

MBZUAI

Digital.Commons@MBZUAI

Computer Vision Faculty Publications

Scholarly Works

4-19-2022

Learning Enriched Features for Fast Image Restoration and Enhancement

Syed Waqas Zamir

Inception Institute of Artificial Intelligence, United Arab Emirates

Aditya Arora

Inception Institute of Artificial Intelligence, United Arab Emirates

Salman Khan

Mohamed bin Zayed University of Artificial Intelligence

Munawar Hayat

University of California at Merced, and Google, United States & Terminus Group, China

Fahad Shahbaz Khan

Mohamed bin Zayed University of Artificial Intelligence

See next page for additional authors
See this and additional works at: <https://dclibrary.mbzuai.ac.ae/cvfp>



Part of the [Artificial Intelligence and Robotics Commons](#)

Preprint: arXiv

Archived with thanks to arXiv

Preprint License: CC by NC SA 4.0

Uploaded 25 August 2022

Recommended Citation

S.W. Zamir et al, "Learning Enriched Features for Fast Image Restoration and Enhancement", 2022, arXiv:2205.01649

This Article is brought to you for free and open access by the Scholarly Works at Digital.Commons@MBZUAI. It has been accepted for inclusion in Computer Vision Faculty Publications by an authorized administrator of Digital.Commons@MBZUAI. For more information, please contact libraryservices@mbzuai.ac.ae.

Authors

Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao

Learning Enriched Features for Fast Image Restoration and Enhancement

Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat,
Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao

Abstract—Given a degraded input image, image restoration aims to recover the missing high-quality image content. Numerous applications demand effective image restoration, e.g., computational photography, surveillance, autonomous vehicles, and remote sensing. Significant advances in image restoration have been made in recent years, dominated by convolutional neural networks (CNNs). The widely-used CNN-based methods typically operate either on full-resolution or on progressively low-resolution representations. In the former case, spatial details are preserved but the contextual information cannot be precisely encoded. In the latter case, generated outputs are semantically reliable but spatially less accurate. This paper presents a new architecture with a holistic goal of maintaining spatially-precise high-resolution representations through the entire network, and receiving complementary contextual information from the low-resolution representations. The core of our approach is a multi-scale residual block containing the following key elements: (a) parallel multi-resolution convolution streams for extracting multi-scale features, (b) information exchange across the multi-resolution streams, (c) non-local attention mechanism for capturing contextual information, and (d) attention based multi-scale feature aggregation. Our approach learns an enriched set of features that combines contextual information from multiple scales, while simultaneously preserving the high-resolution spatial details. Extensive experiments on six real image benchmark datasets demonstrate that our method, named as MIRNet-v2, achieves state-of-the-art results for a variety of image processing tasks, including defocus deblurring, image denoising, super-resolution, and image enhancement. The source code and pre-trained models are available at <https://github.com/swz30/MIRNetv2>.

Index Terms—Multi-scale Feature Representation, Dual-pixel Defocus Deblurring, Image Denoising, Super-resolution, Low-light Image Enhancement, and Contrast Enhancement

1 INTRODUCTION

Owing to the physical limitations of cameras or due to complicated lighting conditions, image degradations of varying severity are often introduced as part of image acquisition. For instance, smartphone cameras come with a narrow aperture and have small sensors with limited dynamic range. Consequently, they frequently generate noisy and low-contrast images. Similarly, images captured under the unsuitable lighting are either too dark or too bright. Image restoration aims to recover the original clean image from its corrupted measurements. It is an ill-posed inverse problem, due to the existence of many possible solutions.

Recent advances in image restoration and enhancement have been led by deep learning models, as they can learn strong (generalizable) priors from large-scale datasets. Existing CNNs typically follow one of the two architecture designs: 1) an encoder-decoder, or 2) high-resolution (single-scale) feature processing. The encoder-decoder models [1], [2], [3], [4] first progressively map the input to a low-resolution representation, and then apply a gradual reverse mapping to the original resolution. Although these approaches learn a broad context by spatial-resolution reduction, on the downside, the fine spatial details are lost, making it extremely hard to recover them in the later stages. On

the other hand, the high-resolution (single-scale) networks [5], [6], [7], [8] do not employ any downsampling operation, and thereby recover better spatial details. However, these networks have limited receptive field and are less effective in encoding contextual information.

Image restoration is a position-sensitive procedure, where pixel-to-pixel correspondence from the input image to the output image is needed. Therefore, it is important to remove only the undesired degraded image content, while carefully preserving the desired fine spatial details (such as true edges and texture). Such functionality for segregating the degraded content from the true signal can be better incorporated into CNNs with the help of large context, e.g., by enlarging the receptive field. Towards this goal, we develop a new *multi-scale* approach that maintains the original high-resolution features along the network hierarchy, thus minimizing the loss of precise spatial details. Simultaneously, our model encodes multi-scale context by using *parallel convolution streams* that process features at lower spatial resolutions. The multi-resolution parallel branches operate in a manner that is complementary to the main high-resolution branch, thereby providing us more precise and contextually enriched feature representations.

One main distinction between our method and the existing multi-scale image processing approaches is how we aggregate contextual information. The existing methods [11], [12], [13] process each scale in isolation. In contrast, we *progressively* exchange and fuse information from coarse-to-fine resolution-levels. Furthermore, different from existing methods that employ a simple concatenation or averaging

- S.W. Zamir, and A. Arora, are with Inception Institute of Artificial Intelligence, UAE. E-mail: waqas.zamir@inceptioniai.org
- S. Khan and F.S. Khan are with Mohammed Bin Zayed University of Artificial Intelligence, UAE.
- M. Hayat is with Monash University, Melbourne, Australia.
- M.-H. Yang is with University of California at Merced, and Google, USA.
- L. Shao is with Terminus Group, China.

TABLE 1: Comparison between MIRNet-v2 and MIRNet [9] under the same experimental settings for image denoising task on the SIDD benchmark dataset [10]. FLOPs and inference times are computed on an image of size 256×256 . When compared to MIRNet [9], MIRNet-v2 is more accurate, while being significantly lighter and faster.

	PSNR	Params (M)	FLOPs (B)	Convs	Activations (M)	Train Time (h)	Inference Time (ms)
MIRNet [9]	39.72	31.79	785	635	1270	139	142
MIRNet-v2 (Ours)	39.84	5.9 (81% ↓)	140 (82% ↓)	406 (36% ↓)	390 (69% ↓)	63 (55% ↓)	39 (72% ↓)

of features coming from multi-resolution branches, we introduce a new *selective kernel* fusion approach that dynamically selects the useful set of kernels from each branch representations using a self-attention mechanism. More importantly, the proposed fusion block combines features with varying receptive fields, while preserving their distinctive complementary characteristics.

The main contributions of this work include:

- A novel feature extraction model that obtains a complementary set of features across multiple spatial scales, while maintaining the original high-resolution features to preserve precise spatial details (Sec. 3).
- A regularly repeated mechanism for information exchange, where the features from coarse-to-fine resolution branches are progressively fused together for improved representation learning (Sec. 3.1).
- A new approach to fuse multi-scale features using a selective kernel network that dynamically combines variable receptive fields and faithfully preserves the original feature information at each spatial resolution (Sec. 3.1.1).

A preliminary version of this work has been published as a conference paper [9]. The MIRNet model [9] is expensive in terms of size and speed. In this work, we make several key modifications to MIRNet [9] that allow us to significantly reduce the computational cost while enhancing model performance (see Table 1). Specifically, in the proposed MIRNet-v2, (a) We demonstrate feature fusion only in the direction from low- to high-resolution streams performs best, and the information flow from high- to low-resolution branches can be removed to improve efficiency. (b) We replace the dual attention unit with a new residual contextual block (RCB). Furthermore, we introduce group convolutions in RCB that are capable of learning unique representations in each filter group, while being more resource efficient than standard convolutions. (c) We employ progressive learning to improve training speed: the network is trained on small image patches in the early epochs and on gradually large patches in the later training epochs. (d) We show the effectiveness of the proposed design on a new task of dual-pixel defocus deblurring [14] alongside the other image processing tasks of image denoising, super-resolution and image enhancement. Our MIRNet-v2 achieves state-of-the-results on *all* six datasets. Furthermore, we extensively evaluate our approach on practical challenges, such as generalization ability across datasets (Sec. 4)

In Table 1, we compare MIRNet-v2 with MIRNet [9] under the same training and inference settings. The results show that MIRNet-v2 is more accurate (improving PSNR from 39.72 dB to 39.84 dB), while reducing the number of parameters and FLOPs by $\sim 81\%$, convolutions by 36% , and

activations by 69% . Furthermore, the training and inference speed is increased by $2.2\times$ and $3.6\times$, respectively.

2 RELATED WORK

Rapidly growing image content necessitates the need to develop effective image restoration and enhancement algorithms. In this paper, we propose a new method capable of performing dual-pixel defocus deblurring, image denoising, super-resolution, and image enhancement. Unlike existing works for these problems, our approach processes features at the original resolution in order to preserve spatial details, while effectively fuses contextual information from multiple parallel branches. Next, we briefly describe the representative methods for each of the studied problems.

2.1 Dual-Pixel Defocus Deblurring

Images captured with wide camera aperture have shallow depth of field (DoF), where the scene regions that lie outside the DoF are out-of-focus. Given an image with defocus blur, the goal of defocus deblurring is to generate an all-in-focus image. Existing defocus deblurring approaches either directly deblur images [14], [15], [16], [17], or first estimate the defocus disparity map and then use it to guide the deblurring procedure [18], [19], [20]. Modern cameras are equipped with dual-pixel sensor that has two photodiodes at each pixel location, thereby generating two sub-aperture views. The phase difference between these views is useful in measuring the amount of defocus blur at each scene point. Recently, Abuolaim *et al.* [14] presented a dual-pixel deblurring dataset (DPDD) and a new method based on encoder-decoder design. In this paper our focus is also on deblurring images directly using the dual-pixel data as in [14], [16]. Previous defocus deblurring works [14], [16] employ the encoder-decoder that repeatedly uses the downsampling operation, thus causing significant fine detail loss. Whereas the architectural design of our approach enables preservation of desired textural details in the restored image.

2.2 Image Denoising

Classic denoising methods are mainly based on modifying transform coefficients [21], [22] or averaging neighborhood pixels [23], [24], [25]. Although the classical approaches perform well, the self-similarity [26] based algorithms, *e.g.*, NLM [27] and BM3D [28], demonstrate promising denoising performance. Numerous patch-based schemes that exploit redundancy (self-similarity) in images are later developed [29], [30], [31], [32]. Recently, deep learning models [6], [9], [33], [34], [35], [36], [37], [38], [39], [40], [41], [42], [43], [44] make significant advances in image denoising, yielding favorable results than those of the hand-crafted methods.

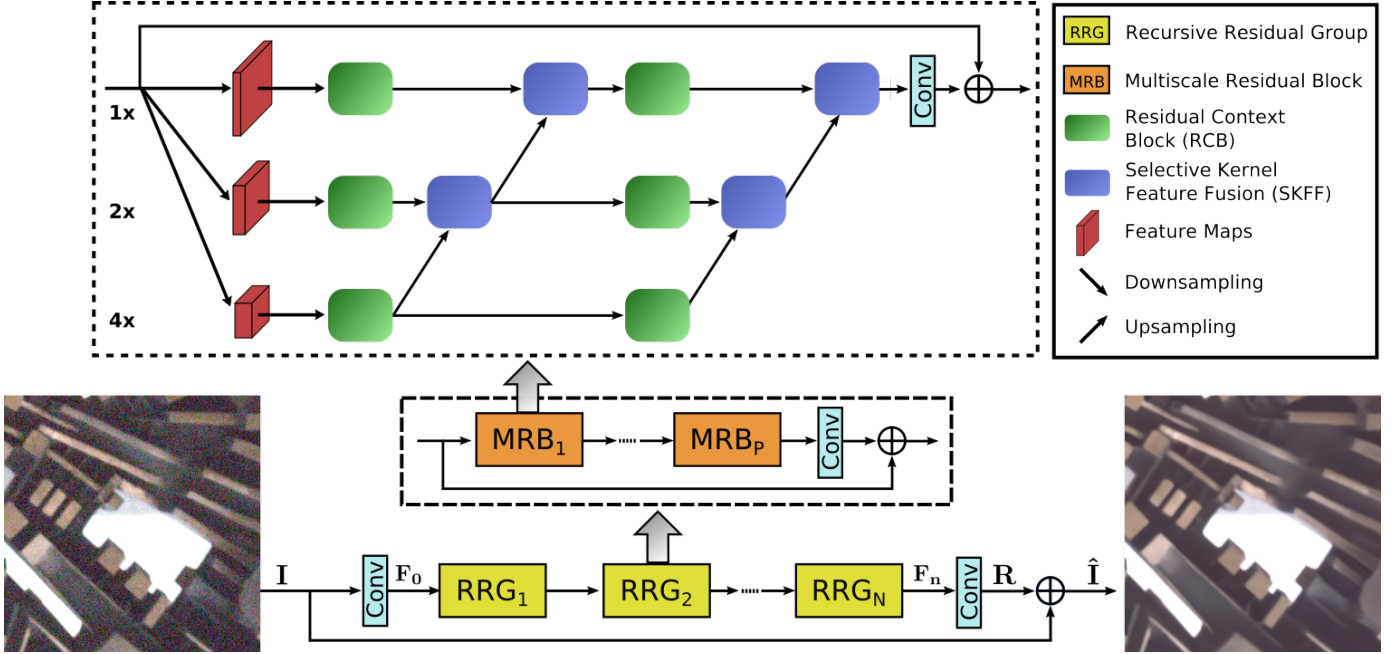


Fig. 1: Framework of the proposed MIRNet-v2 that learns enriched feature representations for image restoration and enhancement. MIRNet-v2 is based on a recursive residual design. In the core of MIRNet-v2 is the multi-scale residual block (MRB) whose main branch is dedicated to maintaining spatially-precise high-resolution representations through the entire network and the complimentary set of parallel branches provide better contextualized features.

2.3 Image Super-Resolution

Prior to the deep-learning era, numerous super-resolution (SR) algorithms have been proposed based on the sampling theory [45], [46], edge-guided interpolation [47], [48], natural image priors [49], [50], patch-exemplars [51], [52] and sparse representations [53], [54]. Currently, deep-learning techniques are being actively explored as they provide dramatically improved results over conventional algorithms. The data-driven SR approaches differ according to their architecture designs [55], [56], [57]. Early methods [5], [58] take a low-resolution (LR) image as input and learn to directly generate its high-resolution (HR) version. In contrast to directly producing a latent HR image, recent SR networks [59], [60], [61], [62] employ the residual learning framework [63] to learn the high-frequency image detail, which is later added to the input LR image to produce the final result. Other networks designed to perform SR include recursive learning [64], [65], [66], progressive reconstruction [67], [68], dense connections [7], [69], [70], attention mechanisms [71], [72], [73], multi-branch learning [68], [74], [75], [76], and generative adversarial networks (GANs) [70], [77], [78], [79].

2.4 Image Enhancement

Oftentimes, cameras generate images that lack vivid details or contrast. A number of factors contribute to the low quality of images, including unsuitable lighting conditions and physical limitations of camera devices. For image enhancement, histogram equalization is the most commonly used approach. However, it frequently produces under- or over-enhanced images. Motivated by the Retinex theory [80], several enhancement algorithms mimicking human vision have been proposed in the literature [81], [82], [83], [84].

Recently, CNNs have been successfully applied to general, as well as low-light, image enhancement problems [85]. Notable works employ Retinex-inspired networks [4], [86], [87], [88], encoder-decoder networks [89], [90], [91], [92], [93], and GANs [94], [95], [96].

3 PROPOSED METHOD

A schematic of the proposed MIRNet-v2 is shown in Fig. 1. We first present an overview of the proposed MIRNet-v2 for image restoration and enhancement. We then provide details of the *multi-scale residual block*, which is the fundamental building block of our method, containing several key elements: (a) parallel multi-resolution convolution streams for extracting (fine-to-coarse) semantically-rich and (coarse-to-fine) spatially-precise feature representations, (b) information exchange across multi-resolution streams, (c) attention-based aggregation of features arriving from different streams, and (d) residual contextual blocks to extract attention-based features.

Overall Pipeline. Given an image $I \in \mathbb{R}^{H \times W \times 3}$, the proposed model first applies a convolutional layer to extract low-level features $F_0 \in \mathbb{R}^{H \times W \times C}$. Next, the feature maps F_0 pass through N number of recursive residual groups (RRGs), yielding deep features $F_n \in \mathbb{R}^{H \times W \times C}$. We note that each RRG contains several multi-scale residual blocks, which is described in Section 3.1. Next, we apply a convolution layer to deep features F_n and obtain a residual image $R \in \mathbb{R}^{H \times W \times 3}$. Finally, the restored image is obtained as $\hat{I} = I + R$. We optimize the proposed network using the Charbonnier loss [97]:

$$\mathcal{L}(\hat{I}, I^*) = \sqrt{\|\hat{I} - I^*\|^2 + \varepsilon^2}, \quad (1)$$

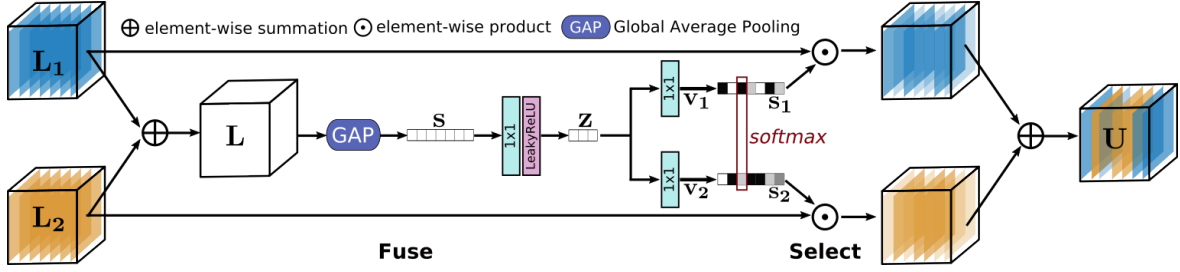


Fig. 2: Schematic for selective kernel feature fusion (SKFF). It operates on features from different resolution streams, and performs aggregation based on self-attention.

where \mathbf{I}^* denotes the ground-truth image, and ε is a constant which we empirically set to 10^{-3} for all the experiments.

3.1 Multi-Scale Residual Block

To encode context, existing CNNs [1], [98], [99], [100], [101], [102] typically employ the following architecture design: (a) the receptive field of neurons is fixed in *each* layer/stage, (b) the spatial size of feature maps is *gradually* reduced to generate a semantically strong low-resolution representation, and (c) a high-resolution representation is *gradually* recovered from the low-resolution representation. However, it is well-understood in vision science that in the primate visual cortex, the sizes of the local receptive fields of neurons in the same region are different [103], [104], [105], [106]. Therefore, a similar mechanism of collecting multi-scale spatial information in the same layer is more effective when incorporated with in CNNs [107], [108], [109], [110]. Motivated by this, we propose the multi-scale residual block (MRB), as shown in Fig. 1. It is capable of generating a spatially-precise output by maintaining high-resolution representations, while receiving rich contextual information from low-resolutions. The MRB consists of multiple (three in this paper) fully-convolutional streams connected in parallel that operate on varying resolution feature maps (ranging from low to high). It allows contextualized-information transfer from the low-resolution streams to consolidate the high-resolution features. Next, we describe the individual components of MRB.

3.1.1 Selective Kernel Feature Fusion

One fundamental property of neurons present in the visual cortex is their ability to change receptive fields according to the stimulus [111]. This mechanism of adaptively adjusting receptive fields can be incorporated in CNNs by using multi-scale feature generation (in the same layer) followed by feature aggregation and selection. The most commonly used approaches for feature aggregation include simple concatenation or summation. However, these choices provide limited expressive power to the network, as reported in [111]. In MRB, we introduce a nonlinear procedure for fusing features coming from different resolution streams using a self-attention mechanism. Motivated by [111], we call it selective kernel feature fusion (SKFF).

The SKFF module performs dynamic adjustment of receptive fields via two operations – *Fuse* and *Select*, as

illustrated in Fig. 2. The *fuse* operator generates global feature descriptors by combining the information from multi-resolution streams. The *select* operator uses these descriptors to recalibrate the feature maps (of different streams) followed by their aggregation. Next, we provide details of both operators. (1) **Fuse**: SKFF receives inputs from two parallel convolution streams carrying different scales of information. We first combine these multi-scale features using an element-wise sum as: $\mathbf{L} = \mathbf{L}_1 + \mathbf{L}_2$. We then apply global average pooling (GAP) across the spatial dimension of $\mathbf{L} \in \mathbb{R}^{H \times W \times C}$ to compute channel-wise statistics $\mathbf{s} \in \mathbb{R}^{1 \times 1 \times C}$. Next, we apply a channel-downscaling convolution layer to generate a compact feature representation $\mathbf{z} \in \mathbb{R}^{1 \times 1 \times r}$, where $r = \frac{C}{8}$ for all our experiments. Finally, the feature vector \mathbf{z} passes through two parallel channel-upscaling convolution layers (one for each resolution stream) and provides us with two feature descriptors \mathbf{v}_1 and \mathbf{v}_2 , each with dimensions $1 \times 1 \times C$. (2) **Select**: This operator applies the softmax function to \mathbf{v}_1 and \mathbf{v}_2 , yielding attention activations \mathbf{s}_1 and \mathbf{s}_2 that we use to adaptively recalibrate multi-scale feature maps \mathbf{L}_1 and \mathbf{L}_2 , respectively. The overall process of feature recalibration and aggregation is defined as: $\mathbf{U} = \mathbf{s}_1 \cdot \mathbf{L}_1 + \mathbf{s}_2 \cdot \mathbf{L}_2$. Note that the SKFF uses $\sim 5\times$ fewer parameters than aggregation with concatenation but generates more favorable results (an ablation study is provided in the experiments section).

3.1.2 Residual Contextual Block

While the SKFF block fuses information across multi-resolution branches, we also need a distillation mechanism to extract useful information from within a feature tensor. Motivated by the advances of recent low-level vision methods [33], [71], [72], [73] which incorporate attention mechanisms [112], [113], [114], [115], we propose the residual contextual block (RCB) to extract features in the convolutional streams. The schematic of RCB is shown in Fig. 3. The RCB suppresses less useful features and only allows more informative ones to pass further. The overall process of RCB is summarized as:

$$\mathbf{F}_{\text{RCB}} = \mathbf{F}_a + W(\text{CM}(\mathbf{F}_b)), \quad (2)$$

where $\mathbf{F}_b \in \mathbb{R}^{H \times W \times C}$ represents feature maps that are obtained by applying two 3×3 group convolution layers to the input features $\mathbf{F}_b \in \mathbb{R}^{H \times W \times C}$ at the beginning of the RCB. These group convolutions are more resource efficient than standard convolutions and capable of learning unique representations in each filter group. W denotes the

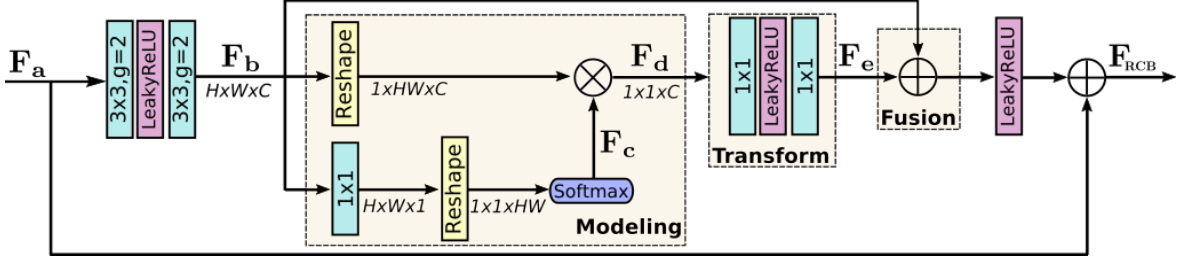


Fig. 3: Architecture of residual contextual block (RCB). In the first two group convolution layers, g represents the number of groups. \otimes denotes matrix multiplication.

last convolutional layer with filter size 1×1 . CM stands for contextual module that is realized in three parts. (1) **Context modeling:** From the original feature maps F_b , we first generate new features $F_c \in \mathbb{R}^{1 \times 1 \times HW}$ by applying 1×1 convolution followed by the reshaping and softmax operations. Next we reshape F_b to $\mathbb{R}^{1 \times HW \times C}$ and perform matrix multiplication with F_c to obtain the global feature descriptor $F_d \in \mathbb{R}^{1 \times 1 \times C}$. (2) **Feature transform:** To capture the inter-channel dependencies we pass the descriptor F_d through two 1×1 convolutions, resulting in new attention features $F_e \in \mathbb{R}^{1 \times 1 \times C}$. (3) **Feature fusion:** We employ element-wise addition operation to aggregate contextual features F_e to each position of the original features F_b .

3.2 Progressive Training Regime

When considering the image patch size for network training, there is a trade-off between the training speed and test-time accuracy [116], [117]. On large patches, CNNs capture fine image details to provide improved results, but they are slower to train. Whereas, training on small image patches is faster, but comes at the cost of accuracy drop. To strike the right balance between the training speed and accuracy, we propose a progressive learning method where the network is trained on smaller image patches in the early epochs and on gradually larger patches in the later training epochs. This approach can also be understood as a curriculum learning process where the network sequentially moves from learning a simpler task to a more complex one (where modeling of fine details is required). The progressive learning strategy on mixed-size image patches not only improves the training speed but also enhances the model performance at test time where the input images can be of different sizes (which is common in image restoration problems).

4 EXPERIMENTS

In this section, we perform qualitative and quantitative assessments of the results produced by our MIRNet-v2 and compare it with the state-of-the-art methods. Next, we describe the datasets, and then provide the implementation details. Finally, we report results for (a) dual-pixel defocus deblurring, (b) image denoising, (c) image super-resolution and (d) image enhancement, on six real image datasets.

4.1 Real Image Datasets

Dual-pixel defocus deblurring. DPDD [14] dataset contains 500 indoor/outdoor scenes captured with a DSLR

camera. Each scene consists of two defocus blurred sub-aperture views captured with a wide camera aperture, and the corresponding all-in-focus ground truth image captured with a narrow aperture. The DDPD dataset is divided into 350 images for training, 74 images for validation and 76 images for testing.

Image denoising. (1) **DND** [118] consists of 50 images captured with four consumer cameras. Since the images are of very high-resolution, the dataset providers extract 20 crops of size 512×512 from each image, yielding 1000 patches in total. All these patches are used for testing (as DND does not contain training or validation sets). The ground-truth noise-free images are not released publicly, therefore the image quality scores in terms of PSNR and SSIM can only be obtained through an online server [119]. (2) **SIDD** [10] is collected with smartphone cameras. Due to the small sensor and high-resolution, the noise levels in smartphone images are much higher than those of DSLRs. SIDD contains 320 image pairs for training and 1280 for validation.

Super-resolution. **RealSR** [120] contains real-world LR-HR image pairs of the same scene captured by adjusting the focal-length of the cameras. RealSR has both indoor and outdoor images taken with two cameras. The number of training image pairs for scale factors $\times 2$, $\times 3$ and $\times 4$ are 183, 234 and 178, respectively. For each scale factor, 30 test images are also provided in RealSR.

Image enhancement. (1) **LoL** [87] is created for low-light image enhancement problem. It provides 485 images for training and 15 for testing. Each image pair in LoL consists of a low-light input image and its corresponding well-exposed reference image. (2) **MIT-Adobe FiveK** [121] contains 5000 images of various indoor and outdoor scenes captured with DSLR cameras in different lighting conditions. The tonal attributes of all images are manually adjusted by five different trained photographers (labelled as experts A to E). Similar to [122], [123], [124], we also consider the enhanced images of expert C as the ground-truth. Moreover, the first 4500 images are used for training and the last 500 for testing.

4.2 Implementation Details

The proposed architecture is end-to-end trainable and requires no pre-training of sub-modules. We train four different networks for four different restoration tasks. For the dual-pixel defocus deblurring, we concatenate the left and right sub-aperture images and feed them as input to the network. The training parameters, common to all experiments, are the following. We use 4 RRGs, each of which further

TABLE 2: Dual-pixel Defocus Deblurring comparisons on the DPDD Dataset [14]. The test set of DPDD contains 37 indoor scenes and 39 outdoor scenes. Best and second best scores are **highlighted** and underlined, respectively.

Method	Indoor Scenes				Outdoor Scenes				Combined			
	PSNR \uparrow	SSIM \uparrow	MAE \downarrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	MAE \downarrow	LPIPS \downarrow	PSNR \uparrow	SSIM \uparrow	MAE \downarrow	LPIPS \downarrow
EBDB [18]	25.77	0.772	0.040	0.297	21.25	0.599	0.058	0.373	23.45	0.683	0.049	0.336
DMENet [20]	25.50	0.788	0.038	0.298	21.43	0.644	0.063	0.397	23.41	0.714	0.051	0.349
JNB [19]	26.73	0.828	0.031	0.273	21.10	0.608	0.064	0.355	23.84	0.715	0.048	0.315
DPDNet [14]	27.48	0.849	0.029	0.189	22.90	0.726	0.052	0.255	25.13	0.786	0.041	0.223
RDPD [16]	28.10	0.843	0.027	0.210	22.82	0.704	0.053	0.298	25.39	0.772	0.040	0.255
MIRNet-v2 (Ours)	28.96	0.881	0.024	0.154	23.59	0.753	0.049	0.205	26.20	0.816	0.037	0.180

contains 2 MRBs. The MRB has 3 parallel streams with channel dimensions of 80, 120, 180 at resolutions $1, \frac{1}{2}, \frac{1}{4}$, respectively. Each stream in MRB has 2 RCBs with shared parameters. The models are trained with the Adam optimizer ($\beta_1 = 0.9$, and $\beta_2 = 0.999$) for 3×10^5 iterations. The initial learning rate is set to 2×10^{-4} . We employ the cosine annealing strategy [125] to steadily decrease the learning rate from initial value to 10^{-6} during training. For progressive training, we use the image patch sizes of 128, 144, 192, and 224. The batch size is set to 64 and, for data augmentation, we perform horizontal and vertical flips.

4.3 Dual-Pixel Defocus Deblurring

We compare the performance of the proposed MIRNet-v2 with the conventional defocus deblurring methods (EBDB [18] and JNB [19]) as well as the learning-based approaches (DMENet [20], DPDNet [14], and RDPD [16]). Table 2 shows that our method achieves state-of-the-art results for both the indoor and outdoor scene categories. In particular, our MIRNet-v2 achieves 0.86 dB PSNR improvement over the previous best method RDPD [16] on indoor images and 0.77 dB on outdoor images. When both scene categories are combined, our method shows performance gains of 0.81 dB over RDPD [14] and 1.07 dB over the second best method DPDNet [14].

In Fig. 4, we provide defocus-deblurred results produced by different methods for both indoor and outdoor scenes. It is noticeable that our method effectively removes the spatially varying defocus blur and produces images that are more sharper and visually faithful to the ground-truth than those of the compared approaches.

4.4 Image Denoising

In this section, we demonstrate the effectiveness of the proposed MIRNet-v2 for image denoising. We train our network only on the training set of the SIDD [10] and directly evaluate it on the test images of both SIDD and DND [118] datasets. Quantitative comparisons in terms of PSNR and SSIM metrics are summarized in Table 3. Our MIRNet-v2 performs favourably against the data-driven, as well as conventional, denoising algorithms. Specifically, when compared to the recent best methods, our algorithm demonstrates a performance gain of 0.32 dB over CycleISP [38] on SIDD and 0.11 dB over DAGL [127] on DND. Furthermore, it is worth noting that CycleISP [38] uses additional training data, yet our method yields considerably better results.

TABLE 3: Denoising comparisons on SIDD [10] and DND [118] datasets. * indicates the methods that use additional training data. Whereas our MIRNet-v2 is only trained on the SIDD images and directly tested on DND.

Method	SIDD [10]		DND [118]	
	PSNR \uparrow	SSIM \uparrow	PSNR \uparrow	SSIM \uparrow
DnCNN [6]	23.66	0.583	32.43	0.790
MLP [126]	24.71	0.641	34.23	0.833
BM3D [28]	25.65	0.685	34.51	0.851
CBDNet* [35]	30.78	0.801	38.06	0.942
DAGL [127]	38.94	0.953	39.77	0.956
RIDNet* [33]	38.71	0.951	39.26	0.953
AINDNet* [42]	38.95	0.952	39.37	0.951
VDN [41]	39.28	0.956	39.38	0.952
DeamNet* [128]	39.47	0.957	39.63	0.953
SADNet* [39]	39.46	0.957	39.59	0.952
DANet+* [40]	39.47	0.957	39.58	0.955
CycleISP* [38]	39.52	0.957	39.56	0.956
MIRNet-v2 (Ours)	39.84	0.959	39.86	0.955

Fig. 5 shows a visual comparisons of our results with those of other competing algorithms. The MIRNet-v2 is effective in removing real noise and produces perceptually-pleasing and sharp images. Moreover, it is can maintain the spatial smoothness of the homogeneous regions without introducing artifacts. In contrast, most of the other methods either yield over-smooth images and thus sacrifice structural content and fine textural details, or produce images with chroma artifacts and blotchy texture.

Generalization capability. The DND and SIDD datasets are acquired with different sets of cameras having different noise characteristics. Since the DND benchmark does not provide training data, setting a new state-of-the-art on DND with our SIDD trained network indicates the good generalization capability of our approach.

4.5 Super-Resolution

We compare our MIRNet-v2 against the state-of-the-art SR algorithms (VDSR [59], SRResNet [79], RCAN [71], LP-KPN [120]) on the testing images of the RealSR [120] for upscaling factors of $\times 2$, $\times 3$ and $\times 4$. Note that all the benchmarked algorithms are trained on the RealSR [120] dataset for a fair comparison. In the experiments, we also include bicubic interpolation [45], which is the most commonly used method for generating super-resolved images. Here, we compute the PSNR and SSIM scores using the Y channel (in YCbCr color space), as it is a common practice

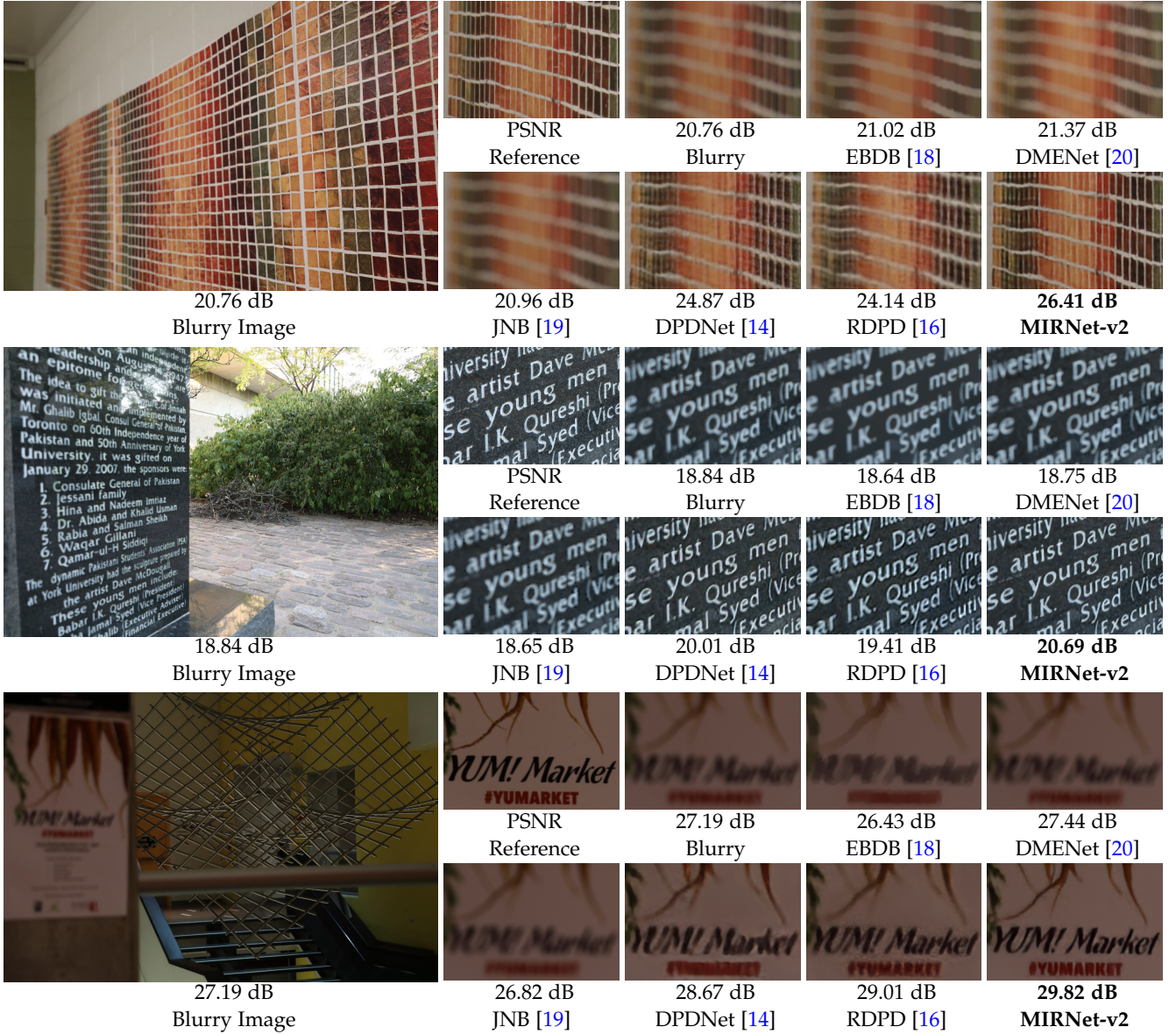


Fig. 4: Visual comparisons for dual-pixel defocus deblurring on the DPDD dataset [14]. Compared to the other approaches, our MIRNet-v2 more effectively removes blur while preserving the fine image details.

in the SR literature [55], [56], [71], [120]. The results in Table 4 show that the bicubic interpolation provides the least accurate results, thereby indicating its low suitability for dealing with real images. Moreover, the same table shows that the recent method LP-KPN [120] achieves marginal improvement of only ~ 0.04 dB over the previous best method RCAN [71]. In contrast, our method significantly advances state-of-the-art and consistently achieves better image quality scores than other approaches for all three scaling factors. Particularly, compared to LP-KPN [120], our method leads to performance gains of 0.48 dB, 0.73 dB, and 0.24 dB for scaling factors $\times 2$, $\times 3$ and $\times 4$, respectively. The trend is similar for the SSIM metric as well.

Visual comparisons in Fig. 6 show that our MIRNet-v2 can effectively recover content structures. In contrast, VDSR [59], SRResNet [79] and RCAN [71] reproduce results with noticeable artifacts. Furthermore, LP-KPN [120] is not able to preserve structures (see near the right edge of the

TABLE 4: Super-resolution evaluation on the RealSR dataset [120]. Compared to the state-of-the-art, our method consistently yields significantly better image quality scores for all three scaling factors.

Method	Scale $\times 2$		Scale $\times 3$		Scale $\times 4$	
	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM
Bicubic	32.61	0.907	29.34	0.841	27.99	0.806
VDSR [59]	33.64	0.917	30.14	0.856	28.63	0.821
SRResNet [79]	33.69	0.919	30.18	0.859	28.67	0.824
RCAN [71]	33.87	0.922	30.40	0.862	28.88	0.826
LP-KPN [120]	33.90	0.927	30.42	0.868	28.92	0.834
MIRNet-v2 (Ours)	34.38	0.934	31.15	0.883	29.16	0.845

crop). Several more examples are provided in Fig. 7 to further compare the image reproduction quality of our method against the previous best method [120]. It can be seen that LP-KPN [120] has a tendency to over-enhance the contrast

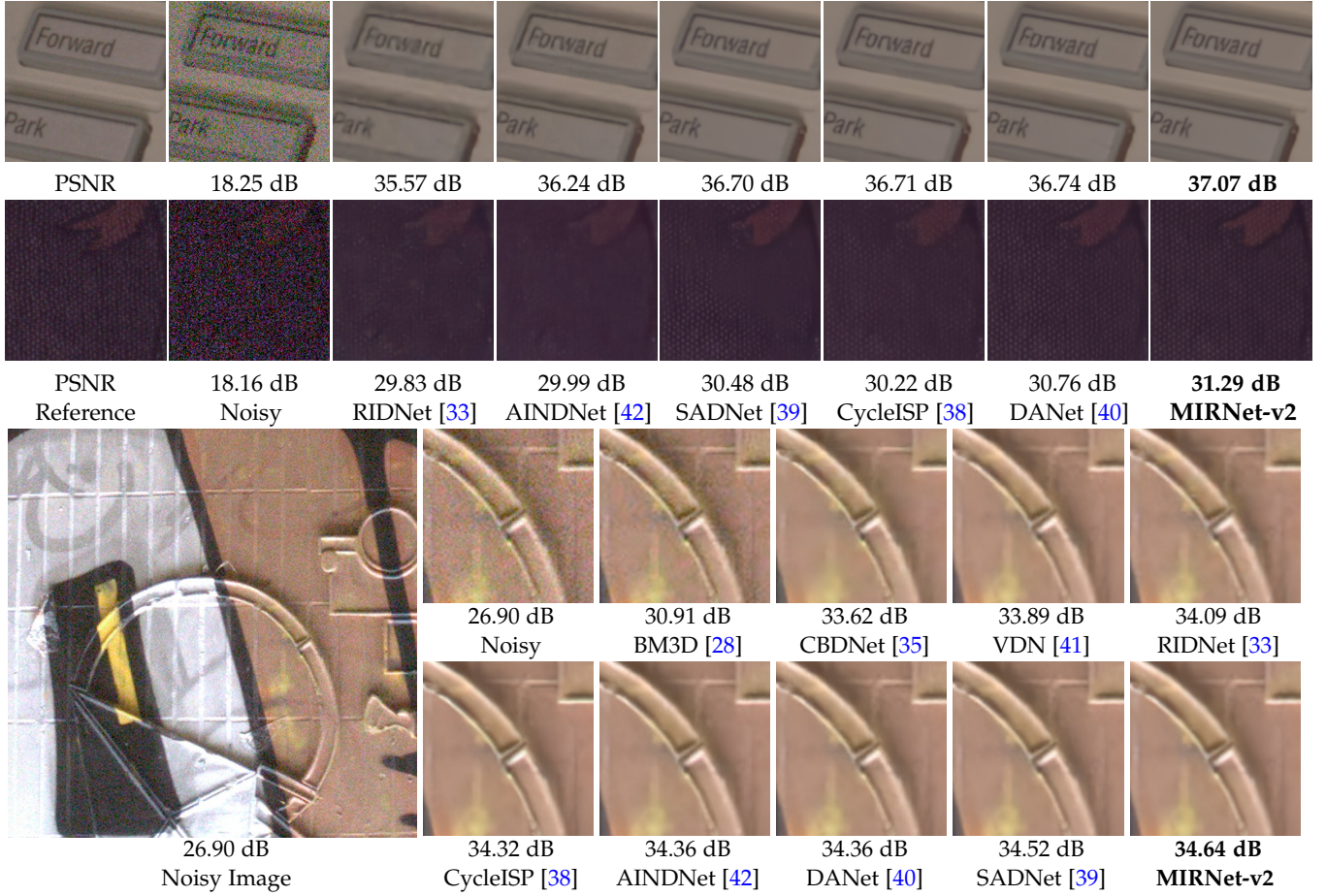


Fig. 5: Image denoising comparisons. First two examples are from SIDD [10] and the last is from DND [118]. The proposed MIRNet-v2 better preserves fine texture and structural patterns in the denoised images.

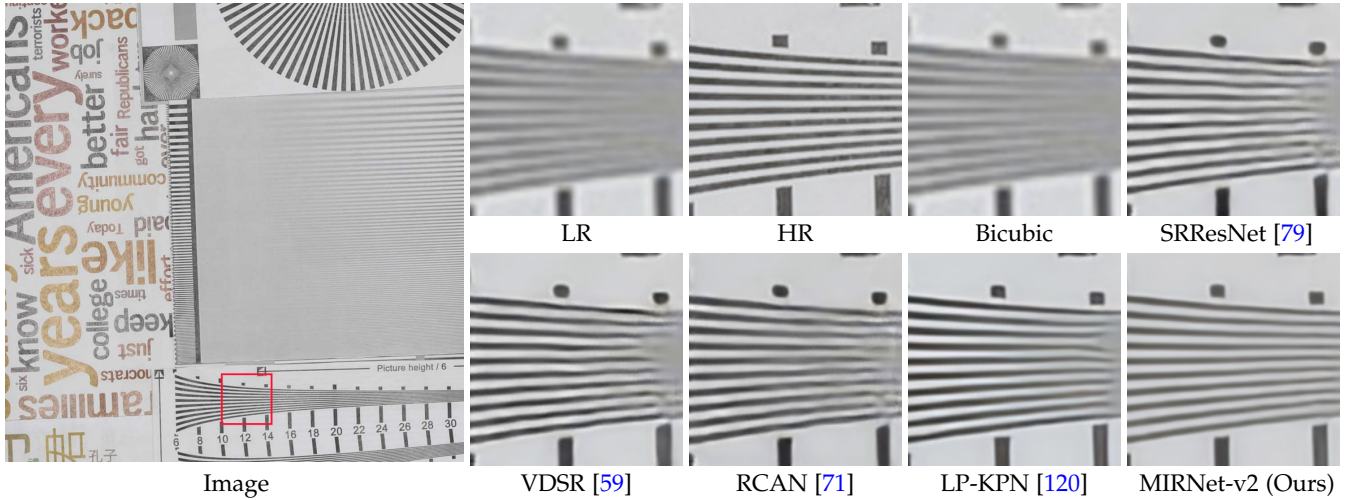


Fig. 6: Comparisons for $\times 4$ super-resolution on the RealSR [120] dataset. The image produced by our MIRNet-v2 is more faithful to the ground-truth than other competing methods (see lines near the right edge of the crops).

(cols. 1, 3, 4) and in turn causes loss of details near dark and high-light areas. In contrast, the proposed MIRNet-v2 successfully reconstructs structural patterns and edges (col. 2) and produces images that are natural (cols. 1, 4) and have better color reproduction (col. 5).

4.6 Image Enhancement

In this section, we demonstrate the effectiveness of our algorithm by evaluating it for the image enhancement task. We report PSNR/SSIM values of our method and several other techniques in Table 5 and Table 6 for the LoL [87] and MIT-Adobe FiveK [121] datasets, respectively. It can be seen

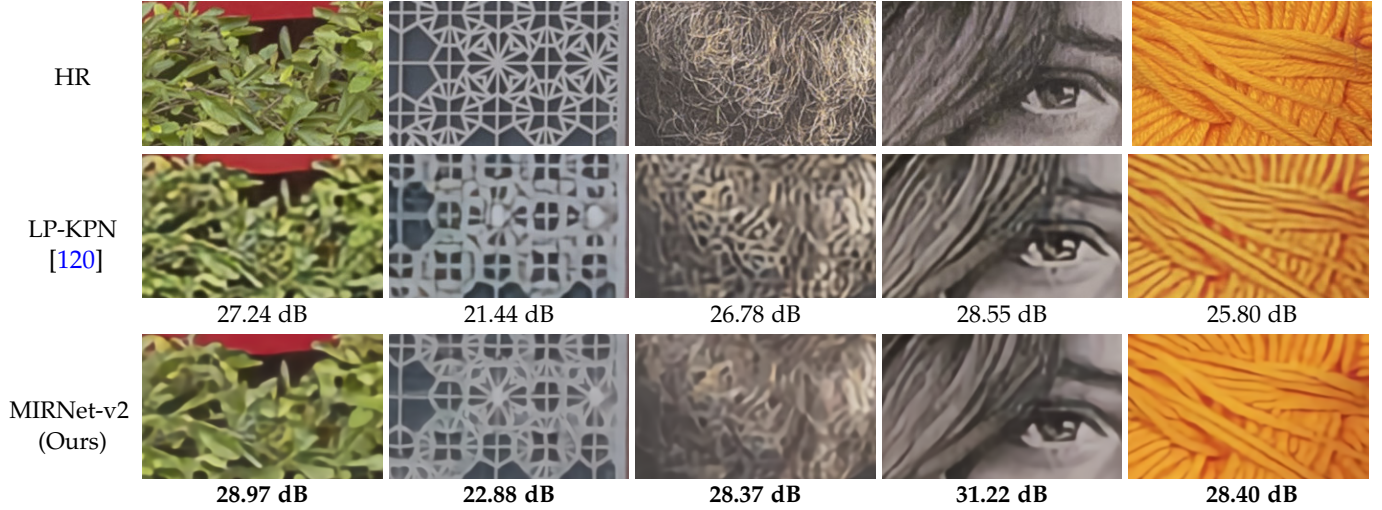


Fig. 7: Additional visual examples for $\times 4$ super-resolution, comparing our MIRNet-v2 against the state-of-the-art approach [120]. Note that all example crops are taken from different images.

TABLE 5: Low-light image enhancement evaluation on the LoL dataset [87]. The proposed method significantly advances the state-of-the-art.

Method	BIMEF [129]	CRM [130]	Dong [131]	LIME [132]	MF [133]	RRM [134]	SRIE [133]	Retinex-Net [87]	MSR [83]	NPE [135]	GLAD [136]	KinD [4]	KinD++ [137]	MIRNet-v2 (Ours)
PSNR	13.86	17.20	16.72	16.76	18.79	13.88	11.86	16.77	13.17	16.97	19.72	20.87	<u>21.30</u>	24.74
SSIM	0.577	0.644	0.582	0.564	0.642	0.658	0.498	0.559	0.479	0.589	0.703	0.810	<u>0.822</u>	0.851

TABLE 6: Image enhancement comparisons on the MIT-Adobe FiveK dataset [121].

Method	HDRNet [138]	W-Box [122]	DR [123]	DPE [94]	DeepUPE [124]	MIRNet-v2 (Ours)
PSNR	21.96	18.57	20.97	22.15	<u>23.04</u>	23.97
SSIM	0.866	0.701	0.841	0.850	<u>0.893</u>	0.931

that our MIRNet-v2 achieves significant improvements over previous approaches. Notably, when compared to the recent best methods, MIRNet-v2 obtains 3.44 dB performance gain over KinD++ [137] on the LoL dataset and 0.93 dB improvement over DeepUPE¹ [124] on the Adobe-Fivek dataset.

We show visual results in Fig. 8 and Fig. 9. Compared to other techniques, our method generates enhanced images that are natural and vivid in appearance and have better global and local contrast.

4.7 Ablation Studies

We study the impact of each of our architectural components and design choices on the final performance. All the ablation experiments are performed for the super-resolution task with $\times 3$ scale factor. The ablation models are trained on image patches of size 128×128 for 10^5 iterations. Table 7 shows that removing skip connections causes the largest performance drop. Without skip connections, the network finds it difficult to converge and yields high training errors, and consequently low PSNR. Furthermore, the information exchange among parallel convolution streams via SKFF is helpful and leads to improved performance. Similarly, RCB contributes positively towards the final image quality.

1. Note that the quantitative results reported in [124] are incorrect. The correct scores are later released by the original authors [link].

Table 8 shows that the proposed RCB provides favorable performance gain over the baseline Resblock from EDSR [74]. Moreover, removing the transform part from RCB causes drop in accuracy. Table 8 also shows that replacing the group convolutions with regular convolutions in RCB increases the PSNR score, but at the cost of significant increase in parameters and FLOPs. Therefore, we opt for RCB with group convolutions ($g=2$) as a balanced choice.

Next, we analyze the feature aggregation strategy in Table 9. It shows that the proposed SKFF generates favorable results compared to summation and concatenation. Note that our proposed SKFF module uses $\sim 5 \times$ fewer parameters than concatenation. Table 10 shows that the progressive learning strategy on mixed-size image patches yields PSNR similar to the model trained on large image patches ($ps=224$), but takes less time for training. Finally, in Table 11 we study how the number of convolutional streams and columns (RCB blocks) of MRB affect the image restoration quality. We note that increasing the number of streams provides significant improvements, thereby justifying the importance of multi-scale features processing. Moreover, increasing the number of columns yields better scores, thus indicating the significance of information exchange among parallel streams for feature consolidation.

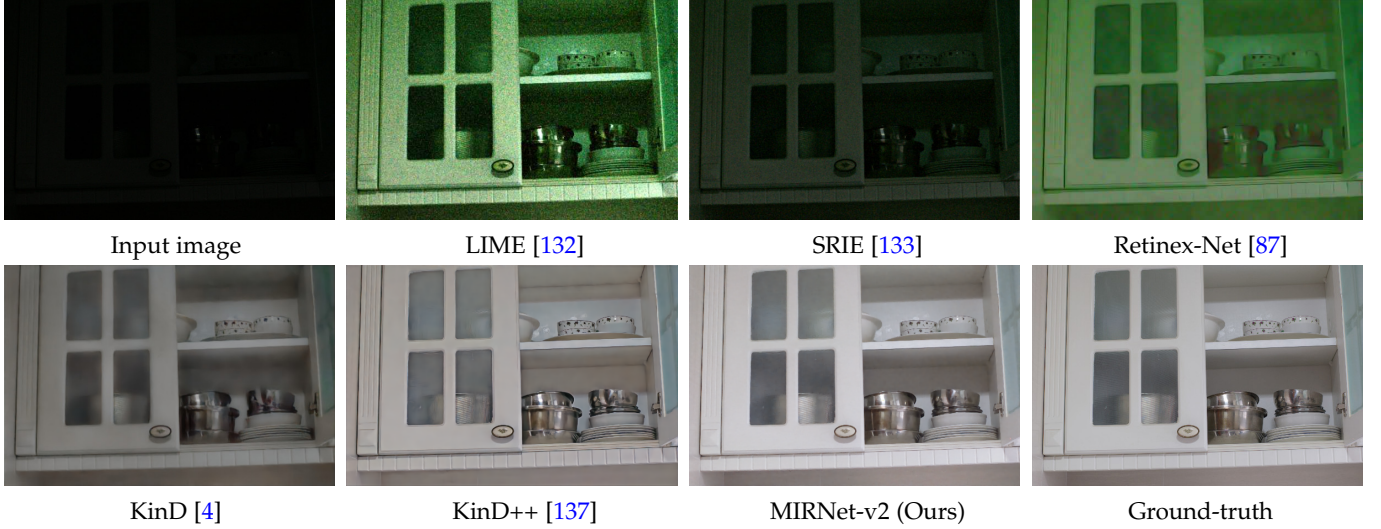


Fig. 8: Visual comparison of low-light enhancement approaches on the LoL dataset [87]. The image produced by our method is visually closer to the ground-truth in terms of brightness and global contrast.

TABLE 7: Impact of individual components of MRB.

Skip connections		✓	✓	✓	✓
RCB	✓			✓	✓
SKFF intermediate	✓	✓			✓
SKFF final	✓	✓	✓	✓	✓
PSNR (in dB)	28.21	30.79	30.85	30.68	30.97

TABLE 8: Effect of individual components of RCB. Resblock from EDSR [74] is taken as baseline. FLOPs are calculated on an image of size 256×256 . ‘g’ represents the number of groups in the group convolutions.

	PSNR	Params (M)	FLOPs (B)
Baseline [74], g=2	30.84	5.0	139.5
+ RCB, g=2	30.97	5.9	139.8
RCB w/o transform, g=2	30.92	5.0	139.7
RCB, g=1	31.05	9.7	253.2

TABLE 9: Feature aggregation. Our SKFF uses $\sim 5 \times$ fewer parameters than ‘Concat’, but generates better results.

	Sum	Concat	SKFF
PSNR (in dB)	30.76	30.83	30.97
Parameters	0	8,192	1,536

TABLE 10: Effect of progressive learning. For progressive training, we gradually increase image patch size from 128×128 to 224×224 .

Patch size	128	144	192	224	Progressive
PSNR (in dB)	30.97	30.99	31.02	31.08	31.06
Train time (h)	14	17	25	33	22

5 CONCLUDING REMARKS

Conventional image restoration and enhancement pipelines either stick to the full resolution features along the net-

TABLE 11: Ablation study on different layouts of MRB. Rows denote the number of parallel resolution streams, and Cols represent the number of columns containing RCBs.

PSNR	Cols = 1	Cols = 2	Cols = 3
Rows = 1	30.01	30.29	30.47
Rows = 2	30.65	30.79	30.85
Rows = 3	30.73	30.97	31.03

work hierarchy or use an encoder-decoder architecture. The first approach helps retain precise spatial details, while the latter one provides better contextualized representations. However, these methods can satisfy only one of the above two requirements, although real-world image restoration tasks demand a combination of both conditioned on the given input sample. In this work, we propose a novel architecture whose main branch is dedicated to full-resolution processing and the complementary set of parallel branches provides better contextualized features. We propose novel mechanisms to learn relationships between features within each branch as well as across multi-scale branches. Our feature fusion strategy ensures that the receptive field can be dynamically adapted without sacrificing the original feature details. Consistent achievement of state-of-the-art results on six datasets for four image restoration and enhancement tasks corroborates the effectiveness of our approach.

ACKNOWLEDGEMENTS

Ming-Hsuan Yang is supported by NSF CAREER grant 1149783. Ling Shao is partially supported by the National Natural Science Foundation of China (grant no. 61929104). Munawar Hayat is supported by the ARC DECRA Fellowship DE200101100.

REFERENCES

- [1] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: convolutional networks for biomedical image segmentation. In *MICCAI*, 2015. 1, 4

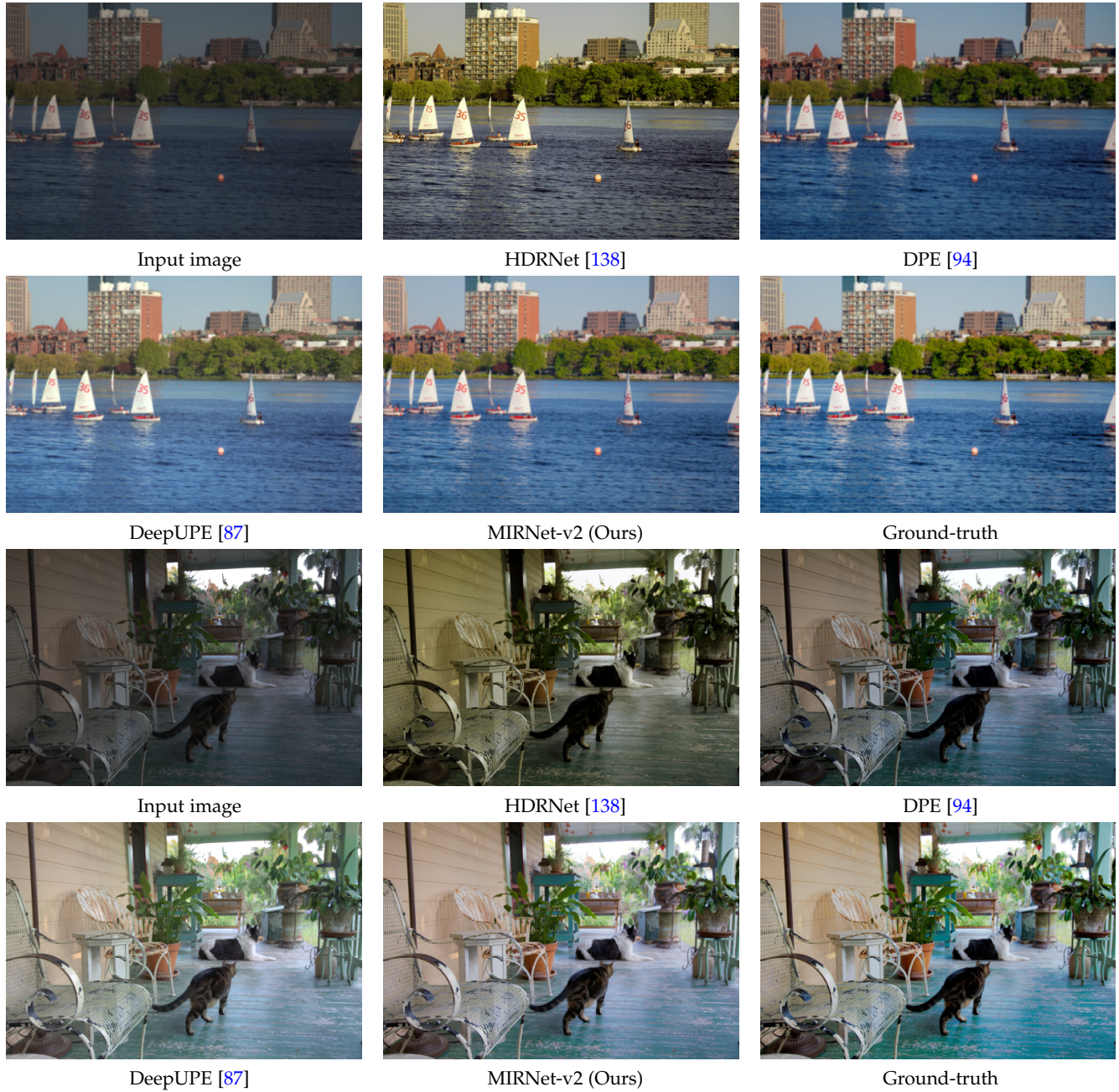


Fig. 9: Visual results of image enhancement on the MIT-Adobe FiveK [121] dataset. Compared to the state-of-the-art, our MIRNet-v2 makes better color and contrast adjustments and produces images that appear vivid, natural and pleasant.

- [2] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *ICCV*, 2019. 1
- [3] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *CVPR*, 2018. 1
- [4] Yonghua Zhang, Jiawan Zhang, and Xiaojie Guo. Kindling the darkness: A practical low-light image enhancer. In *MM*, 2019. 1, 3, 9, 10
- [5] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *TPAMI*, 2015. 1, 3
- [6] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *TIP*, 2017. 1, 2, 6
- [7] Yulun Zhang, Yapeng Tian, Yu Kong, Bineng Zhong, and Yun Fu. Residual dense network for image restoration. *TPAMI*, 2020. 1, 3
- [8] Andrey Ignatov, Nikolay Kobyshev, Radu Timofte, Kenneth Vanhoey, and Luc Van Gool. DSLR-quality photos on mobile devices with deep convolutional networks. In *ICCV*, 2017. 1
- [9] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In *ECCV*, 2020. 2
- [10] Abdelrahman Abdelhamed, Stephen Lin, and Michael S Brown. A high-quality denoising dataset for smartphone cameras. In *CVPR*, 2018. 2, 5, 6, 8
- [11] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Jiaya Jia. Scale-recurrent network for deep image deblurring. In *CVPR*, 2018. 1
- [12] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, 2017. 1
- [13] Shuhang Gu, Yawei Li, Luc Van Gool, and Radu Timofte. Self-guided network for fast image denoising. In *ICCV*, 2019. 1
- [14] Abdullah Abuolaim and Michael S Brown. Defocus deblurring using dual-pixel data. In *ECCV*, 2020. 2, 5, 6, 7

- [15] Laurent D'Andrès, Jordi Salvador, Axel Kochale, and Sabine Süsstrunk. Non-parametric blur map regression for depth of field extension. *TIP*, 2016. 2
- [16] Abdullah Abuolaim, Mauricio Delbracio, Damien Kelly, Michael S. Brown, and Peyman Milanfar. Learning to reduce defocus blur by realistically modeling dual-pixel data. In *ICCV*, 2021. 2, 6, 7
- [17] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, 2022. 2
- [18] Ali Karaali and Claudio Rosito Jung. Edge-based defocus blur estimation with adaptive scale selection. *TIP*, 2017. 2, 6, 7
- [19] Jianping Shi, Li Xu, and Jiaya Jia. Just noticeable defocus blur detection and estimation. In *CVPR*, 2015. 2, 6, 7
- [20] Junyong Lee, Sungkil Lee, Sunghyun Cho, and Seungyong Lee. Deep defocus map estimation using domain adaptation. In *CVPR*, 2019. 2, 6, 7
- [21] Leonid P Yaroslavsky. Local adaptive image restoration and enhancement with the use of DFT and DCT in a running window. In *Wavelet Applications in Signal and Image Processing IV*, 1996. 2
- [22] Eero P Simoncelli and Edward H Adelson. Noise removal via bayesian wavelet coring. In *ICIP*, 1996. 2
- [23] Carlo Tomasi and Roberto Manduchi. Bilateral filtering for gray and color images. In *ICCV*, 1998. 2
- [24] Pietro Perona and Jitendra Malik. Scale-space and edge detection using anisotropic diffusion. *TPAMI*, 1990. 2
- [25] Leonid I Rudin, Stanley Osher, and Emad Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 1992. 2
- [26] Alexei A Efros and Thomas K Leung. Texture synthesis by non-parametric sampling. In *ICCV*, 1999. 2
- [27] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *CVPR*, 2005. 2
- [28] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising by sparse 3-D transform-domain collaborative filtering. *TIP*, 2007. 2, 6, 8
- [29] Weisheng Dong, Guangming Shi, and Xin Li. Nonlocal image restoration with bilateral variance estimation: a low-rank approach. *TIP*, 2012. 2
- [30] Shuhang Gu, Lei Zhang, Wangmeng Zuo, and Xiangchu Feng. Weighted nuclear norm minimization with application to image denoising. In *CVPR*, 2014. 2
- [31] Julien Mairal, Francis Bach, Jean Ponce, Guillermo Sapiro, and Andrew Zisserman. Non-local sparse models for image restoration. In *ICCV*, 2009. 2
- [32] Rachid Hedjam, Reza Farrahi Moghaddam, and Mohamed Cheriet. Markovian clustering for the non-local means image denoising. In *ICIP*, 2009. 2
- [33] Saeed Anwar and Nick Barnes. Real image denoising with feature attention. *ICCV*, 2019. 2, 4, 6, 8
- [34] Tim Brooks, Ben Mildenhall, Tianfan Xue, Jiawen Chen, Dillon Sharlet, and Jonathan T Barron. Unprocessing images for learned raw denoising. In *CVPR*, 2019. 2
- [35] Shi Guo, Zifei Yan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Toward convolutional blind denoising of real photographs. In *CVPR*, 2019. 2, 6, 8
- [36] Tobias Plötz and Stefan Roth. Neural nearest neighbors networks. In *NeurIPS*, 2018. 2
- [37] Kai Zhang, Wangmeng Zuo, and Lei Zhang. FFDNet: Toward a fast and flexible solution for CNN-based image denoising. *TIP*, 2018. 2
- [38] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. CycleISP: Real image restoration via improved data synthesis. In *CVPR*, 2020. 2, 6, 8
- [39] Meng Chang, Qi Li, Huajun Feng, and Zhihai Xu. Spatial-adaptive network for single image denoising. In *ECCV*, 2020. 2, 6, 8
- [40] Zongsheng Yue, Qian Zhao, Lei Zhang, and Deyu Meng. Dual adversarial network: Toward real-world noise removal and noise generation. In *ECCV*, 2020. 2, 6, 8
- [41] Zongsheng Yue, Hongwei Yong, Qian Zhao, Deyu Meng, and Lei Zhang. Variational denoising network: Toward blind noise modeling and removal. In *NeurIPS*, 2019. 2, 6, 8
- [42] Yoonsik Kim, Jae Woong Soh, Gu Yong Park, and Nam Ik Cho. Transfer learning from synthetic to real-noise denoising with adaptive instance normalization. In *CVPR*, 2020. 2, 6, 8
- [43] Faming Fang, Juncheng Li, Yiting Yuan, Tiejong Zeng, and Guixu Zhang. Multilevel edge features guided network for image denoising. *IEEE TNNLS*, 2020. 2
- [44] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *CVPR*, 2021. 2
- [45] Robert Keys. Cubic convolution interpolation for digital image processing. *TASSP*, 1981. 3, 6
- [46] Michal Irani and Shmuel Peleg. Improving resolution by image registration. *CVGIP*, 1991. 3
- [47] Jan Allebach and Ping Wah Wong. Edge-directed interpolation. In *ICIP*, 1996. 3
- [48] Lei Zhang and Xiaolin Wu. An edge-guided image interpolation algorithm via directional filtering and data fusion. *TIP*, 2006. 3
- [49] Kwang In Kim and Younghee Kwon. Single-image super-resolution using sparse regression and natural image prior. *TPAMI*, 2010. 3
- [50] Zhiwei Xiong, Xiaoyan Sun, and Feng Wu. Robust web image/video super-resolution. *TIP*, 2010. 3
- [51] Hong Chang, Dit-Yan Yeung, and Yimin Xiong. Super-resolution through neighbor embedding. In *CVPR*, 2004. 3
- [52] Gilad Freedman and Raanan Fattal. Image and video upscaling from local self-examples. *TOG*, 2011. 3
- [53] Jianchao Yang, John Wright, Thomas S Huang, and Yi Ma. Image super-resolution via sparse representation. *TIP*, 2010. 3
- [54] Jianchao Yang, John Wright, Thomas Huang, and Yi Ma. Image super-resolution as sparse representation of raw image patches. In *CVPR*, 2008. 3
- [55] Zhihao Wang, Jian Chen, and Steven CH Hoi. Deep learning for image super-resolution: A survey. *TPAMI*, 2019. 3, 7
- [56] Saeed Anwar, Salman Khan, and Nick Barnes. A deep journey into super-resolution: A survey. *arXiv*, 2019. 3, 7
- [57] Jianrui Cai, Shuhang Gu, Radu Timofte, and Lei Zhang. Ntire 2019 challenge on real image super-resolution: Methods and results. In *CVPRW*, 2019. 3
- [58] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Learning a deep convolutional network for image super-resolution. In *ECCV*, 2014. 3
- [59] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Accurate image super-resolution using very deep convolutional networks. In *ICCV*, 2016. 3, 6, 7, 8
- [60] Ying Tai, Jian Yang, Xiaoming Liu, and Chunyan Xu. Memnet: A persistent memory network for image restoration. In *ICCV*, 2017. 3
- [61] Ying Tai, Jian Yang, and Xiaoming Liu. Image super-resolution via deep recursive residual network. In *CVPR*, 2017. 3
- [62] Zheng Hui, Xiumei Wang, and Xinbo Gao. Fast and accurate single image super-resolution via information distillation network. In *CVPR*, 2018. 3
- [63] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3
- [64] Jiwon Kim, Jung Kwon Lee, and Kyoung Mu Lee. Deeply-recursive convolutional network for image super-resolution. In *CVPR*, 2016. 3
- [65] Wei Han, Shiyu Chang, Ding Liu, Mo Yu, Michael Witbrock, and Thomas S Huang. Image super-resolution via dual-state recurrent networks. In *CVPR*, 2018. 3
- [66] Namhyuk Ahn, Byungkoo Kang, and Kyung-Ah Sohn. Fast, accurate, and lightweight super-resolution with cascading residual network. In *ECCV*, 2018. 3
- [67] Zhaowen Wang, Ding Liu, Jianchao Yang, Wei Han, and Thomas Huang. Deep networks for image super-resolution with sparse prior. In *ICCV*, 2015. 3
- [68] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate superresolution. In *CVPR*, 2017. 3
- [69] Tong Tong, Gen Li, Xiejie Liu, and Qinquan Gao. Image super-resolution using dense skip connections. In *ICCV*, 2017. 3
- [70] Xintao Wang, Ke Yu, Shixiang Wu, Jinjin Gu, Yihao Liu, Chao Dong, Yu Qiao, and Chen Change Loy. ESRGAN: enhanced super-resolution generative adversarial networks. In *ECCVW*, 2018. 3

- [71] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *ECCV*, 2018. 3, 4, 6, 7, 8
- [72] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *CVPR*, 2019. 3, 4
- [73] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. In *ICLR*, 2019. 3, 4
- [74] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *CVPRW*, 2017. 3, 9, 10
- [75] Ryan Dahl, Mohammad Norouzi, and Jonathon Shlens. Pixel recursive super resolution. In *ICCV*, 2017. 3
- [76] Juncheng Li, Faming Fang, Kangfu Mei, and Guixu Zhang. Multi-scale residual network for image super-resolution. In *ECCV*, 2018. 3
- [77] Seong-Jin Park, Hyeonseoek Son, Sunghyun Cho, Ki-Sang Hong, and Seungyong Lee. SFEAT: Single image super-resolution with feature discrimination. In *ECCV*, 2018. 3
- [78] Mehdi SM Sajjadi, Bernhard Scholkopf, and Michael Hirsch. Enhancenet: Single image super-resolution through automated texture synthesis. In *ICCV*, 2017. 3
- [79] Christian Ledig, Lucas Theis, Ferenc Huszár, Jose Caballero, Andrew Cunningham, Alejandro Acosta, Andrew Aitken, Alykhan Tejani, Johannes Totz, Zehan Wang, et al. Photo-realistic single image super-resolution using a generative adversarial network. In *CVPR*, 2017. 3, 6, 7, 8
- [80] Edwin H Land. The retinex theory of color vision. *Scientific american*, 1977. 3
- [81] Marcelo Bertalmío, Vicent Caselles, Edoardo Provenzi, and Alessandro Rizzi. Perceptual color correction through variational techniques. *TIP*, 2007. 3
- [82] R. Palma-Amestoy, E. Provenzi, M. Bertalmío, and V. Caselles. A perceptually inspired variational framework for color enhancement. *TPAMI*, 2009. 3
- [83] Daniel J Jobson, Zia-ur Rahman, and Glenn A Woodell. A multiscale retinex for bridging the gap between color images and the human observation of scenes. *TIP*, 1997. 3, 9
- [84] Alessandro Rizzi, Carlo Gatta, and Daniele Marini. From retinex to automatic color equalization: issues in developing a new algorithm for unsupervised color equalization. *Journal of Electronic Imaging*, 2004. 3
- [85] Andrey Ignatov and Radu Timofte. Ntire 2019 challenge on image enhancement: Methods and results. In *CVPRW*, 2019. 3
- [86] Liang Shen, Zihan Yue, Fan Feng, Quan Chen, Shihao Liu, and Jie Ma. Msr-net: Low-light image enhancement using deep convolutional network. *arXiv*, 2017. 3
- [87] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. *BMVC*, 2018. 3, 5, 8, 9, 10, 11
- [88] Huibin Chang, Michael K Ng, Wei Wang, and Tiejong Zeng. Retinex image enhancement via a learned dictionary. *Optical Engineering*, 2015. 3
- [89] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *ECCV*, 2018. 3
- [90] Kin Gwn Lore, Adedotun Akintayo, and Soumik Sarkar. LLNet: a deep autoencoder approach to natural low-light image enhancement. *Pattern Recognition*, 2017. 3
- [91] Wenqi Ren, Sifei Liu, Lin Ma, Qianqian Xu, Xiangyu Xu, Xiaochun Cao, Junping Du, and Ming-Hsuan Yang. Low-light image enhancement via a deep hybrid network. *TIP*, 2019. 3
- [92] Kangfu Mei, Juncheng Li, Jiajie Zhang, Haoyu Wu, Jie Li, and Rui Huang. Higher-resolution network for image demosaicing and enhancing. In *ICCVW*, 2019. 3
- [93] Jiaqian Li, Juncheng Li, Faming Fang, Fang Li, and Guixu Zhang. Luminance-aware pyramid network for low-light image enhancement. *IEEE Transactions on Multimedia*, 2020. 3
- [94] Yu-Sheng Chen, Yu-Ching Wang, Man-Hsin Kao, and Yung-Yu Chuang. Deep photo enhancer: Unpaired learning for image enhancement from photographs with gans. In *CVPR*, 2018. 3, 9, 11
- [95] Andrey Ignatov, Nikolay Kobyshev, Radu Timofte, Kenneth Vanhoey, and Luc Van Gool. Wespe: weakly supervised photo enhancer for digital cameras. In *CVPRW*, 2018. 3
- [96] Yubin Deng, Chen Change Loy, and Xiaoou Tang. Aesthetic-driven image enhancement by adversarial learning. In *ACM Multimedia*, 2018. 3
- [97] Pierre Charbonnier, Laure Blanc-Feraud, Gilles Aubert, and Michel Barlaud. Two deterministic half-quadratic regularization algorithms for computed imaging. In *ICIP*, 1994. 3
- [98] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hourglass networks for human pose estimation. In *ECCV*, 2016. 4
- [99] Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han. Learning deconvolution network for semantic segmentation. In *ICCV*, 2015. 4
- [100] Bin Xiao, Haiping Wu, and Yichen Wei. Simple baselines for human pose estimation and tracking. In *ECCV*, 2018. 4
- [101] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. SegNet: a deep convolutional encoder-decoder architecture for image segmentation. *TPAMI*, 2017. 4
- [102] Xi Peng, Rogerio S Feris, Xiaoyu Wang, and Dimitris N Metaxas. A recurrent encoder-decoder network for sequential face alignment. In *ECCV*, 2016. 4
- [103] David H Hubel and Torsten N Wiesel. Receptive fields, binocular interaction and functional architecture in the cat's visual cortex. *The Journal of physiology*, 1962. 4
- [104] Maximilian Riesenhuber and Tomaso Poggio. Hierarchical models of object recognition in cortex. *Nature neuroscience*, 1999. 4
- [105] Thomas Serre, Lior Wolf, Stanley Bileschi, Maximilian Riesenhuber, and Tomaso Poggio. Robust object recognition with cortex-like mechanisms. *TPAMI*, 2007. 4
- [106] Chou P Hung, Gabriel Kreiman, Tomaso Poggio, and James J DiCarlo. Fast readout of object identity from macaque inferior temporal cortex. *Science*, 2005. 4
- [107] Gao Huang, Danlu Chen, Tianhong Li, Felix Wu, Laurens van der Maaten, and Kilian Q Weinberger. Multi-scale dense networks for resource efficient image classification. In *ICLR*, 2018. 4
- [108] Ke Sun, Bin Xiao, Dong Liu, and Jingdong Wang. Deep high-resolution representation learning for human pose estimation. In *CVPR*, 2019. 4
- [109] Damien Fourure, Rémi Emonet, Élisabeth Fromont, Damien Muselet, Alain Trémeau, and Christian Wolf. Residual conv-deconv grid network for semantic segmentation. In *BMVC*, 2017. 4
- [110] C. Szegedy, Wei Liu, Yangqing Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. In *CVPR*, 2015. 4
- [111] Xiang Li, Wenhao Wang, Xiaolin Hu, and Jian Yang. Selective kernel networks. In *CVPR*, 2019. 4
- [112] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In *CVPR*, 2018. 4
- [113] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018. 4
- [114] Yue Cao, Jiarui Xu, Stephen Lin, Fangyuan Wei, and Han Hu. Global context networks. *TPAMI*, 2020. 4
- [115] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM Computing Surveys*, 2021. 4
- [116] Elad Hoffer, Berry Weinstein, Itay Hubara, Tal Ben-Nun, Torsten Hoefler, and Daniel Soudry. Mix & match: training convnets with mixed image sizes for improved accuracy, speed and scale resiliency. *arXiv:1908.08986*, 2019. 5
- [117] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *ICML*, 2019. 5
- [118] Tobias Plotz and Stefan Roth. Benchmarking denoising algorithms with real photographs. In *CVPR*, 2017. 5, 6, 8
- [119] <https://noise.visinf.tu-darmstadt.de/benchmark/>, 2017. [Online; accessed 29-Feb-2020]. 5
- [120] Jianrui Cai, Hui Zeng, Hongwei Yong, Zisheng Cao, and Lei Zhang. Toward real-world single image super-resolution: A new benchmark and a new model. In *ICCV*, 2019. 5, 6, 7, 8, 9
- [121] Vladimir Bychkovsky, Sylvain Paris, Eric Chan, and Frédo Durand. Learning photographic global tonal adjustment with a database of input/output image pairs. In *CVPR*, 2011. 5, 8, 9, 11
- [122] Yuanming Hu, Hao He, Chenxi Xu, Baoyuan Wang, and Stephen Lin. Exposure: A white-box photo post-processing framework. *TOG*, 2018. 5, 9
- [123] Jongchan Park, Joon-Young Lee, Donggeun Yoo, and In So Kweon. Distort-and-recover: Color enhancement using deep reinforcement learning. In *CVPR*, 2018. 5, 9

- [124] Ruixing Wang, Qing Zhang, Chi-Wing Fu, Xiaoyong Shen, Wei-Shi Zheng, and Jiaya Jia. Underexposed photo enhancement using deep illumination estimation. In *CVPR*, 2019. 5, 9
- [125] Ilya Loshchilov and Frank Hutter. SGDR: Stochastic gradient descent with warm restarts. In *ICLR*, 2017. 6
- [126] Harold C Burger, Christian J Schuler, and Stefan Harmeling. Image denoising: Can plain neural networks compete with BM3D? In *CVPR*, 2012. 6
- [127] Chong Mou, Jian Zhang, and Zhuoyuan Wu. Dynamic attentive graph learning for image restoration. In *ICCV*, 2021. 6
- [128] Chao Ren, Xiaohai He, Chuncheng Wang, and Zhibo Zhao. Adaptive consistency prior based deep network for image denoising. In *CVPR*, 2021. 6
- [129] Zhenqiang Ying, Ge Li, and Wen Gao. A bio-inspired multi-exposure fusion framework for low-light image enhancement. *arXiv preprint arXiv:1711.00591*, 2017. 9
- [130] Zhenqiang Ying, Ge Li, Yurui Ren, Ronggang Wang, and Wenmin Wang. A new image contrast enhancement algorithm using exposure fusion framework. In *CAIP*, 2017. 9
- [131] Xuan Dong, Guan Wang, Yi Pang, Weixin Li, Jiangtao Wen, Wei Meng, and Yao Lu. Fast efficient algorithm for enhancement of low lighting video. In *ICME*, 2011. 9
- [132] Xiaojie Guo, Yu Li, and Haibin Ling. Lime: Low-light image enhancement via illumination map estimation. *TIP*, 2016. 9, 10
- [133] Xueyang Fu, Delu Zeng, Yue Huang, Xiao-Ping Zhang, and Xinghao Ding. A weighted variational model for simultaneous reflectance and illumination estimation. In *CVPR*, 2016. 9, 10
- [134] Yong Liu, Ruiping Wang, Shiguang Shan, and Xilin Chen. Structure inference net: Object detection using scene-level context and instance-level relationships. In *CVPR*, 2018. 9
- [135] Shuhang Wang, Jin Zheng, Hai-Miao Hu, and Bo Li. Naturalness preserved enhancement algorithm for non-uniform illumination images. *TIP*, 2013. 9
- [136] Wenjing Wang, Chen Wei, Wenhan Yang, and Jiaying Liu. Glad-net: Low-light enhancement network with global awareness. In *FG*, 2018. 9
- [137] Yonghua Zhang, Xiaojie Guo, Jiayi Ma, Wei Liu, and Jiawan Zhang. Beyond brightening low-light images. *IJCV*, 2021. 9, 10
- [138] Michaël Gharbi, Jiawen Chen, Jonathan T Barron, Samuel W Hasinoff, and Frédo Durand. Deep bilateral learning for real-time image enhancement. *TOG*, 2017. 9, 11



Syed Waqas Zamir received the Ph.D. degree from University Pompeu Fabra, Spain, in 2017. He is a Research Scientist at Inception Institute of Artificial Intelligence in UAE. His research interests include low-level computer vision, computational imaging, image and video processing, color vision and image restoration and enhancement.



Aditya Arora is a Research Engineer at Inception Institute of Artificial Intelligence in UAE. His research interests include image and video processing, computational photography and low-level vision.



received his Ph.D. degree from the University of Western Australia in 2016. His thesis received an honorable mention on the Dean's List Award. His research interests include computer vision and machine learning.



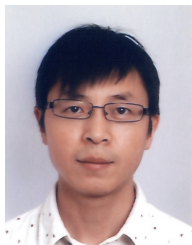
his field, including TPAMI, IJCV, CVPR, ECCV and ICCV. His research interests are in computer vision and machine/deep learning.



Fahad Khan is a faculty member at MBZUAI, United Arab Emirates and Linköping University, Sweden. From 2018 to 2020 he worked as a Lead Scientist at the Inception Institute of Artificial Intelligence (IIAI), Abu Dhabi, United Arab Emirates. He received the M.Sc. degree in Intelligent Systems Design from Chalmers University of Technology, Sweden and a Ph.D. degree in Computer Vision from Autonomous University of Barcelona, Spain. He has achieved top ranks on various international challenges (Visual Object Tracking VOT: 1st 2014 and 2018, 2nd 2015, 1st 2016; VOT-TIR: 1st 2015 and 2016; OpenCV Tracking: 1st 2015; 1st PASCAL VOC 2010). His research interests include a wide range of topics within computer vision and machine learning, such as object recognition, object detection, action recognition and visual tracking. He has published articles in high-impact computer vision journals and conferences in these areas. He serves as a regular program committee member for leading computer vision conferences such as CVPR, ICCV, and ECCV.



Ming-Hsuan Yang is affiliated with Google, UC Merced, and Yonsei University. Yang serves as a program co-chair of IEEE International Conference on Computer Vision (ICCV) in 2019, program co-chair of Asian Conference on Computer Vision (ACCV) in 2014, and general co-chair of ACCV 2016. Yang served as an associate editor of the IEEE Transactions on Pattern Analysis and Machine Intelligence, and is an associate editor of the International Journal of Computer Vision, Image and Vision Computing and Journal of Artificial Intelligence Research. He received the NSF CAREER award and Google Faculty Award. He is a Fellow of the IEEE.



Ling shao is the Chief Scientist of Terminus Group and the President of Terminus International. He was the founding CEO and Chief Scientist of the Inception Institute of Artificial Intelligence, Abu Dhabi, UAE. His research interests include computer vision, deep learning, medical imaging and vision and language. He is a fellow of the IEEE, the IAPR, the BCS and the IET.