate details from the hazy image. In this paper, we pay attention to the feature extraction and utilization of the input image itself. To achieve this, we propose a Multi-scale Topological Network (MSTN) to fully explore the features at different scales. Meanwhile, we design a Multi-scale Feature Fusion Module (MFFM) and an Adaptive Feature Selection Module (AFSM) to achieve the selection and fusion of features at different scales, so as to achieve progressive image dehazing. This topological network provides a large number of search paths that enable the network to extract abundant image features as well as strong fault tolerance and robustness. In addition, AFSM and MFFM can adaptively select important features and ignore interference information when fusing different scale representations. Extensive experiments are conducted to demonstrate the superiority of our method compared with state-of-the-art methods.

**Index Terms**—Image dehazing, multi-scale topological network, feature fusion, adaptive feature selection.

## I. INTRODUCTION

**H**AZE is a common atmospheric phenomenon produced by small floating particles. Particulate matter floating in the air causes light scattering and attenuation, thereby reducing the visibility of distant objects. However, hazy images will cause difficulties with their processing and analysis, which will seriously affect the performance of downstream tasks such as image classification, image segmentation, object detection, crowd counting, and other high-level computer vision tasks. This is not conducive to the construction of safe and stable artificial intelligence systems, such as video surveillance systems and unmanned driving systems. In order to solve this problem, the task of image dehazing, especially single image dehazing came into being and has drawn significant attention in the last few years.

Single image dehazing is an extremely hot topic in computer vision, which aims to reconstruct a haze-free image from the hazy one (Fig. 1). However, due to the absorption and reflection of the haze, the captured scene image will suffer

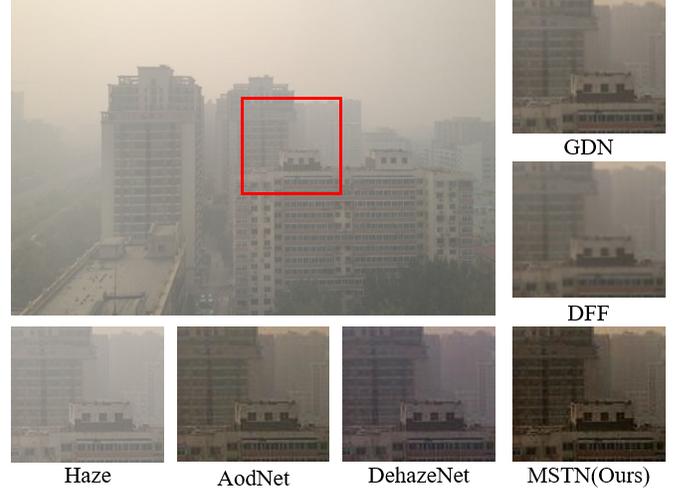


Fig. 1. An example of single image dehazing. Obviously, the haze-free image reconstructed by our MSTN shows better visual effect.

from color distortion, blur, and low contrast, which causes the quality of the image to deteriorate. Therefore, single image dehazing still is a challenging task and many methods have been proposed to try to solve this task.

The atmosphere scattering model provides a theoretical basis for hazy imaging and is also a basic method for image dehazing. As shown in Eq. (1), the atmosphere scattering model can be defined as:

$$I_i(x) = J_i(x)t(x) + A(1 - t(x)), \quad i = 1, 2, 3 \quad (1)$$

where  $I(x)$  is the observed hazy image,  $J(x)$  is the clear image, and  $A$  represents the global atmospheric light intensity. Meanwhile,  $i$  denotes  $R, G, B$  channels in a RGB image,  $x$  represents the pixel locations,  $I_i(x)$  and  $J_i(x)$  represents the value of the  $i$ -th channel of the hazy or clear image in location  $x$ .  $t(\cdot)$  is the medium transmission function,  $t(x) = e^{-\beta d(x)}$ ,  $\beta$  and  $d(x)$  represent the atmosphere scattering parameter and the scene depth, respectively. The atmosphere scattering model shows that single image dehazing is an ill-posed problem, which is a challenging task without the priors of  $A$  and  $t(x)$

$$J(x) = \frac{I(x) - A}{t(x)} + A. \quad (2)$$

In the past, in order to better deal with the problem of single dehazing, many methods have been proposed to learned different prior knowledge to estimate  $A$  and  $t(x)$ , then obtain the haze-free image based on the atmosphere scattering model. For example, dark channel prior [1], color attenuation prior

<sup>†</sup>: Equal contribution. <sup>\*</sup>: The corresponding author.

Q. Yi, J. Li, F. Fang, and G. Zhang are with the Shanghai Key Laboratory of Multidimensional Information Processing, East China Normal University, Shanghai, China, and also with the school of Computer Science and Technology, East China Normal University, Shanghai, China. (E-mail: qiaosiyijoyies@gmail.com, cvjunchengli@gmail.com, fmfang@cs.ecnu.edu.cn, gxzhang@cs.ecnu.edu.cn)

A. Jiang is with the School of Computer and Information Engineering, Jiangxi Normal University, Nanchang, China. (E-mail: jiangaiwen@jxnu.edu.cn)

[2], and non-local prior [3] are proposed for transmission function  $t(\cdot)$  estimation. Meanwhile, some works focus on estimating the atmospheric light  $A$ , such as [4], [5]. Based on the estimated  $\hat{t}(x)$  and  $\hat{A}$ , the clear image  $\hat{J}$  can be recovered by the following formulation

$$J(x) = \frac{I(x) - \hat{A} \cdot (1 - \hat{t}(x))}{\hat{t}(x)} = \frac{1}{\hat{t}(x)} I(x) - \frac{\hat{A}}{\hat{t}(x)} + \hat{A}. \quad (3)$$

However, due to the complexity of the real environment, the prior may be easily violated in practice. Therefore, the methods based on the atmosphere scattering model may not be able to accurately estimate the transmission map and the global atmospheric light intensity, resulting in the inability to obtain clear haze-free images. This will greatly limit the model speed, versatility, and performance.

Recently, convolutional neural networks (CNNs) have achieved remarkable success in many computer vision tasks and also greatly promoted the development of image dehazing. With the powerful feature extraction capabilities of CNN, more and more CNN-based image dehazing methods have been proposed for  $A$  and  $t(x)$  estimation or directly learn the mapping between hazy and clear images. For example, Cai et al. proposed the first CNN model (Dehazenet [6]) to directly remove haze from the hazy image. Li et al. proposed a all-in-one dehazing Network (AODNet [7]), which based on a re-formulated atmospheric scattering model and directly generates the clean image through a light-weight CNN. After that, CNN-based image dehazing models have been blooming and refreshing the best results, including PFF-Net [8], DCPDN [9], EPDN [10], PDR-Net [11], GDN [12], and DFF [13]. Although the aforementioned methods have made a big breakthrough in image dehazing. However, most existing image dehazing models have the following shortcomings:

- 1) Most existing methods focus on microstructure design, that is, build the network and achieve image dehazing by stacking the carefully designed feature extraction modules. This modular design strategy ignores the connectivity of the model. In addition, this cascaded structural design greatly reduces the possible topological paths, which is not conducive to building an effective model.
- 2) Most existing methods ignore the morphological difference of hazy images at different scales. Therefore, these models do not pay attention to the extraction, propagation, fusion, and utilization of the multi-scale image features.
- 3) The structure of these models is getting bigger, deeper, and more complex, which is not conducive to building an efficient and real-time dehazing model.

According to [14], the lower-level features have higher resolution and more texture details but lower semantic information, and the higher-level features have more semantic information but fewer texture details. Therefore, the core of this work is to build an effective network that can fully extract and utilize image features at different stages. Specifically, we propose a Multi-scale Topological Network (MSTN) to progressively remove the haze in the hazy image. MSTN is a topological network that can promote the transmission and utilization

of image feature flows. Meanwhile, these topological sub-nets enable the network to detect rich image features while increasing the fault tolerance of the model. Considering the effectiveness of multi-scale image features, we also introduce the multi-scale strategy into the model. Therefore, MSTN can be considered as a multi-branch network and each branch is used to extract images features at different scales. However, if these branches are independent of each other, they cannot form a topological network, which will greatly reduce the model performance. In order to solve this problem, we take the output of the lower-resolution branch as the input of the previous branch. In addition, we design an Adaptive Feature Selection Module (AFSM) and a Multi-scale Feature Fusion Module (MFFM) to realize automatic selection and fusion of multi-scale image features, which helps to make full use of the features of the image itself.

In summary, the main contributions of this work include:

- We reveal the importance of topology for deep network design and proposed a Multi-scale Topological Network (MSTN) for image dehazing, which shows stronger robustness and versatility. Compared with existing models, MSTN achieves better results with less execution time.
- We design a Multi-scale Feature Fusion Module (MFFM), which can promote the interaction and fusion between different scale features, thereby improving the utilization of multi-scale features.
- We propose an Adaptive Feature Selection Module (AFSM) to automatically select image features at different scales. Compared with directly adding all different scales features together, this module can effectively remove redundant features to achieve better feature extraction and utilization.

The remaining parts of this paper are organized as follows. Section II reviews related works including prior-based and learning-based image dehazing methods, topological network, multi-scale feature extraction, and Attention Mechanism. A detailed explanation of the proposed MSTN is given in Section III. The experimental results and ablation analysis are presented in Section IV and V, respectively. Finally, we draw a conclusion in Section VI.

## II. RELATED WORK

### A. Single Image Dehazing

1) *Prior-based Methods*: The prior-based methods use the characteristics of the image to estimate  $A$  or  $t(x)$  and recover the clear haze-free image according to the atmospheric scattering model. For example, Fattle [15] added a surface shadow factor to the atmospheric scattering model to estimate the transmission map; He et al. [1] proposed a dehazing algorithm based on dark channels prior (DCP), which estimates the transmission map by the DCP; Fattle [16] proposed a color-line prior dehazing method based on the observation that the color of a small image patch exhibits a one-dimensional distribution in the RGB color space; Although these prior-based methods have achieved varying degrees of success, their performance depends on the accuracy and validity of the proposed priors.

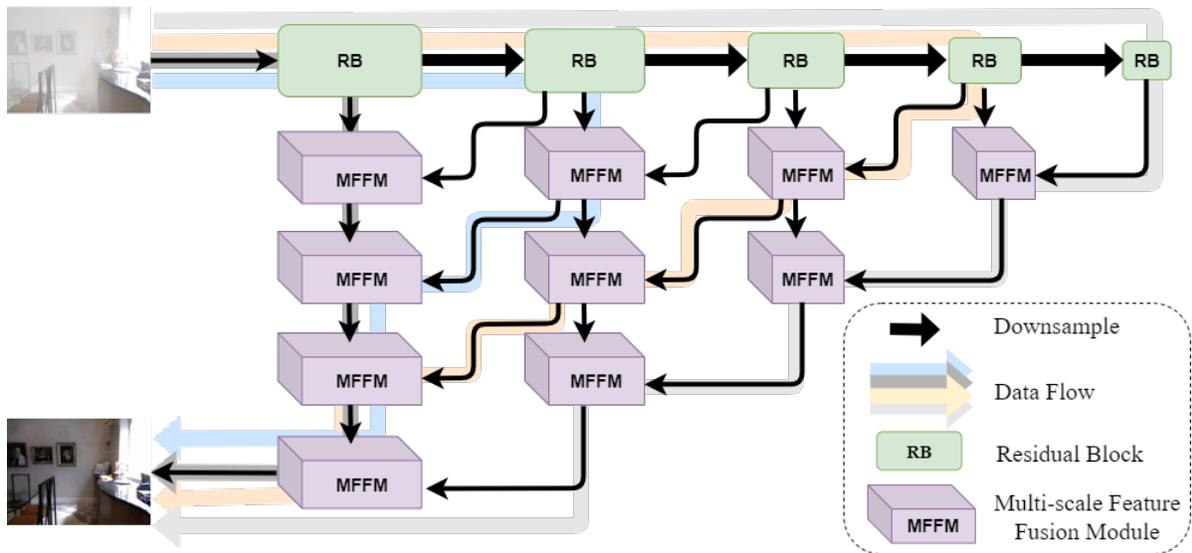


Fig. 2. The architecture of our proposed Multi-scale Topological Network (MSTN), which is a multi-branch network that contains five rows and five columns.

2) *Learning-based Methods*: With the rise of CNN, learning-based methods have become the mainstream method and have achieved tremendous development. These methods usually use CNN to estimate  $A$  and  $t(x)$  or directly recover the clear haze-free image by an efficient CNN. For example, Cai et al. [6] adopt a three-layer convolutional neural network (Dehazenet) to estimate the transmission map; Zhang et al. [9] proposed a density dehazing network (DCPDN) that can simultaneously estimate the transmission map and atmospheric light intensity; Liu et al. [12] proposed a grid network (GDN) to directly reconstruct clear images; Dong et al. [13] proposed multi-scale based deep network which works on strengthen-operate-subtract-boosting strategy for image dehazing (DFF). Although these learning-based methods have made great progress, they did not fully extract the features of the hazy image itself, resulting in sub-optimal reconstruction results.

3) *Topological Network*: Topology is a discipline that studies the properties of geometric figures or spaces that can remain unchanged after continuously changing shapes. It only considers the positional relationship between objects without considering their shape and size. The biggest characteristic of the topological architecture is its invariance under local deformation, which can simplify the network design. For example, Attara et al. [17] proposed a method that based on supervised machine learning algorithms and utilizes the topological similarities of networks for the classification task. Li et al. [18] proposed a recursive fractal network that can construct an infinite variety of topological structures through a simple basic component. At the same time, this structure has been proven to be more fault-tolerant, stable, and robust. Thus, the topological structure will be the focus of our research.

4) *Multi-scale Feature Extraction and Utilization*: Plenty of researches have pointed out that the image will exhibit different characteristics at different scales. Therefore, making

full use of the features of the input image itself can further improve model performance. In recent years, many works have been proposed for multi-scale features extraction and utilization, which can be roughly divided into two categories: (i) The most widely used method is to obtain images with different resolutions after multiple downsampling operations, and then extract the features separately. This type of method is commonly used in image segmentation and object detection tasks, such as FPN [14], FPT [19], and PyConvResNet [20]. (ii) Another method is to extract image features by different convolutional kernels. This type of method adjusts the size of the receptive field through different convolutional kernels, so as to achieve multi-scale feature extraction. The most famous methods includes VGG [21], MSRN [22], MSIN [23], MSFFRN [24], and MDCN [25]. In this work, we aim to introduce the multi-scale strategy into the topological network to better mine and utilize multi-scale image features.

5) *Attention Mechanism*: Recent years, the attention mechanism [26]–[31] has been shown significant advantages in a range of tasks, from neural machine translation in natural language processing to image captioning in image understanding. The important information is highlighted by the attention mechanism and the less useful information is suppressed. Attention has been widely used in recent applications such as person Re-ID, image recovery, and segmentation. To boost the performance of image classification, SENet [26] brings an effective, lightweight gating mechanism to self-recalibrate the feature map via channel-wise importances. Beyond channel, CBAM [27] introduce spatial attention in a similar way. Furthermore, SKNet [28] focus on the adaptive receptive field size of neurons by introducing the attention mechanisms. Different them, the core of our method, which takes different scales features as input for learning and output the selected and fused features, is to automatically respond and select features from different scale inputs. In the image dehazing

task, the GDN [12] learn a coefficient for adding different scale feature as the attention mechanism. Moreover, the coefficient is learned by global learning and do not depend on the different scale feature. But, our method utilize the attention mechanism to learn the relationship between the different scale feature and highlight the most informative feature expressions in the different scale feature.

### III. MULTI-SCALE TOPOLOGICAL NETWORK (MSTN)

In this paper, we propose a Multi-scale Topological Network (MSTN) for single image dehazing. As shown in Fig. 2, MSTN is essentially a multi-branch network, which contains  $i$  rows and  $j$  columns. Among them,  $i$  denotes the depth of the network and  $j$  represents the scale of the model, respectively. Meanwhile, we can clearly observe that at each branch the model contains one Residual Block (RB [32]) and several Multi-scale Feature Fusion Module (MFFM). It is worth noting that all these branches are used to extract image features at different scales and the input of each branch is obtained from the output of the previous branch through downsample operation. In addition, RBs are used for feature extraction and MFFMs are the core module of MSTN, which are used for multi-scale feature selection and fusion. However, if these branches are independent of each other, multi-scale features cannot interact together, which will greatly reduce the model performance. In order to solve this problem, we introduce skip connection between the adjacent branches. In other words, the outputs of the current branch are sent to the previous branch. Therefore, image features at different scales can be transferred, interacted, and merged via the MFFM.

Define  $I_{hazy}$  and  $I_{clear}$  as the input hazy image and the reconstructed haze-free image, the model can be defined as

$$I_{clear} = F_{MSTN}(I_{hazy}), \quad (4)$$

where  $F_{MSTN}(\cdot)$  represents the proposed MSTN. As mentioned above, MSTN is a multi-branch network that contains  $i$  rows and  $j$  columns, each row denotes the depth of the network and each column denotes the different scales of the model. We define the first row and first column as  $i = 0$  and  $j = 0$ , respectively. Therefore, the outputs ( $R_{i,j}$ ) of each RB or MFFM can be defined as

$$R_{i,j} = \begin{cases} F_{i,j}(I_{hazy}) & \text{when } i = 0, j = 0 \\ F_{i,j}(R'_{i,j-1}) & \text{when } i = 0, j > 0 \\ M_{i,j}(R_{i-1,j}, R_{i-1,j+1}) & \text{when } i > 0 \end{cases} \quad (5)$$

where  $F_{i,j}(\cdot)$  and  $M_{i,j}(\cdot)$  denote the operation of RB and MFFM in the  $i$ -th row and  $j$ -th column, respectively. Meanwhile,  $R'$  is the result obtained by the downsample operation

$$R' = R \downarrow_2. \quad (6)$$

It is worth noting that the downsampling operation is realized using a convolutional layer instead of traditional methods such as bilinear or bicubic interpolation. This is because bilinear and bicubic will cause a lot of information to be lost, which is not conducive to image reconstruction, so we use

convolutional layer to let the model automatically learn the redundant features that need to be removed.

During training, MSTN is optimized with  $L_1$  loss function. Therefore, given a training dataset  $\left\{ I_{hazy}^n, I_{clear}^n \right\}_n^N$ , we solve

$$\hat{\theta} = \arg \min_{\theta} \frac{1}{N} \sum_{n=1}^N \|F_{\theta}(I_{hazy}^n) - I_{clear}^n\|_1, \quad (7)$$

where  $\theta$  denotes the parameter set of our model and  $F(\cdot)$  represents the proposed MSTN. Each module of the network will be described in the following sections.

#### A. Multi-scale Topological Architecture

In this paper, we propose a multi-scale topological architecture as the backbone of MSTN. Similar to the Feature Pyramid Network [14], MSTN also adopts the pyramid-like structure to obtain multi-scale image features. In other words, we gradually downsample the resolution of the image and extract image features at different scales. After that, we progressively restore the resolution of the image and use the extracted multi-scale features to reconstruct the final haze-free image. This strategy can fully mine the potential features of the input image itself, improve the model performance, and reduce the memory consumed during runtime. However, most of the previous works are simply add all the features extracted from each scale branch, which is not conducive to the interaction between different scales features. In order to solve this problem, we introduce skip connection between the adjacent branches. Therefore, image features with different scales can be interacted and merged through MFFM. It is worth noting that the intermediate results of each branch are sent to the corresponding position of the previous branch for feature selection and fusion rather than just the final output. This allows the hierarchical features to be fully utilized, which can further improve the quality of reconstructed images. Meanwhile, these skip connections make the model constitute a topological network, which provide a large number of search paths that enable the network to extract abundant image features to reconstruct high-quality haze-free images. In Fig. 2, the color data flows represent some examples of the path of the model. Among them, the gray one represents the complete pathway, the dark gray represents the path of the model without multi-scale strategy. In addition, the blue and orange lines represent two intermediate paths. This topological architecture makes the network contains multiple sub-networks and all these sub-networks complement each other, greatly improving the stability and fault tolerance of the model.

#### B. Multi-scale Feature Fusion Module (MFFM)

MFFM is the core module of MSTN, which is designed for multi-scale image feature selection, interaction, and fusion. As shown in Fig. 3, MFFM takes  $R_{i-1,j}$  and  $R_{i-1,j+1}$  as inputs and output the fused image features  $R_{i,j}$ . According to the figure, we can clearly observe that the module contains a adaptive feature selection module (AFSM), a convolutional layer, a deconvolutional layer, a residual block (RB), and a residual skip connection. Firstly, we apply downsampling

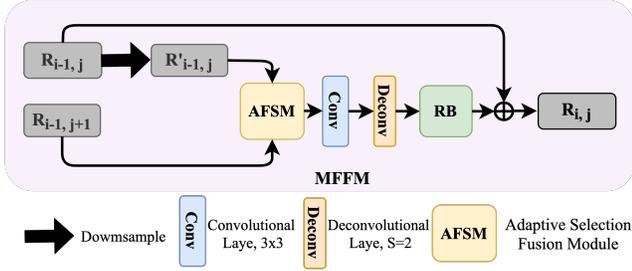


Fig. 3. The complete architecture of MFFM.

operation on  $R_{i-1,j}$  to obtain  $R'_{i-1,j}$ . Then,  $R'_{i-1,j}$  and the extracted features from the next branch  $R_{i-1,j+1}$  are sent to the AFSM for feature selection. This is a crucially important step that used to automatically select useful features and remove redundant features. After that, a convolutional layer, a deconvolutional layer, and a residual block are applied to the selected multi-scale image features to obtain new representations. Finally, we introduce local residual learning strategy to further improve the information flow. The introduced residual learning strategy make the module only needs to learn the different areas between the input and output features, which can greatly accelerate the learning process. Meanwhile, this allows the ASFM to selectively select the missing features from different scale branches. With the help of MFFM, the model has enough flexibility for selecting important features from different scale representations and can expand the representation capabilities of CNN.

### C. Adaptive Feature Selection Module (AFSM)

According to Lin et al. [14], we know that image features with different scales have different semantic information. Making full use of multi-scale image features can effectively improve the quality of reconstructed images. However, most existing methods directly cascade or add all multi-scale image features for image reconstruction, it will bring a lot of redundant features that not conducive to building a efficient and accurate model. In 2019, Li et al. [28] proposed a Selective Kernel Networks (SKNet), which can adaptively adjust its receptive field size based on multiple scales of input information. Inspired by this, we introduce the attention mechanism to the model and propose an Adaptive Feature Selection Module (AFSM) for different image feature selection and fusion. As shown in Fig. 4, AFSM takes different scales features  $R'_{i-1,j}$  and  $R_{i-1,j+1}$  as inputs for learning, and output the selected and fused  $R''_{i,j}$ . Specifically, we first fuse the results from different branches via an element-wise summation:

$$R'_{i,j} = R'_{i-1,j} + R_{i-1,j+1}. \quad (8)$$

Then we generate channel-wise statistics  $\mathbf{s} \in \mathbb{R}^C$  by using global average pooling. The  $c$ -th element of  $\mathbf{s}$  is calculated by shrinking  $R'_{i,j}$  through spatial dimensions  $H \times W$ :

$$s^c = f_g(R'_{i,j}{}^c) = \frac{1}{H \times W} \sum_{x=1}^H \sum_{y=1}^W R'_{i,j}{}^c(x, y) \quad (9)$$

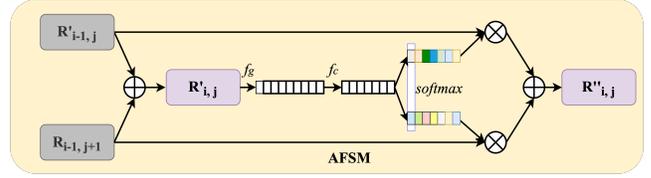


Fig. 4. The complete architecture of AFSM.

After that, we applied a fully connected layer to generate compressed features  $\mathbf{z} \in \mathbb{R}^{d \times 1}$  for precise and adaptive selection

$$\mathbf{z} = f_c(\mathbf{s}) \quad (10)$$

Finally, a soft attention across channels is used to adaptively select important information from different branches.

$$a = \frac{e^{\mathbf{A}\mathbf{z}}}{e^{\mathbf{A}\mathbf{z}} + e^{\mathbf{B}\mathbf{z}}}, b = \frac{e^{\mathbf{B}\mathbf{z}}}{e^{\mathbf{A}\mathbf{z}} + e^{\mathbf{B}\mathbf{z}}}, \quad (11)$$

where  $\mathbf{A}, \mathbf{B} \in \mathbb{R}^{C \times d}$ , and  $a, b$  denote the attention vector of  $R'_{i-1,j}$  and  $R_{i-1,j+1}$ , respectively. Specifically, the softmax function is used on  $a$  and  $b$ , so  $a + b = 1$ . After getting  $a$  and  $b$ , the output  $R''_{i,j}$  can be calculated as follow:

$$R''_{i,j} = a \cdot R'_{i-1,j} + b \cdot R_{i-1,j+1}. \quad (12)$$

With the help of this module, our MSTN can efficiently and automatically select and fuse multi-scale image features. This provides a new solution for image restoration task which based on multi-scale architecture.

## IV. EXPERIMENTS

### A. Dataset

In this paper, we use RESIDE [33], Middlebury [36], and NH-HAZE [37] to prove the effectiveness of our proposed MSTN on image dehazing task. Moreover, we also adopt the derain dataset (DID-MDN [38]) further verify the effectiveness of the proposed network on other image restoration tasks, thereby verifying the scalability and versatility of MSTN.

1) *RESIDE*: RESIDE [33] is a large-scale image dehazing dataset, which includes synthetic hazy images of indoor and outdoor and real-world hazy images. In RESIDE, the atmospheric scattering model is adopted to generate the synthetic hazy images. In this work, we use Indoor Training Set (ITS) and Outdoor Training Set (OTS) as the training dataset and select Synthetic Objective Testing Set (SOTS) and Hybrid Subjective Testing Set (HSTS) as the test dataset, respectively. Among them, ITS contains 1,399 clear images and 13,990 hazy images with the size of  $620 \times 460$ . Each clear image generates 10 hazy images with  $\beta \in [0.6, 1.8]$  and  $A \in [0.7, 1.0]$ , and the depth map  $d(x)$  comes from the NYU Depth V2 [39] and Middlebury Stereo datasets [40]. Similar to ITS, OTS also contains a large number of images, but the depth map  $d(x)$  of OTS is estimated by using the algorithm developed in [41] and  $\beta \in [0.04, 0.2]$ ,  $A \in [0.8, 1.0]$ . The SOTS contains 500 indoor hazy images and 500 outdoor hazy images, and their generation methods are the same as ITS and OTS, respectively. In addition, both synthetic hazy images and real-world hazy

TABLE I

QUANTITATIVE (PSNR/SSIM) COMPARISONS WITH SOTA IMAGE DEHAZING METHODS ON RESIDE-SOTS [33] (INDOOR AND OUTDOOR) AND RESIDE-HSTS [33] (SYNTHETIC). THE BEST AND SECOND BEST RESULTS ARE HIGHLIGHTED WITH RED AND BLUE FONTS, RESPECTIVELY.

Method	DCP [1]	CAP [2]	DehazeNet [6]	MSCNN [34]	NLD [3]	AODNet [7]	DCPDN [9]	GFN [35]	GDN [12]	DFF [13]	MSTN (Ours)	
SOTS (Indoor)	PSNR↑	16.61	19.05	21.14	17.12	17.29	19.06	15.85	22.30	32.16	33.75	35.37
	SSIM↑	0.855	0.836	0.847	0.796	0.778	0.850	0.818	0.880	0.984	0.985	0.987
SOTS (Outdoor)	PSNR↑	19.13	18.12	22.46	19.48	17.97	20.29	19.93	21.55	30.86	32.21	32.61
	SSIM↑	0.815	0.758	0.851	0.839	0.821	0.877	0.845	0.844	0.982	0.979	0.981
HSTS (Synthetic)	PSNR↑	14.48	21.57	24.48	18.64	18.92	20.55	22.94	22.06	32.75	32.72	35.48
	SSIM↑	0.761	0.873	0.915	0.817	0.741	0.897	0.874	0.847	0.983	0.9781	0.987

TABLE II

QUANTITATIVE (PSNR/SSIM) COMPARISONS WITH SOTA IMAGE DEHAZING METHODS ON MIDDLEBURY [36] AND NH-HAZE [37]. THE BEST AND SECOND BEST RESULTS ARE HIGHLIGHTED WITH RED AND BLUE FONTS, RESPECTIVELY

Method	DCP	AODNet	DCPDN	GFN	GDN	DFF	MSTN (Ours)	
MiddleBury	PSNR↑	11.94	13.94	12.23	14.01	14.21	15.82	17.50
	SSIM↑	0.762	0.764	0.725	0.754	0.778	0.868	0.863
NH-HAZE	PSNR↑	10.57	15.41	17.42	15.17	15.23	16.21	18.42
	SSIM↑	0.52	0.57	0.61	0.52	0.56	0.58	0.63

TABLE III

QUANTITATIVE (SSEQ/BLIINDS-II) COMPARISONS ON RESIDE-HSTS [33]. THE BEST AND SECOND BEST RESULTS ARE HIGHLIGHTED WITH RED AND BLUE FONTS, RESPECTIVELY

Method	HSTS			
	Synthetic		Real	
	SSEQ↓	BLIINDS-II↓	SSEQ↓	BLIINDS-II↓
DCP [1]	86.15	90.70	68.65	69.35
CAP [2]	85.32	85.75	67.67	63.55
DehazeNet [6]	86.01	87.15	68.34	60.35
MSCNN [34]	85.56	88.70	68.44	60.35
NLD [3]	86.28	85.30	67.96	70.80
AODNet [7]	86.75	87.50	70.05	74.75
DCPDN [9]	33.36	31.89	43.18	49.30
GDN [12]	29.59	22.89	-	-
DFF [13]	31.24	25.67	37.27	34.25
MSTN (Ours)	28.76	21.56	35.74	32.55

images are included in the HSTS. It is worth noting that the real real-world hazy images in these datasets can be used to verify the dehazing ability of MSTN in real scenes.

2) *Middlebury Stereo Dataset*: Middlebury [36] is a high-resolution stereo indoor dataset with subpixel-accurate ground truth. Similar to ITS, the atmospheric scattering model is adopted to generate synthetic hazy images. Considering its high-resolution characteristics, we adopt Middlebury as an assistant testing dataset to demonstrate the robustness of our proposed MSTN. In the experiment, we use the model pre-trained on ITS and directly applied the model to the Middlebury dataset to show the model performance.

3) *NH-HAZE*: NH-HAZE dataset [37] was proposed in the NTIRE2020 Image Dehazing Challenge [42], which is a non-homogeneous realistic dataset that contains 55 outdoor scenes. In NH-HAZE, the haze was be introduced in the scene by using a professional haze generator, which can simulates the real conditions of hazy scenes. Moreover, the hazy and haze-free corresponding scenes contain the same visual content captured under the same illumination parameters. Following the challenge setting, we use images 1 ~ 50 as the training dataset and images 51 ~ 55 as the testing dataset.

4) *DID-MDN*: DID-MDN [38] is a derain dataset, which includes three different rain-density images, that is light, medium, and heavy rain-density, respectively. In DID-MDN, the training dataset includes 12,000 images and the test dataset includes 12,00 images.

### B. Implementation Details

**Model setting**: In the final version of MSTN, the value of  $i$  and  $j$  are set as 5, which means that MSTN has 5 branches and the maximum depth of the first branch is 5. This also means that the model contains 5 branches with different scales.

**Training setting**: During training, we use RGB image as input and augment the image with random rotation(90°, 180°, 270°), horizontal flip, and vertical flip. In addition, we randomly crop 240 × 240 image patch on the image as the input and set batch-size as 16. The initial learning rate is  $1 \times 10^{-4}$  and cosine annealing strategy [43] is applied to adjust the learning rate. We implement our model with the PyTorch framework and update it with Adam optimizer ( $5 \times 10^7$  iterations). All our experiments are performed on GTX TitanXP GPUs.

### C. Comparison with SOTA Image Dehazing Methods

We compare MSTN with 10 SOTA image dehazing methods, including DCP [1], CAP [2], DehazeNet [6], MSCNN [34], NLD [3], AODNet [7], DCPDN [9], GFN [35], GDN [12], and DFF [13]. In addition, we use PSNR, SSIM, SSEQ [44], and BLIINDS-II [45] to evaluate the quality of dehazed images. Among them, the larger the PSNR and SSIM value, the better the result. Contrary, the smaller the SSEQ and BLIINDS-II value, the better the result. It worth noting that SSEQ and BLIINDS-II are no-reference image quality assessment methods, which can effectively reflect the visual effect of the reconstructed images.

1) *Quantitative Comparison on Synthesis Images*: In this work, we trained two different versions of the model for indoor and outdoor scenarios. These two models were trained on ITS and OTS dastates, respectively. In TABLE I, we provide

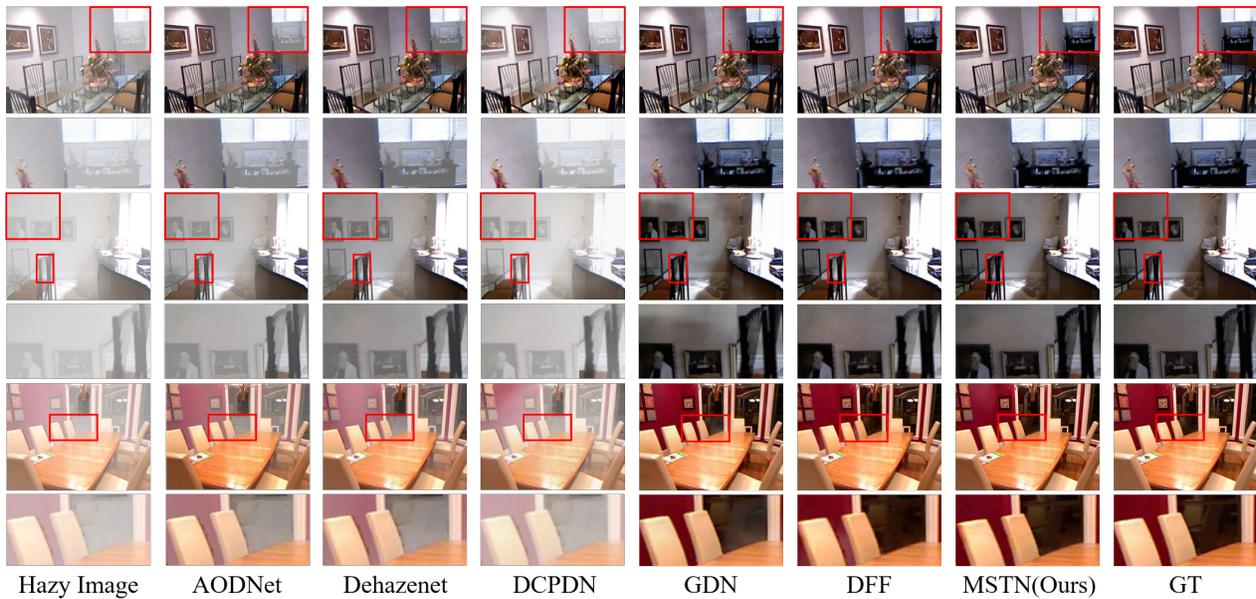


Fig. 5. Visual comparison with SOTA image dehazing methods on the RESIDE-SOTS [33] (Indoor) dataset. **Please zoom in to view details.**

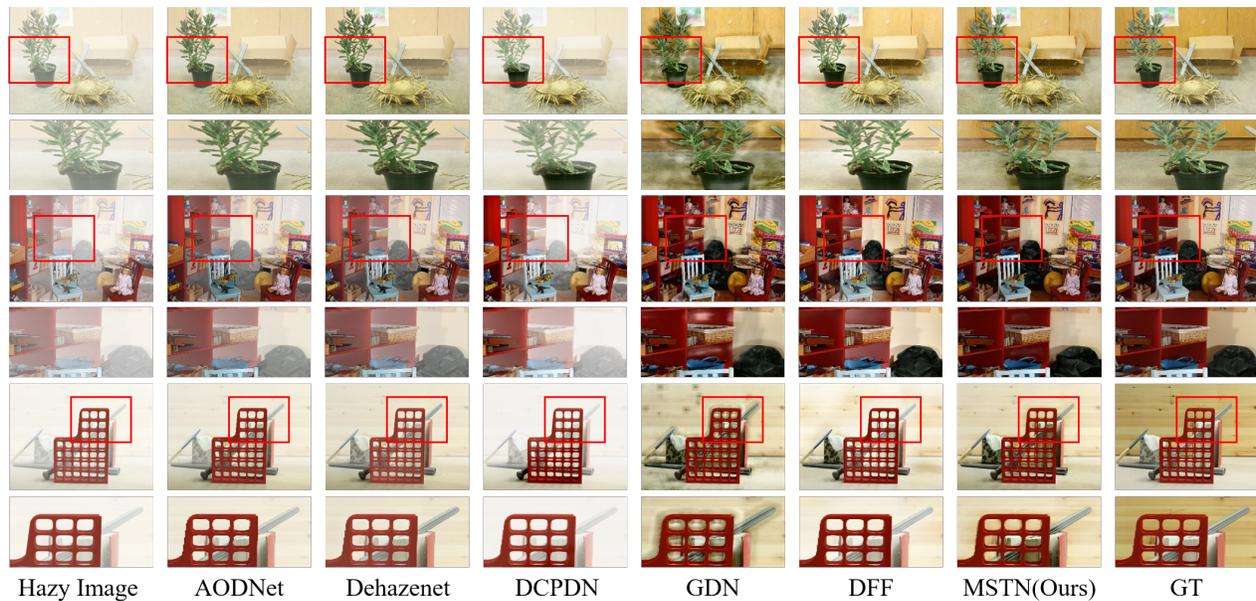


Fig. 6. Visual comparison with the SOTA image dehazing methods on the MiddleBury [36] dataset. **Please zoom in to view details.**

the PSNR/SSIM comparisons with SOTA image dehazing methods on SOTS [33] (Indoor and Outdoor). Obviously, our MSTN achieves the best results whether in the indoor or outdoor scenes. Among these methods, GDN [12] and DFF [13] are the latest methods and achieved the SOTA results at the time. Despite this, compared with them, our MSTN still achieved better results with a great advantage. Specifically, compared to the second-best model, the average PSNR results of MSTN in Indoor and Outdoor scenarios has increased 1.62dB and 0.40dB, respectively. This is a significant improvement and provide a new SOTA results on the image dehazing task. This is because the proposed multi-

scale topological architecture can extract rich features from the input image, so that the model can reconstruct high-quality haze-free images. Meanwhile, in order to verify the generalization ability of the model, we directly use the pre-trained MSTN on OTS and ITS to reconstruct haze-free images on HSTS [33] (Synthetic) and MiddleBury [36], respectively. According to TABLE I and II, we can observe that our MSTN still achieves the best results on all of these two datasets. It is worth mentioning that compared to the second-best model, the results of MSTN on these two datasets are improved by 1.68dB and 2.73dB, respectively. Moreover, we provide the SSEQ and BLINDS-II comparison of these methods on HSTS [33]

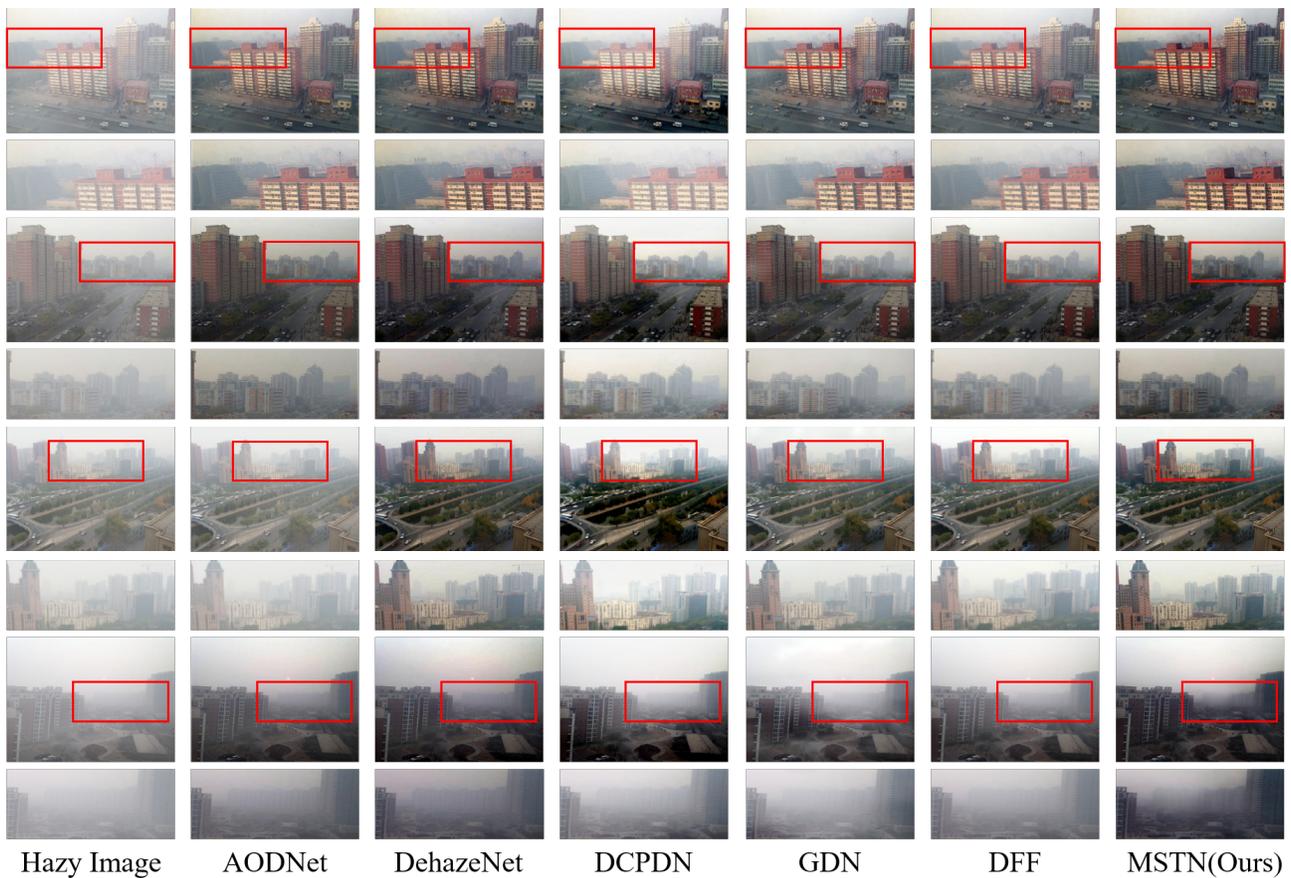


Fig. 7. Visual comparison with SOTA image dehazing methods on the RESIDE-SOTS [33] (Outdoor) dataset. **Please zoom in to view details.**

(Synthetic) in TABLE III. Obviously, our MSTN still achieves the best results. This further verified the effectiveness and powerful generalization capabilities of MSTN.

2) **Visual Comparison on Synthesis Images:** In Figs. 5, 6, and 7, we show the visual comparison with other image dehazing on SOTS [33] (Indoor and Outdoor) and MiddleBury [36] datasets. Among them, the images in the SOTS (Indoor) and MiddleBury datasets contain relatively low haze density while the images in the SOTS (Outdoor) dataset contain high haze concentrations. According to Figs. 5 and 6, we can clearly observe that the image reconstructed by AODNet, Dehazenet, and DCPDN still contains a lot of haze. Compared with these methods, GDN and DFF can reconstruct more clear haze-free images. However, carefully observing these reconstructed images, we find that these images contain a lot of artifacts and false edges, especially on walls, doors and flat areas. This greatly limits the promotion and application of these models. On the contrary, our MSTN can reconstruct high-quality haze-free images without artifacts. In Fig. 7, we show the dehazing effect of the model in the outdoor scenes. Obviously, outdoor scenes have higher concentrations of haze and the distribution of these haze is uneven. Therefore, it is more challenging to recover haze-free images from these images. According to the figure, we can found that all of the compared methods are failed to restoration high-quality images. Compared with these methods, our MSTN can reconstruct more clear images.

Although the image reconstructed by our MSTN also contains some haze, the performance of MSTN has been greatly improved compared with the previous methods. This fully demonstrates the effectiveness of MSTN.

3) **Results on Real-World Images:** The concentration and distribution of haze in natural scenes are more diverse and complex than the simulated images. Therefore, the task of real image dehazing is more difficult. In this part, real hazy image datasets (RESIDE-HSTS [33] (Real-world) and NH-HAZE [37]) are used to further assess the practicality of our MSTN. For HSTS (Real-world), following the setup in RESIDE, we use the model pre-trained on the OTS to reconstruct haze-free images. For NH-HAZE, we follow the setup of the NTIRE 2020 Challenge on Image Dehazing to train and test our model. The qualitative results of these two datasets are presented in TABLE III (right) and II (bottom), respectively. According to the results, we can clearly observe that MSTN achieves the best results in all evaluation indicators. In addition, we provide some dehazing results on real-world images in Fig. 8. According to the figure, we can found that (i). The haze-free images reconstructed by other methods still contains varying degrees of haze; (ii). The haze-free images reconstructed by DCPDN are over-exposed, causing the color of the reconstructed image to deviate; 3) The haze-free images reconstructed by GDN contains a lot of artifacts and the distribution of haze is uneven. All of these phenomena expose

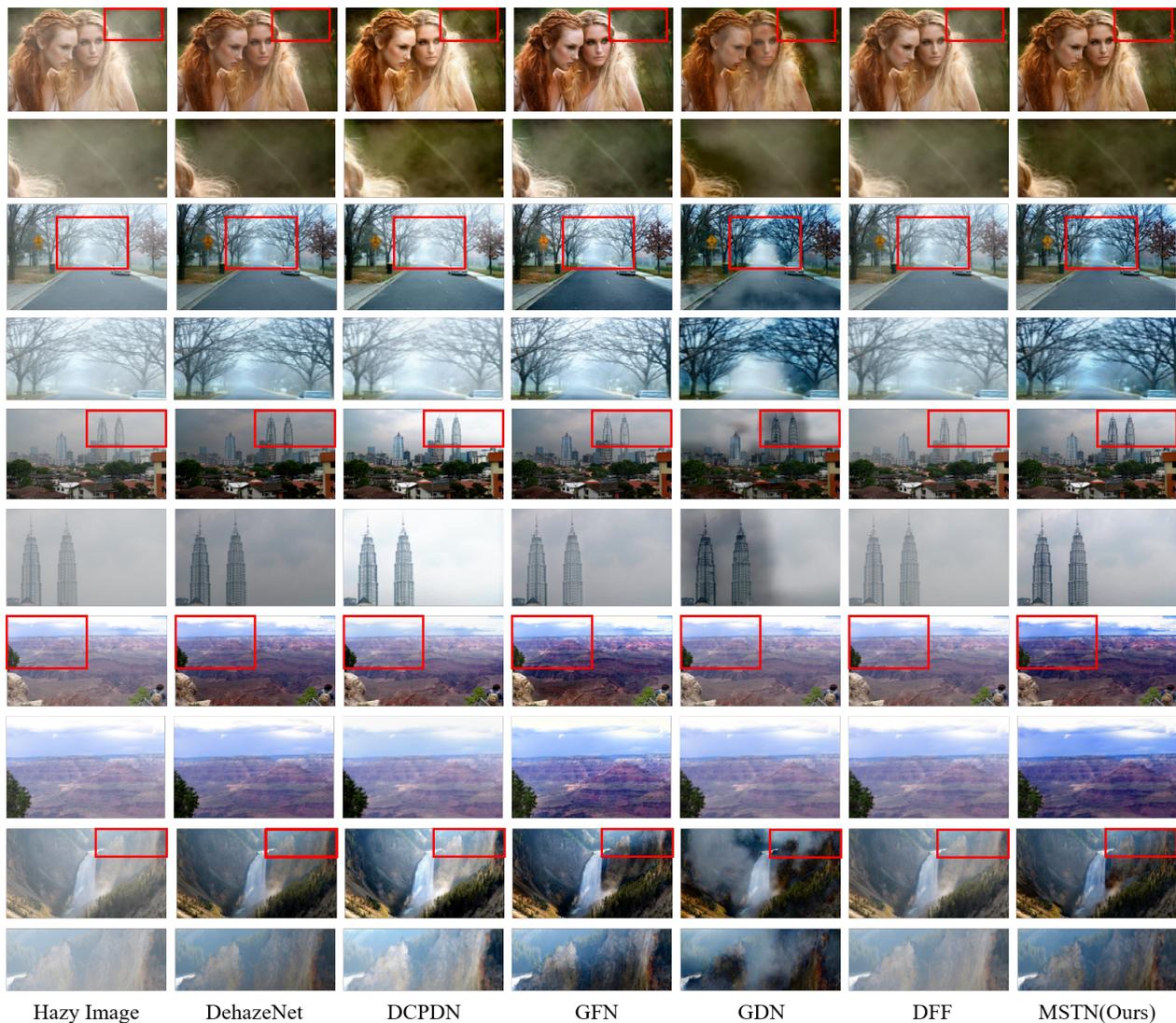


Fig. 8. Visual comparison with SOTA image dehazing methods on real-world hazy images. **Please zoom in to view details.**

the flaws of these models. In contrast, MSTN can reconstruct more clear and realistic haze-free images, which further proves the effectiveness of MSTN in practical applications.

## V. ANALYSIS AND DISCUSSION

### A. Study of Model Architecture

In this article, we propose a Multi-scale Topological Network (MSTN) for image dehazing. In order to study the effectiveness of the proposed architecture, we provide a series of ablation studies in this section. It is worth noting that in order to quickly verify the effectiveness of each module, the training settings in this section are as follows: batch size = 8, patch size =  $128 \times 128$ , and  $1 \times 10^6$  iterations.

1) **Effectiveness of AFSM:** AFSM is designed to select and fuse different image features, which also serves as the core component of MFFM. In order to verify the effectiveness of AFSM, we designed two simplified model, named MSTN (baseline) and MSTN (w/o AFSM). Among them, MSTN

TABLE IV  
STUDY OF AFSM AND MFFM ON RESIDE-SOTS [33] (INDOOR).

Methods	MSTN (Baseline)	MSTN (w/o AFSM)	MSTN (w/o MFFM)
PSNR	<b>31.37</b>	31.02	31.03
SSIM	<b>0.975</b>	0.973	0.974

(baseline) was restrained by the new training settings and MSTN (w/o AFSM) has the same architecture with MSTN (baseline) but replaces all AFSMs in the model with the element-wise addition operation. According to TABLE IV and Fig. 9, we can clearly observe that when the element-wise addition operation is used to replace the AFSM, the PSNR result of the model drops by 0.35dB. Meanwhile, the modified model is unstable during the training, which make the model difficult to converge.

2) **Effectiveness of MFFM:** As we mentioned before, MFFM is the core module of the proposed MSTN, which is de-

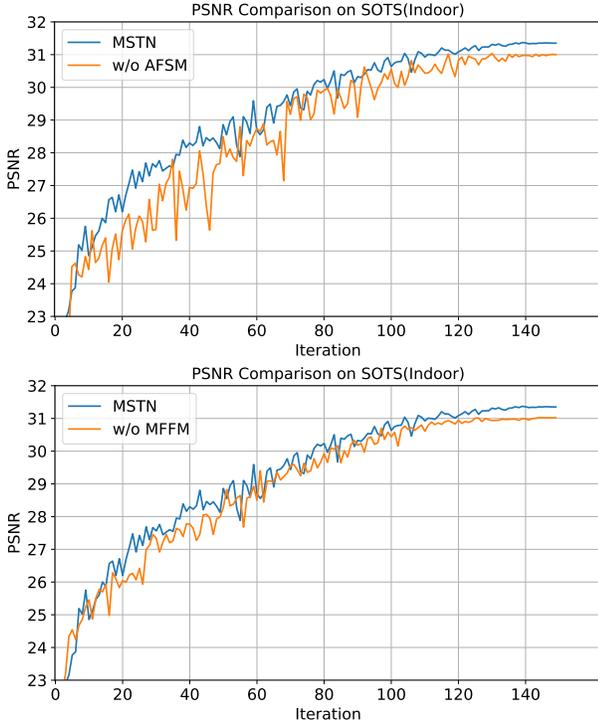


Fig. 9. Study on the effectiveness of AFSM and MFFM.

signed for multi-scale image feature selection, interaction, and fusion. According to the cross-scale skip connections, MFFM receives two feature maps with different scales as inputs and output the selected and fused features. In order to verify the effectiveness of MFFM, we designed a new model, named MSTN (w/o MFFM). MSTN (w/o MFFM) is a new model that remove all skip connections between different scales and replace all MFFMs in the MSTN (baseline) with residual blocks (RBs). In TABLE IV and Fig. 9, we provide the PSNR results and training curves of MSTN (baseline) and MSTN (w/o MFFM). According to the results, we can find that when MFFMs are replaced by RBs, the performance of the model drops by 0.34dB. This greatly illustrates the effectiveness of MFFM. Meanwhile, this illustrates the importance of multi-scale features and the rationality of MFFM design.

3) **Effectiveness of Multi-scale Architecture:** As shown in Fig. 2, MSTN adopts the pyramid-like structure to obtained multi-scale image features. In this paper, the final version of MSTN set  $i = 5$  and  $j = 5$ . This means that MSTN can extract image features with 5 different scales. In order to show the performance of the model under different scales, we designed a new set of models, and set  $i = 2, j = 2, i = 3, j = 3, i = 4, j = 4, i = 6, j = 6$ , respectively. This setting makes these models can extract different number of scales image features. In Fig. 10, we show the performance and parameters changes of these models. Obviously, the PSNR result increases as the scale number increases. Meanwhile, we can observe that when the number of scales continues to increase (such as  $i = 6$  and  $j = 6$ ), the model performance can be further improved. This means that the results reported in this paper are not the best results of MSTN. However, it cannot be ignored that the

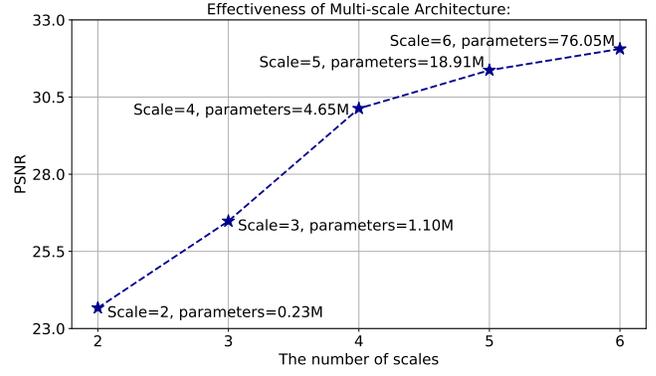


Fig. 10. Study the effectiveness of multi-scale architecture on SOTS (indoor).

TABLE V  
STUDY THE EFFECTIVENESS OF MULTI-SCALE ARCHITECTURE ON SOTS (INDOOR). THE BEST RESULT ARE HIGHLIGHTED.

PSNR	SSIM	Dark gray	Blue	Orange	Gray
23.05	0.881	✓			
27.03	0.943		✓		
29.56	0.951			✓	
<b>30.45</b>	<b>0.960</b>				✓

parameter quantity will increase as the scale number increases. Therefore, the number of scale can be selected according to actual demands. We set  $i = 5, j = 5$  in this paper to achieve a good balance between the model size and performance.

As shown in Fig. 2, we marked 4 roadmaps (Data Flow) on MSTN with different colors. This represents four simplified versions of MSTN with different structures. It is worth noting that, except for the modules marked with "data flow", the modules in these 4 models have been removed. Meanwhile, these 4 models all have 4 MFFMs. The only difference between these 4 models is they can extract different types of multi-scale image features. Specifically, the case 1 model (dark gray data flow) is a flat model, which can only extract image features with one scale. The case 4 model (gray data flow) is a multi-scale model, which can extract rich multi-scale image features. According to the TABLE V, we can clearly observe that when the model can extract more different

TABLE VI  
INVESTIGATIONS OF THE MODEL SIZE AND EXECUTION TIME.

Method	Param.	Platform	Times (s)
DCP	-	Matlab	1.532
CAP	-	matlab	0.808
DehazeNet	0.08M	Matlab	1.102
MSCNN	0.08	matlab	2.48
NLD	-	matlab	9.89
AodNet	0.02M	Mat-Caffe	0.402
GFN	0.51M	Mat-Caffe	1.373
DCPDN	66.89M	Pytorch	0.248
GDN	0.96M	Pytorch	0.0150
DFF	31.35M	Pytorch	0.0202
Ours	18.91M	Pytorch	<b>0.0139</b>

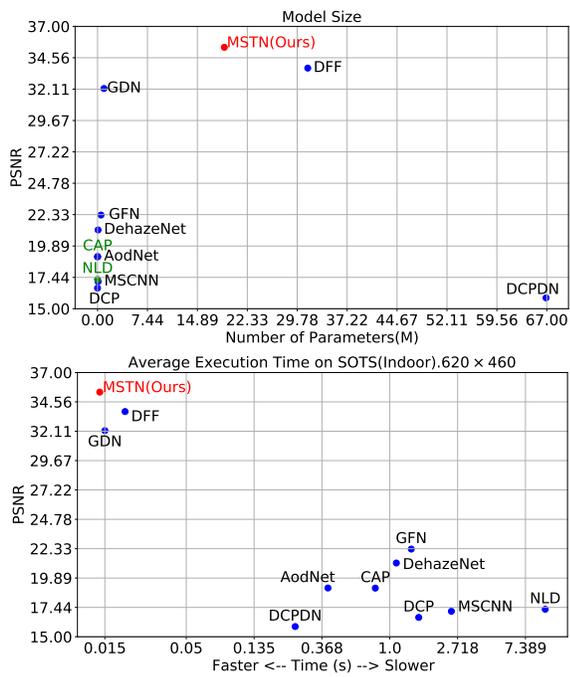


Fig. 11. Investigations of the model size and execution time.

scales image features, the model can achieve better results. All these experiments proved the importance of multi-scale image features and the effectiveness of the designed multi-scale architecture.

### B. Study of Model Model Size and Execution Time

Various large size image dehazing models have been proposed in recent years. These models always accompanied by numerous parameters, which means that these models require more storage space, computing resources, and execution time. In this paper, we aim to explore an efficient and accurate image dehazing model. Therefore, we need a more efficient network structure, not just increase the model parameters and depth. In TABLE VI we show the comparison of model parameters and execution time. Notice that all reported models use the released code and test on the same workstation. The time is the average time required for recovering 500 images of the size of  $620 \times 460$ . In Fig. 11, we intuitively display the comparison of model size, execution time, and performance of each models in the form of dot chart. According to the figure, we can draw the following conclusions: (1). Compared with lightweight models (e.g., DCP, MSCNN, AODNet, DehazeNet, GFN, GDN), the performance of MSTN is greatly improved; (2). Compared with large models (e.g., DFF and DCPDN), MSTN achieves better results with fewer parameters; (3). Compared with all reported image dehaizng models, MSTN achieves better results with less execution time. In summary, MSTN achieves a well balance between model performance, size, and execution time, which provide new solution for real-time image dehazing.

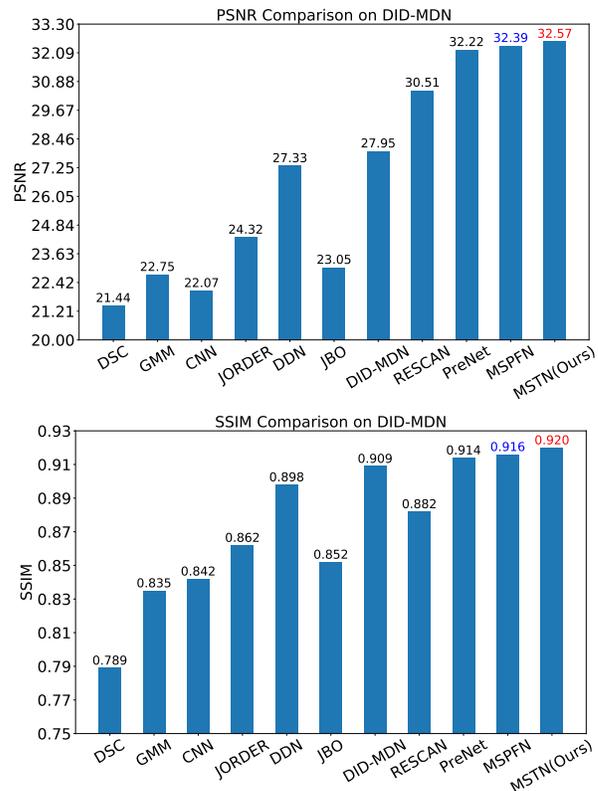


Fig. 12. Quantitative comparisons on DID-MDN [38]. The best and second best results are highlighted with red and blue fonts, respectively.

### C. Exploring on Other Image Restoration Task

In this paper, MSTN is proposed for the task of single image dehazing. According to our observation, MSTN is an efficient and accurate multi-scale topological network that can not only suitable for the image dehazing task. In order to explore the performance of MSTN on other image restoration tasks, we transfer MSTN to the task of single image deraining. Similar to image dehazing, the task of image deraining aims to reconstruct a clean image from the rain image. Following previous works, we use DID-MDN [38] to retrain our MSTN and compare it with 10 image deraining models, including DSC [46], GMM [47], CNN [48], JORDER [49], DDN [50], JBO [51], DID-MDN [38], RESCAN [52], PreNet [53], and MSPFN [54]. PSNR and SSIM results are provide in Fig. 12. According to the figure, we can clearly observe that MSTN achieves the best results in both PSNR and SSIM. This further proves the effectiveness of MSTN. This also means that MSTN is a highly scalable model that can be applied to other image restoration tasks. In future works, we will further verify the versatility and robustness of MSTN on other image restoration tasks like image desnowing and image denoising.

### D. Study on Hazy Images

In Fig. 7, we provide the image reconstructed by MSTN on the RESIDE-SOTS [33] (Outdoor) dataset. According to the figure, we can clearly observe that MSTN can reconstruct more clear images compared to other models. Moreover, we found that the image reconstructed by our MSTN is



Fig. 13. Study on hazy images (RESIDE-SOTS [33] (Outdoor)). Obviously, MSTN can reconstruct clear and high-quality haze-free images.

even clearer than the GT (Ground-Truth) image (Fig. 13). This phenomenon attracted our attention. Therefore, we re-investigated the RESIDE-SOTS [33] (Outdoor) dataset. This dataset contains 500 indoor hazy iamges and all these images are synthetic image. We investigated these images and found that these images contain different concentrations of haze itself. Therefore, part of the GT images in this dataset are hazy. However, there are plenty of clear GT images in the training dataset, thus the powerful learning ability of MSTN can learn how to reconstruct haze-free images from the hazy image. Therefore, our MSTN can reconstruct clearer images than GT images. This is because the dehazing ability learned by MSTN can remove the haze from the GT image itself. In Fig. 8, we provide the reconstruction results of MSTN on real hazy images. Obviously, MSTN can reconstruct high-quality haze-free images on real hazy images. This further proves the effectiveness and practicality of MSTN.

## VI. CONCLUSIONS

In this paper, we proposed an efficient and accurate Multi-scale Topological Network (MSTN) for single image dehaizng, which achieved competitive results on multiple datasets. MSTN adopts a new type of multi-scale topological architecture, which provides a large number of search paths and topological sub-networks that can fully extract image features from the input hazy image and improve the model stability and robustness. Meanwhile, we proposed a Multi-scale Feature Fusion Module (MFFM) and an Adaptive Feature Selection Module (AFSM) to realize the automatic transmission, selection, and fusion of multi-scale image features. Extensive experiments show that this special structure makes our model can extract rich image features to reconstruct high-quality haze-free images with texture details. Additionally, we achieved promising results by applying the model to other image restoration tasks such as image deraining. This further proves the effectiveness and versatility of the model. In future works, we will further verify the performance of the proposed model in more image restoration tasks.

## REFERENCES

- [1] K. He, J. Sun, and X. Tang, "Single image haze removal using dark channel prior," in *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, 2009, pp. 1956–1963.
- [2] Q. Zhu, J. Mai, and L. Shao, "A fast single image haze removal algorithm using color attenuation prior," *IEEE transactions on image processing*, vol. 24, no. 11, pp. 3522–3533, 2015.
- [3] D. Berman, S. Avidan *et al.*, "Non-local image dehazing," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 1674–1682.
- [4] M. Sulami, I. Glatzer, R. Fattal, and M. Werman, "Automatic recovery of the atmospheric light in hazy images," in *2014 IEEE International Conference on Computational Photography (ICCP)*. IEEE, 2014, pp. 1–11.
- [5] W. Wang, X. Yuan, X. Wu, and Y. Liu, "Fast image dehazing method based on linear transformation," *IEEE Transactions on Multimedia*, vol. 19, no. 6, pp. 1142–1155, 2017.
- [6] B. Cai, X. Xu, K. Jia, C. Qing, and D. Tao, "Dehazenet: An end-to-end system for single image haze removal," *IEEE Transactions on Image Processing*, vol. 25, no. 11, pp. 5187–5198, 2016.
- [7] B. Li, X. Peng, Z. Wang, J. Xu, and D. Feng, "Aod-net: All-in-one dehazing network," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 4770–4778.
- [8] K. Mei, A. Jiang, J. Li, and M. Wang, "Progressive feature fusion network for realistic image dehazing," in *Asian conference on computer vision*, 2018, pp. 203–215.
- [9] H. Zhang and V. M. Patel, "Densely connected pyramid dehazing network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 3194–3203.
- [10] Y. Qu, Y. Chen, J. Huang, and Y. Xie, "Enhanced pix2pix dehazing network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8160–8168.
- [11] C. Li, C. Guo, J. Guo, P. Han, H. Fu, and R. Cong, "Pdr-net: Perception-inspired single image dehazing network with refinement," *IEEE Transactions on Multimedia*, vol. 22, no. 3, pp. 704–716, 2020.
- [12] X. Liu, Y. Ma, Z. Shi, and J. Chen, "Griddehazenet: Attention-based multi-scale network for image dehazing," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 7314–7323.
- [13] H. Dong, J. Pan, L. Xiang, Z. Hu, X. Zhang, F. Wang, and M.-H. Yang, "Multi-scale boosted dehazing network with dense feature fusion," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2157–2167.
- [14] T.-Y. Lin, P. Dollár, R. Girshick, K. He, B. Hariharan, and S. Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [15] R. Fattal, "Single image dehazing," *ACM transactions on graphics (TOG)*, vol. 27, no. 3, pp. 1–9, 2008.
- [16] —, "Dehazing using color-lines," *ACM transactions on graphics (TOG)*, vol. 34, no. 1, pp. 1–14, 2014.
- [17] N. Attar and S. Aliakbary, "Classification of complex networks based on similarity of topological network features," *Chaos: An Interdisciplinary Journal of Nonlinear Science*, vol. 27, no. 9, p. 091102, 2017.

- [18] J. Li, Y. Yuan, K. Mei, and F. Fang, "Lightweight and accurate recursive fractal network for image super-resolution," in *Proceedings of the IEEE International Conference on Computer Vision Workshops*, 2019, pp. 0–0.
- [19] D. Zhang, H. Zhang, J. Tang, M. Wang, X. Hua, and Q. Sun, "Feature pyramid transformer," *Proceedings of the European Conference on Computer Vision*, 2020.
- [20] I. C. Duta, L. Liu, F. Zhu, and L. Shao, "Pyramidal convolution: Rethinking convolutional neural networks for visual recognition," *arXiv preprint arXiv:2006.11538*, 2020.
- [21] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [22] J. Li, F. Fang, K. Mei, and G. Zhang, "Multi-scale residual network for image super-resolution," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 517–532.
- [23] Z. Li, S. Li, N. Zhang, L. Wang, and Z. Xue, "Multi-scale invertible network for image super-resolution," in *Proceedings of the ACM Multimedia Asia*, 2019, pp. 1–6.
- [24] J. Qin, Y. Huang, and W. Wen, "Multi-scale feature fusion residual network for single image super-resolution," *Neurocomputing*, vol. 379, pp. 334–342, 2020.
- [25] J. Li, F. Fang, J. Li, K. Mei, and G. Zhang, "Mdcn: Multi-scale dense cross network for image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, 2020.
- [26] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [27] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.
- [28] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 510–519.
- [29] Z. Zhang, Q. Wu, Y. Wang, and F. Chen, "High-quality image captioning with fine-grained and semantic-guided visual attention," *IEEE Transactions on Multimedia*, vol. 21, no. 7, pp. 1681–1693, 2019.
- [30] D. Li, T. Yao, L. Duan, T. Mei, and Y. Rui, "Unified spatio-temporal attention networks for action recognition in videos," *IEEE Transactions on Multimedia*, vol. 21, no. 2, pp. 416–428, 2019.
- [31] J. Li, J. Li, F. Fang, F. Li, and G. Zhang, "Luminance-aware pyramid network for low-light image enhancement," *IEEE Transactions on Multimedia*, pp. 1–1, 2020.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [33] B. Li, W. Ren, D. Fu, D. Tao, D. Feng, W. Zeng, and Z. Wang, "Benchmarking single-image dehazing and beyond," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 492–505, 2018.
- [34] W. Ren, S. Liu, H. Zhang, J. Pan, X. Cao, and M.-H. Yang, "Single image dehazing via multi-scale convolutional neural networks," in *European conference on computer vision*. Springer, 2016, pp. 154–169.
- [35] W. Ren, L. Ma, J. Zhang, J. Pan, X. Cao, W. Liu, and M.-H. Yang, "Gated fusion network for single image dehazing," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3253–3261.
- [36] D. Scharstein, H. Hirschmüller, Y. Kitajima, G. Krathwohl, N. Nešić, X. Wang, and P. Westling, "High-resolution stereo datasets with subpixel-accurate ground truth," in *German conference on pattern recognition*, 2014, pp. 31–42.
- [40] D. Scharstein and R. Szeliski, "High-accuracy stereo depth maps using structured light," in *2003 IEEE Computer Society Conference on Computer Vision and Pattern Recognition, 2003. Proceedings.*, vol. 1, 2003, pp. I–I.
- [37] C. O. Ancuti, C. Ancuti, and R. Timofte, "Nh-haze: An image dehazing benchmark with non-homogeneous hazy and haze-free images," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 444–445.
- [38] H. Zhang and V. M. Patel, "Density-aware single image de-raining using a multi-stream dense network," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 695–704.
- [39] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgb-d images," in *European conference on computer vision*, 2012, pp. 746–760.
- [41] F. Liu, C. Shen, G. Lin, and I. Reid, "Learning depth from single monocular images using deep convolutional neural fields," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 10, pp. 2024–2039, 2015.
- [42] C. O. Ancuti, C. Ancuti, F.-A. Vasluianu, and R. Timofte, "Ntire 2020 challenge on nonhomogeneous dehazing," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 490–491.
- [43] T. He, Z. Zhang, H. Zhang, Z. Zhang, J. Xie, and M. Li, "Bag of tricks for image classification with convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 558–567.
- [44] L. Liu, B. Liu, H. Huang, and A. C. Bovik, "No-reference image quality assessment based on spatial and spectral entropies," *Signal Processing: Image Communication*, vol. 29, no. 8, pp. 856–863, 2014.
- [45] M. A. Saad, A. C. Bovik, and C. Charrier, "Blind image quality assessment: A natural scene statistics approach in the dct domain," *IEEE transactions on Image Processing*, vol. 21, no. 8, pp. 3339–3352, 2012.
- [46] Y. Luo, Y. Xu, and H. Ji, "Removing rain from a single image via discriminative sparse coding," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 3397–3405.
- [47] Y. Li, R. T. Tan, X. Guo, J. Lu, and M. S. Brown, "Rain streak removal using layer priors," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2736–2744.
- [48] X. Fu, J. Huang, X. Ding, Y. Liao, and J. Paisley, "Clearing the skies: A deep network architecture for single-image rain removal," *IEEE Transactions on Image Processing*, vol. 26, no. 6, pp. 2944–2956, 2017.
- [49] W. Yang, R. T. Tan, J. Feng, J. Liu, Z. Guo, and S. Yan, "Deep joint rain detection and removal from a single image," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 1357–1366.
- [50] X. Fu, J. Huang, D. Zeng, Y. Huang, X. Ding, and J. Paisley, "Removing rain from single images via a deep detail network," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 3855–3863.
- [51] L. Zhu, C.-W. Fu, D. Lischinski, and P.-A. Heng, "Joint bi-layer optimization for single-image rain streak removal," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 2526–2534.
- [52] X. Li, J. Wu, Z. Lin, H. Liu, and H. Zha, "Recurrent squeeze-and-excitation context aggregation net for single image deraining," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 254–269.
- [53] D. Ren, W. Zuo, Q. Hu, P. Zhu, and D. Meng, "Progressive image deraining networks: A better and simpler baseline," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 3937–3946.
- [54] K. Jiang, Z. Wang, P. Yi, C. Chen, B. Huang, Y. Luo, J. Ma, and J. Jiang, "Multi-scale progressive fusion network for single image deraining," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 8346–8355.