# COLA-Net: Collaborative Attention Network for Image Restoration

Chong Mou, Jian Zhang, Xiaopeng Fan, Hangfan Liu, Ronggang Wang

*Abstract*—Local and non-local attention-based methods have been well studied in various image restoration tasks while leading to promising performance. However, most of the existing methods solely focus on one type of attention mechanism (local or non-local). Furthermore, by exploiting the self-similarity of natural images, existing pixel-wise non-local attention operations tend to give rise to deviations in the process of characterizing long-range dependence due to image degeneration. To overcome these problems, in this paper we propose a novel collaborative attention network (COLA-Net) for image restoration, as the first attempt to combine local and non-local attention mechanisms to restore image content in the areas with complex textures and with highly repetitive details respectively. In addition, an effective and robust patch-wise non-local attention model is developed to capture long-range feature correspondences through 3D patches. Extensive experiments on synthetic image denoising, real image denoising and compression artifact reduction tasks demonstrate that our proposed COLA-Net is able to achieve state-of-the-art performance in both peak signal-to-noise ratio and visual perception, while maintaining an attractive computational complexity. The source code is available on https://github.com/MC-E/COLA-Net.

*Index Terms*—Image restoration, deep neural network, image denoising, non-local attention, feature fusion.

## I. INTRODUCTION

IMAGE restoration aims to recover the underlying high-quality image $\mathbf{x}$ from its degraded measurement $\mathbf{y} = \mathbf{Ax} + \mathbf{n}$, where $\mathbf{A}$ is a linear degradation matrix, and $\mathbf{n}$ represents additive noise. It is typically an ill-posed problem due to the irreversible degradation process [1], [2], [3]. Since deep learning methods have been successfully applied in various computer vision tasks, many deep-learning-based methods [4], [5], [6], [7] have been proposed to solve this ill-posed problem. Most of them focus on local processing by combining convolutional layers and element-wise operations. To aggregate more useful information from a large receptive field, some strategies such as hourglass-shaped architecture [8], [9], [10], [11], dilated convolutions [12], [13], [6] or stacking more convolutional layers [14], [7], [15], [16] are applied in image restoration tasks. Furthermore, inspired by some traditional excellent image restoration algorithms such as BM3D [17], group sparse representation [18], and non-local

means [19], which take advantage of self-similarity within the whole image to recover a local image content. Some recent works tried to implement this idea via deep networks. One way is to establish the long-range dependence based on pixels, *e.g.*, NLRN [20] and RNAN [14], which applied non-local neural networks [21] in image restoration tasks. However, due to the fact that pixels are usually noisy in image restoration tasks, establishing relationships between pixels is prone to be biased and unreliable, especially when images are heavily corrupted. In addition to non-local attention mechanism, local attention mechanism is also an important strategy in computer vision community and has been studied in many image restoration tasks [22], [12], [23]. However, the main drawback of these methods is that the receptive field during image restoration is relatively small.

Through the above analysis, both local and non-local attentions have their unique drawbacks, but they can complement each other in many aspects. For instance, when an image contains sufficient repetitive details, non-local operations will be more useful. But when an image has a lot of complex textures, directly applying non-local operations will cause over-smooth artifacts, while local operations may be an appropriate choice. How to make a trade-off between local and non-local operations in image restoration has not been explored.

To address the above issues, we present the first attempt to exploit both local attention and non-local attention to restore image content in areas with complex textures and highly repetitive details, respectively. It is important to note that this combination is learnable and self-adaptive. Part of our previous work has been reported in [24]. To be concrete, for local attention operation, we apply local channel-wise attention on different scales to enlarge the size of receptive field of local operation, while for non-local attention operation, we develop a novel and robust patch-wise non-local attention model for constructing long-range dependence between image patches to restore every patch by aggregating useful information (self-similarity) from the whole image. **The main contributions of this paper are summarized as follows:**

- We propose a novel **COL**laborative **A**ttention **Net**work, dubbed **COLA-Net**, which incorporates both local and non-local attention mechanisms into deep networks for image restoration tasks. To the best of our knowledge, COLA-Net is the first attempt to combine local and non-local operations to restore complex textures and repetitive details distinguishingly.
- We propose an effective patch-wise non-local attention model to establish a more reliable long-range dependence during image restoration.

C. Mou, J. Zhang, and R. Wang are with the School of Electronic and Computer Engineering, Peking University Shenzhen Graduate School, Shenzhen, China. (e-mail: eechongm@stu.pku.edu.cn; zhangjian.sz@pku.edu.cn; rgwang@pkusz.edu.cn).

X. Fan is with the School of Computer Science and Technology, Harbin Institute of Technology, China (e-mail: fxp@hit.edu.cn).

H. Liu is with Center for Biomedical Image Computing & Analytics, University of Pennsylvania, USA (e-mail: hfliu@upenn.edu).

- We carry out extensive experiments on three typical image restoration tasks, i.e., synthetic image denoising, real image denoising, and compression artifact reduction, showing that our proposed COLA-Net achieves state-of-the-art results while maintaining an attractive computational complexity.

## II. RELATED WORK

In what follows, we give a brief review of both local and non-local attention mechanisms for image restoration and focus on the specific methods most relevant to our own.

**Non-local Attention.** Self-similarity is an important prior to image restoration especially for image denoising, and this prior has been widely used in several traditional non-local image restoration methods [19], [17], [18], [25], [26], [27]. In general, non-local attention models take a degraded input $\mathbf{y} = \{\mathbf{y}_i | i \in \mathbb{I}\}$, where $\mathbb{I}$ denotes the set of indices of pixels/patches in the whole image, and $\mathbf{y}_i$ stands for the $i$-th pixel or the patch centered on the $i$-th pixel. Each corrupted element $\mathbf{y}_i$ can be restored by a set of similar items $\mathbf{y}_j$ from a search region $\mathbb{Q} \subset \mathbb{I}$. Thus, non-local operations can be formally defined as:

$$\hat{\mathbf{x}}_i = \frac{1}{z_i} \sum_{j \in \mathbb{Q}} \phi(\mathbf{y}_i, \mathbf{y}_j) G(\mathbf{y}_j), \forall i, \tag{1}$$

where $\hat{\mathbf{x}}_i$ and $z_i$ represent the estimated clean value and the normalizing constant calculated by $z_i = \sum_{j \in \mathbb{Q}} \phi(\mathbf{y}_i, \mathbf{y}_j)$. The function $\phi$ can compute the similarity between query items $\mathbf{y}_i$ and key items $\mathbf{y}_j$. $G$ is an embedding function to transform $\mathbf{y}_j$ to another representation.

The seminal work non-local means [19] searched similar patches within a local region and averaged central pixels weighted by similarity. This method applied weighted Euclidean distance with Gaussian kernel to compute the similarity between two patches, which can be formulated as:

$$\phi(\mathbf{y}_i, \mathbf{y}_j) = e^{-\frac{\|\mathbf{y}_i - \mathbf{y}_j\|_{2,\alpha}^2}{h^2}}, \tag{2}$$

where $\alpha$ and $h$ refer to the standard deviation of Gaussian kernel and the filter factor, respectively. Identity mapping: $G(\mathbf{y}_j) = \mathbf{y}_j$ is directly used as the embedding function in this model. Rather than simply averaging similar pixels, the popular method BM3D [17] generated a stack of 3D matching patches and utilized a 3D filter to restore corrupted patches.

Recently, deep neural networks (DNN) have been prevalent in the community of computer vision, and non-local neural network [21] has been proposed for high-level vision tasks such as object detection and classification. In [21], embedding function $G$ is a convolutional layer that can be viewed as a linear embedding function, and dot product, i.e., $\phi(\mathbf{y}_i, \mathbf{y}_j) = \mathbf{y}_i^T \mathbf{y}_j$ is applied to compute the similarity between query and keys. Based on the success of non-local neural networks, NLRN [20] and RNAN [14] were proposed for image restoration tasks. In NLRN [20], the distance matrix of non-local layers is shared to enable the feature correlation to be propagated along with adjacent recurrent states. In RNAN [14], a residual non-local attention learning was proposed to train very deep networks by preserving more low-level

features. However, pixel-wise non-local attention is unreliable due to the image degeneration. Another way is to establish the long-range dependence based on patch matching and restore local patch or central pixel through similar patches in other places, e.g., [28], [29], [30]. Nevertheless, the patching matching step is isolated from the training process. Therefore it can not be jointly trained with image restoration networks. In [31], a light-weight model was proposed to utilize non-local attention and sparsity principles among feature patches for image restoration. In addition, N3Net [32] and [33] proposed learnable patch matching. Since they can only match a small number of blocks within a limited area, these two methods are not efficient enough as pixel matching mechanisms.

**Local Attention.** Compared with non-local attention, local attention is originally designed for high-level vision tasks [34], [35], [36]. RCAN [37] proposed a very deep residual channel attention networks for highly accurate image super-resolution. RIDNet [22] was the first attempt to combine local attention in image denoising tasks. Later on, MIRNet [23] expanded local attention to a wider range of image restoration tasks. To enlarging the receptive field, ADNet [13] applied dilated convolution for denoising tasks. However, compared with non-local methods, the main drawback of local attention mechanism during image restoration is the relatively small receptive field.

**Image Restoration Architectures.** Stacking convolutional layers is the most well-known CNN-based strategy for image restoration. Dong *et al.* proposed ARCNN [38] for image compression artifact reduction with several stacked convolutional layers. Zhang *et al.* proposed DnCNN [4] for image restoration with the help of residual learning and batch normalization. Based on DnCNN, Zhang *et al.* proposed FFDNet [5] for blind image denoising, which improves the generalization of DNN-based image restoration methods. Zhang *et al.* further proposed IRCNN [6], which applied dilated convolution in image restoration tasks. Recently, great progress has been made in the image restoration community, and diverse novel models have been proposed. Tai *et al.* proposed MemNet [7] to apply dense connection in convolutional layers, and He *et al.* [39] applied this idea to a light-weight design. Tian et al. [40] proposed a coarse-to-fine architecture to perform image super-resolution. Guo *et al.* proposed CBDNet [8], making a sufficient improvement in dealing with real-noise image corruption problems, and its hourglass-shaped architecture is applied in many image restoration works [9], [10], [11] afterward.

## III. COLLABORATIVE ATTENTION NETWORK (COLA-NET)

In this section, we will elaborate our proposed COLlaborative Attention Network (COLA-Net), which is able to adjust the contribution of local operation and non-local operation adaptively in the process of image restoration.

### A. Framework

An important concept presented in our work is that local attention operation and non-local attention operation are not
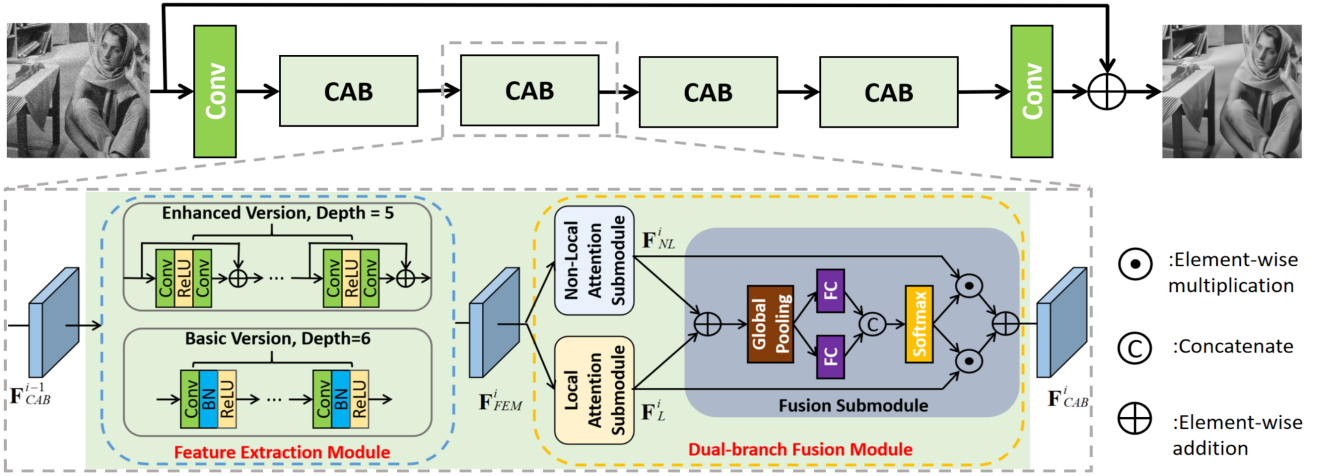
Fig. 1. Illustration of the network architecture of our proposed COLA-Net. COLA-Net is mainly cascaded by several collaborative attention blocks (CAB). The second row further shows the detailed architecture of each CAB.

all-powerful to restore any kind of image content. Using these two operations without distinction will produce passive effects. To rectify this weakness, we propose a novel collaborative attention network (COLA-Net) by combining local attention submodule and non-local attention submodule to restore complex textures and repetitive details, respectively. Note that this combination is not a simple series connection or parallel connection but is self-adaptive and learnable.

As illustrated in Fig. 1, our proposed COLA-Net is mainly composed of several collaborative attention blocks (CAB) (four by default) that are cascaded together. Let's denote $\mathbf{I}_{LQ}$ and $\mathbf{I}_{REC}$ as the low-quality (LQ) input and the reconstruction output of COLA-Net. We first use one convolutional layer to extract the shallow features $\mathbf{F}_S$ from $\mathbf{I}_{LQ}$:

$$\mathbf{F}_S = \mathcal{H}_{SF}(\mathbf{I}_{LQ}), \tag{3}$$

where $\mathcal{H}_{SF}(\cdot)$ denotes the shallow feature extraction. $\mathbf{F}_S$ is then used for feature restoration with collaborative attention blocks (CAB). So we can further have:

$$\mathbf{F}_{CAB}^i = \mathcal{H}_{CAB}^i(...\mathcal{H}_{CAB}^1(\mathbf{F}_S)), \tag{4}$$

where $\mathcal{H}_{CAB}^i(\cdot)$ and $\mathbf{F}_{CAB}^i$ denote the function of the $i$-th CAB and its corresponding restoration result, $i = 1, 2, 3, 4$. The architecture of CAB is shown in the second row of Fig. 1, which is composed of two parts, i.e., feature exaction module (FEM) and dual-branch fusion module (DFM). In this paper, we provide two versions of FEM. In consideration of the computational complexity, we develop a basic version of FEM that applies a DnCNN [4] architecture, containing six convolutional blocks. Each block consists of a $3 \times 3$ convolutional layer with 64 filters, a batch normalization layer, and a ReLU activation function. Considering the depth of neural networks is important for image restoration, an enhanced version of FEM in the residual domain is derived to further improve restoration performance with moderate growth of parameters. The COLA-Net equipped with the basic FEM is named COLA-B, and the COLA-Net equipped with the enhanced FEM is named COLA-E. All convolutional layers in FEM have 64 filters and $3 \times 3$ kernel size. In the experiment,

we set the depth of the basic FEM to six and the enhanced FEM to five. The output of FEM in the $i$-th CAB is denoted by $\mathbf{F}_{FEM}^i$. DFM is the core of COLA-Net, which performs local attention operation and non-local attention operation in two parallel branches and fuses them adaptively. More details of DFM will be given in the next subsection.

At the end of the network, we apply one convolutional layer to transform the output of the last CAB from feature domain to image domain and adopt residual learning to facilitate network training. Thus, we can get the output $\mathbf{I}_{REC}$ as:

$$\mathbf{I}_{REC} = \mathbf{I}_{LQ} + \mathcal{H}_{DF}(\mathbf{F}_{CAB}^4) = \mathcal{H}_{COLA}(\mathbf{I}_{LQ}), \tag{5}$$

where $\mathcal{H}_{DF}(\cdot)$ and $\mathcal{H}_{COLA}(\cdot)$ represent the functions of the last convolutional layer and the whole COLA-Net, respectively.

To show the effectiveness of COLA-Net, we choose the commonly used $L_2$ loss as the objective function. Given a training set $\{\mathbf{I}_{LQ}^b, \mathbf{I}_{HQ}^b\}_{b=1}^B$ that contains $B$ corrupted low-quality (LQ) inputs and their high-quality (HQ) labels. The goal of training can be defined as:

$$L(\mathbf{\Theta}) = \frac{1}{B} \sum_{b=1}^B \left\| \mathbf{I}_{HQ}^b - \mathcal{H}_{COLA}(\mathbf{I}_{LQ}^b) \right\|_2^2, \tag{6}$$

where $\mathbf{\Theta}$ refers to the learnable parameters of COLA-Net.

### B. Dual-branch Fusion Module

Our proposed dual-branch fusion module (DFM), labeled with a yellow box in Fig. 1, comprises three parts: local attention submodule, non-local attention submodule, and fusion submodule. The role of DFM is to restore complex textures through local operation, recover repetitive details based on self-similarity, and fuse dual-branch restoration results in an adaptive way.

As illustrated in Fig. 1, given $\mathbf{F}_{FEM}^i$, the local attention submodule and non-local attention submodule generate their outputs $\mathbf{F}_L^i$ and $\mathbf{F}_{NL}^i$ in a paralleled manner:

$$\begin{cases} \mathbf{F}_L^i = \mathcal{H}_L(\mathbf{F}_{FEM}^i) \\ \mathbf{F}_{NL}^i = \mathcal{H}_{NL}(\mathbf{F}_{FEM}^i), \end{cases} \tag{7}$$
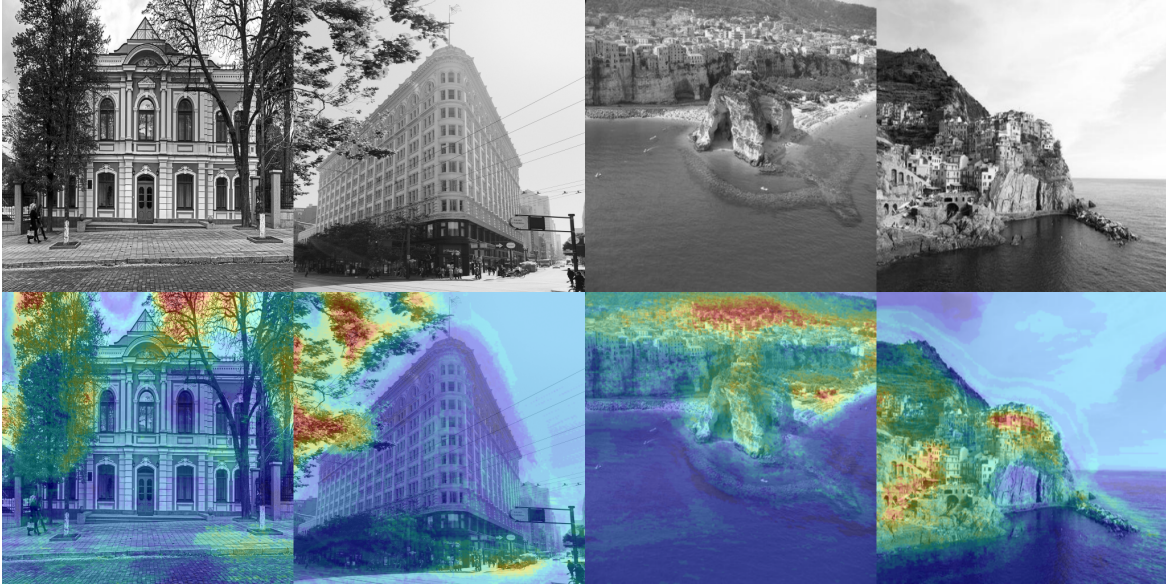
Fig. 2. Visualization of attention dependence during image restoration. The first row represents the original images (HQ) and the second row represents attention dependence during image restoration. For illustration purposes, the heat-maps directly overlap with original images in the second row. The deeper the color the more local operations are demanded, on the contrary, more non-local operations are demanded.

where $\mathcal{H}_L(\cdot)$ and $\mathcal{H}_{NL}(\cdot)$ refer to the functions of the local and non-local attention submodules. The architecture of these two submodules will be elaborated in the following subsections.

After getting $\mathbf{F}_{NL}^i$ and $\mathbf{F}_L^i$, and motivated by [36], the final output $\mathbf{F}_{CAB}^i$ is obtained by the fusion submodule in the following three steps:

$$\begin{cases} \mathbf{v} = GlobalPooling(\mathbf{F}_{NL}^i + \mathbf{F}_L^i) \\ \mathbf{w}_{NL}, \mathbf{w}_L = softmax([\mathcal{H}_{FC}^1(\mathbf{v}), \mathcal{H}_{FC}^2(\mathbf{v})]) \\ \mathbf{F}_{CAB}^i = \mathbf{F}_{NL}^i \cdot \mathbf{w}_{NL} + \mathbf{F}_L^i \cdot \mathbf{w}_L, \end{cases} \quad (8)$$

where $\mathcal{H}_{FC}^1(\cdot)$ and $\mathcal{H}_{FC}^2(\cdot)$ refer to two independent fully-connected layers.

To be concrete, as shown in Eq. 8, we first merge the restoration results from dual branches via an element-wise addition. Then, we apply a global pooling to produce a global feature vector $\mathbf{v} \in \mathbb{R}^C$, which is used as the guidance of adaptive and accurate selection between local operation and non-local operation. Next, we apply two fully connected layers ($\mathcal{H}_{FC}^1$ and $\mathcal{H}_{FC}^2$) to generate two weight vectors ($\mathbf{w}_{NL}$ and $\mathbf{w}_L$) to perform channel-wise selection between two restoration results. Obviously, $\mathbf{w}_{NL}$ and $\mathbf{w}_L$ are content-aware, thus self-adaptive. In this way, local and non-local attention branches can restore image content selectively based on their characteristics.

In order to further verify the necessity of local and non-local attention results, we make a visualization of the attention dependence during image restoration by heat-maps. The heat value $h$ in Fig. 2 is computed as follows:

$$h = \frac{1}{M} \sum_{m=1}^{M} \mathbf{a}_m, \quad (9)$$

where $\mathbf{a}_m = 1$ for $\mathbf{w}_L^m >= \mathbf{w}_{NL}^m$ and $0$ otherwise, and $M$ is the length of dependence weights ($\mathbf{w}_{NL}$ or $\mathbf{w}_L$). Thus, $h$
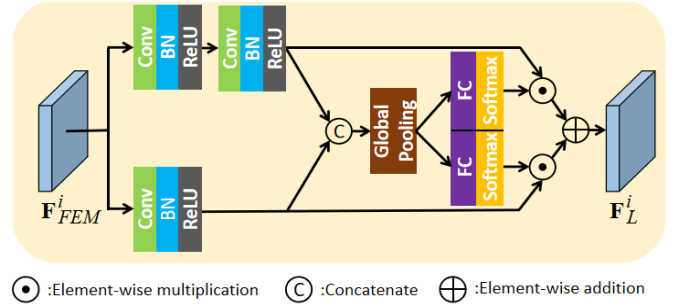


Fig. 3. Illustration of our proposed local attention submodule.

can represent the preference or focus between local and non-local attention operations in the process of image restoration in specific places. In Fig. 2, the deeper the color, the more local attention operations are demanded. On the contrary, more non-local attention operations are needed. For illustration purposes, we overlap the heat-maps on original images in the second row of Fig. 2. We can find that the heated parts are mainly concentrated in areas with complex textures or without sufficiently repetitive details indicating that these parts demand more local attention operations to restore. Conversely, light-colored areas contain a large amount of repetitive image content. The results of this experiment are entirely consistent with our motivation.

## C. Local Attention Submodule

Inspired by [36], [35], we employ the channel-wise attention mechanism to design the local attention submodule to perform channel selection with multiple sizes of receptive field. The detailed architecture of our local attention submodule is shown in Fig. 3. In this submodule, we adopt two branches to carry different numbers of convolutional layers to generate feature
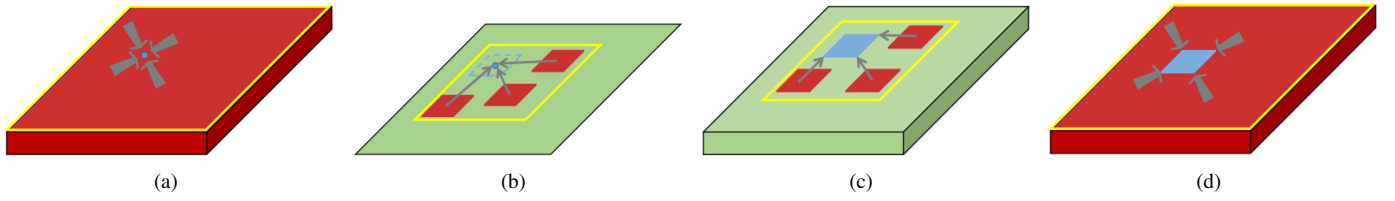
Fig. 4. Visual comparison of various non-local operations. The blue parts and the red parts in images represent restored image content and corresponding long-range dependent items. The yellow boundary refers to the search region. (a) Common non-local neural networks [21] applies pixel matching and pixel updating within the whole image. (b) Classic patch-wise non-local operation takes non-local means [19] for example, applying patch matching and pixel updating within a limited region. (c) Learnable patch-wise non-local operation takes N3Net [32] for example, applying learnable patch matching and patch updating with a limited region and limited matching items. (d) Our proposed patch-wise non-local model applies learnable patch matching and patch updating within the whole image.
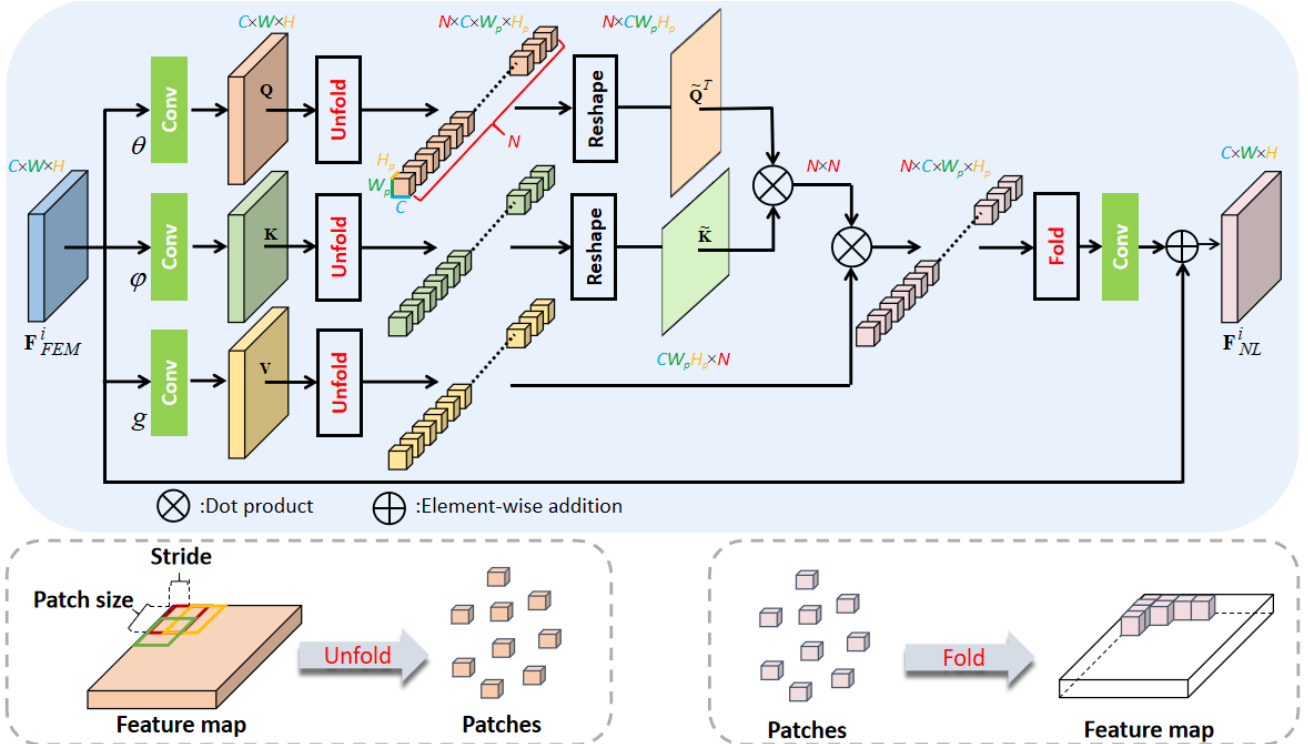


Fig. 5. Illustration of our proposed non-local attention submodule. In particular, a novel effective and robust patch-wise non-local attention model is developed. The **unfold** operation is to extract sliding local patches from a batched input feature map, while the **fold** operation is to combine an array of sliding local patches into a large feature map.

maps with different sizes of receptive field. The channel-wise attention is independently performed on these two outputs, and the results are added together. The whole process of local attention submodule can be represented by $\mathbf{F}_L^i = \mathcal{H}_L(\mathbf{F}_{FEM}^i)$.

### D. Non-local Attention Submodule

In this part, we will introduce our proposed patch-wise non-local attention model. To make a clear distinction, we first conduct a simple comparison and analysis of typical non-local methods, including patch-wise non-local operations and non-local neural networks [21] (pixel-wise). Then we will give the details of our proposed patch-wise non-local attention model.

**Comparison and Analysis.** Fig. 4 shows a visual comparison of various non-local operations. The blue parts and the red parts in Fig. 4 refer to the restored image content and its corresponding long-range dependent items, respectively, while the

yellow boundary refers to the search region of pixels/patches matching. In Fig. 4(a), the non-local neural networks [21] constructed the long-range dependence between pixels, and it applied a learnable embedding function to make the matching process adaptive. However, the long-range dependence based on pixels is unreliable in image restoration tasks due to image corruption. In Fig. 4(b), the classic non-local means [19] applied patch matching and pixel updating (only applied central pixels) to restore corrupted image content. This type of traditional method suffers from the drawback that parameters are fixed. Fig. 4(c) shows the operation of learnable patch-wise non-local methods, and we take N3Net [32] for example. This kind of method can perform patch matching and updating adaptively, but they are not efficient enough. The illustration of our method is shown in Fig. 4(d). In what follows, we will present our solution to make the patch-wise non-local
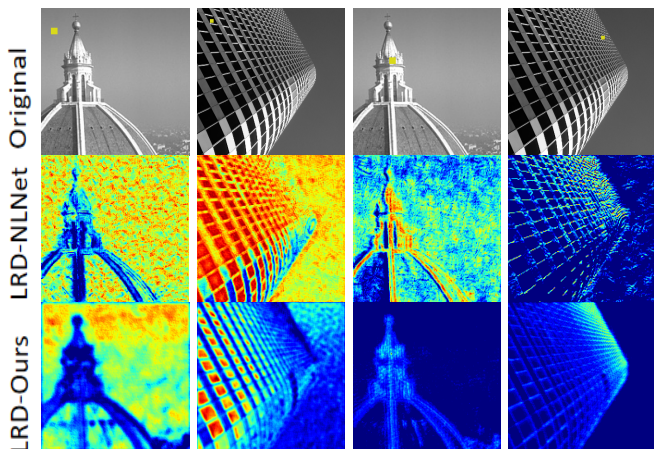
Fig. 6. Visualization of our patch-wise non-local attention operation and its counterpart (non-local neural networks [21]) which is represented as NLNet. The original images are shown in the first row and the query pixel/patch is labeled with a yellow dotted box. Long-range dependence (LRD) is shown below in the form of heat-maps. The brighter color indicates higher engagement.

operation more effective and efficient.

**Patch-wise Non-local Attention Model.** Based on the non-local neural networks [21] and existing patch-wise non-local operations, we propose a novel patch-wise non-local attention model. As illustrated in Fig. 5, given the input feature maps $\mathbf{F}_{FEM}^i \in \mathbb{R}^{C \times W \times H}$, we use three independent $1 \times 1$ convolutional layers with trainable parameters $\mathbf{W}_g$, $\mathbf{W}_\theta$, and $\mathbf{W}_\varphi$ as the embedding functions. These three embedding functions are represented as $g$, $\theta$, and $\varphi$. They are used to generate the query ($\mathbf{Q}$), key ($\mathbf{K}$) and value ($\mathbf{V}$). The embedding process does not change the size of feature maps, which can be defined as follows:

$$\begin{cases} \mathbf{Q} = \theta(\mathbf{F}_{FEM}^i) \\ \mathbf{K} = \varphi(\mathbf{F}_{FEM}^i) \\ \mathbf{V} = g(\mathbf{F}_{FEM}^i). \end{cases} \tag{10}$$

Rather than directly reshaping $\mathbf{Q}, \mathbf{K}$ and calculating the relationship in pixel level as [21], we propose to utilize the **unfold** operation to extract sliding local patches from these transformed feature maps with stride $s$ and patch-size of $W_p \times H_p$. Then, as shown in Fig. 5, we obtain three groups of 3D patches. Each group has $N$ patches with the size of each patch being $C \times W_p \times H_p$. The novelty of our proposed non-local attention model is to calculate the relationship in the unit of 3D patch, which is more effective and robust.

Next, we reshape each 3D patch unfolded from $\mathbf{Q}$ and $\mathbf{K}$ into a 1D feature vector, obtaining $\tilde{\mathbf{Q}}$ and $\tilde{\mathbf{K}}$. The distance matrix denoted by $\mathbf{M} \in \mathbb{R}^{N \times N}$ can be efficiently calculated by dot product:

$$\phi(\tilde{\mathbf{Q}}, \tilde{\mathbf{K}}) = softmax(\tilde{\mathbf{Q}}^T \tilde{\mathbf{K}}). \tag{11}$$

Each 3D patch in $\mathbf{V}$ can aggregate useful information to update itself guided by this distance matrix $\mathbf{M}$, which is essentially a weighted sum updating process in the unit of patch.

Finally, we utilize the **fold** operation to combine this array of updated sliding local patches into a feature map with

the size of $C \times W \times H$. This process can be viewed as the inverse process of the unfold operation. In consideration of overlapping between patches, we apply the averaging to deal with the overlapped areas. Compared with [33] applying a prediction neural network to infer the similarity between pairwise patches and [32] using a differentiable $k$-Nearest method to perform patch matching in a limited search region, our proposed patch-wise non-local attention model is more effective to fully utilize useful information from the whole image to restore each patch.

The pixel-wise non-local method [21] (corresponding to Fig. 4 (a)) is the counterpart of our method, which has been widely used in many image restoration tasks [14], [32]. To demonstrate that our proposed patch-wise non-local method can capture more reliable long-range dependence than [21] in the same case, we visualize their non-local attention maps in image denoising task ($\sigma = 30$), and the comparison is presented in Fig. 6. The heat-maps of our method are extracted from the second CAB of COLA-Net, and the results of its counterpart are extracted from the same place by replacing our patch-wise non-local method with the non-local neural networks [21]. Specifically, in image degradation, reliable correlations can aggregate more useful information to improve restoration results. From Fig. 6, it is obvious that the long-range dependence built based on pixels is susceptible to noisy signals. In comparison, our learnable patch-wise non-local attention model is scarcely influenced by noisy signals and can capture more reliable long-range dependence.

TABLE I
COMPUTATIONAL COMPLEXITY AND PARAMETER NUMBER
COMPARISON. PSNR AND SSIM VALUES ARE OBTAINED FROM
URBAN100 TEST SET [41] WITH NOISE LEVEL EQUALING TO 50.

| Algorithm | #Parameter ↓ | Running Time ↓ | PSNR/SSIM ↑ |
|---|---|---|---|
| DnCNN [4] | 0.56M | **0.03**s | 26.28/0.7874 |
| ADNet [13] | **0.52**M | 0.06s | 26.64/0.8072 |
| N3Net [32] | 0.72M | 0.08s | 26.82/0.8184 |
| NLRN [20] | 0.35M | 38.17s | 27.49/0.8279 |
| RNAN [14] | 8.96M | 1.25s | 27.65/0.8348 |
| COLA-B | 1.10M | 0.51s | 27.76/0.8373 |
| COLA-E | 1.88M | 0.68s | **27.84/0.8392** |

## IV. EXPERIMENTAL RESULTS

### A. Training Details and Evaluation

To verify the effectiveness of our proposed COLA-Net, we apply the basic version (COLA-B) and the enhanced version (COLA-E) of COLA-Net into three typical image restoration tasks: synthetic image denoising, real image denoising, and compression artifact reduction. For synthetic image denoising and compression artifact reduction tasks, we train COLA-Net with DIV2K [42] dataset, which contains 800 high-quality images. For real image denoising, we adopt training data containing 400 images corrupted by synthetic noise from BSD500 [43] dataset with the noise-level randomly selected from [10,50] and 120 images corrupted by realistic noise from RENOIR [44] dataset, which is similar to CBDNet [8]. Our model is trained on a GTX1080Ti GPU with the initial learning rate $lr = 1 \times 10^{-3}$ and performs halving per 200

TABLE II
QUANTITATIVE RESULTS (PSNR AND SSIM) ABOUT GRAY-SCALE IMAGE DENOISING. BEST RESULTS ARE **HIGHLIGHTED**.

| Dataset | $\sigma$ | BM3D [17] | DnCNN [4] | FFDNet [5] | ADNet [13] | N3Net [32] | NLRN [20] | COLA-B | COLA-E |
|---|---|---|---|---|---|---|---|---|---|
| Set12 | 15 | 32.37/0.8952 | 32.86/0.9031 | 32.75/0.9027 | 32.98/0.9044 | 33.03/0.9056 | 33.16/0.9070 | 33.20/0.9088 | **33.27/0.9097** |
| | 25 | 29.96/0.8504 | 30.44/0.8622 | 30.43/0.8634 | 30.58/0.8649 | 30.55/0.8648 | 30.80/0.8689 | 30.85/0.8701 | **30.90/0.8716** |
| | 50 | 26.70/0.7676 | 27.19/0.7829 | 27.31/0.7903 | 27.37/0.7903 | 27.43/0.7948 | 27.64/0.7980 | 27.73/0.8020 | **27.77/0.8032** |
| | 70 | 25.21/0.7176 | 25.56/0.7273 | 25.81/0.7451 | 25.63/0.7364 | 25.90/0.7510 | -/- | 26.18/0.7572 | **26.25/0.7606** |
| BSD68 | 15 | 31.07/0.8717 | 31.73/0.8907 | 31.63/0.8902 | 31.74/0.8910 | 31.78/0.8927 | 31.88/0.8932 | 31.88/0.8938 | **31.92/0.8968** |
| | 25 | 28.57/0.8013 | 29.23/0.8278 | 29.19/0.8289 | 29.25/0.8288 | 29.30/0.8321 | 29.41/0.8331 | 29.42/0.8339 | **29.46/0.8368** |
| | 50 | 25.62/0.6864 | 26.23/0.7189 | 26.29/0.7345 | 26.29/0.7210 | 26.39/0.7293 | 26.47/0.7298 | 26.48/0.7308 | **26.52/0.7340** |
| | 70 | 24.46/0.6323 | 24.85/0.6567 | 25.04/0.6700 | 24.95/0.6625 | 25.14/0.6753 | -/- | 25.17/0.6735 | **25.25/0.6788** |
| Urban100 | 15 | 32.35/0.9220 | 32.68/0.9255 | 32.43/0.9273 | 32.87/0.9304 | 33.08/0.9333 | 33.45/0.9354 | 33.60/0.9376 | **33.73/0.9387** |
| | 25 | 29.71/0.8777 | 29.97/0.8797 | 29.92/0.8887 | 30.24/0.8920 | 30.19/0.8925 | 30.94/0.9018 | 31.17/0.9062 | **31.33/0.9086** |
| | 50 | 25.95/0.7791 | 26.28/0.7874 | 26.52/0.8057 | 26.64/0.8072 | 26.82/0.8184 | 27.49/0.8279 | 27.76/0.8373 | **27.84/0.8392** |
| | 70 | 24.27/0.7165 | 24.34/0.7178 | 24.87/0.7495 | 24.53/0.7304 | 25.15/0.7658 | -/- | 25.99/0.7848 | **26.15/0.7910** |
| Parameters | | - | 0.56M | 0.49M | 0.52M | 0.72M | 0.35M | 1.10M | 1.88M |

epochs. During training, we employ Adam optimizer and each mini-batch contains 32 images with size of $64 \times 64$ randomly cropped from training data. A data augmentation method, which is the same as RNAN [14] is also applied in the training process. For each task, we apply commonly used test sets for testing and report PSNR (dB) and/or SSIM [45] to evaluate the performance of each method.

TABLE III
PSNR AND SSIM COMPARISON WITH RNAN [14] ON GRAY-SCALE IMAGE DENOISING. BEST RESULTS ARE **HIGHLIGHTED**.

| Dataset | $\sigma$ | RNAN [14] | COLA-B | COLA-E |
|---|---|---|---|---|
| BSD68 | 10 | 34.04/0.9295 | 34.03/0.9290 | **34.08/0.9299** |
| | 30 | 28.61/0.8094 | 28.62/0.8096 | **28.64/0.8110** |
| | 50 | 26.48/0.7306 | 26.48/0.7308 | **26.52/0.7340** |
| | 70 | 25.18/0.6746 | 25.17/0.6735 | **25.25/0.6788** |
| Urban100 | 10 | 35.52/0.9553 | 35.56/0.9554 | **35.64/0.9559** |
| | 30 | 30.20/0.8902 | 30.26/0.8910 | **30.41/0.8936** |
| | 50 | 27.65/0.8348 | 27.76/0.8373 | **27.84/0.8392** |
| | 70 | 25.89/0.7835 | 25.99/0.7848 | **26.15/0.7910** |
| Parameters | | 8.96M | 1.10M | 1.88M |

### B. Comparisons with State-of-the-Art Methods

We first make a brief complexity analysis of some representative methods and then compare our proposed COLA-Net with some recent state-of-the-art methods in synthetic image denoising, real image denoising, and compression artifact reduction applications.

*1) Computational Complexity:* Realizing that the number of trainable parameters can not completely reflect the complexity of a model, especially for deep non-local methods. Thus, we employ both the number of trainable parameters and the running time that a model takes to process a $256 \times 256$ image to evaluate the complexity of different methods. Note that all running times are tested on GPU. We also provide the denoising performance ($\sigma = 50$) on the Urban100 [41] dataset of various methods. The complexity analysis of several representative methods is reported in Table I. Compared with some competitive non-local-based denoiser, *e.g.*, RNAN [14], which stacked a lot of convolutional layers with the help of residual connection, and NLRN [20], which applied a recurrent strategy to cycle through one input many times, both the basic

version (COLA-B) and enhanced version (COLA-E) of our COLA-Net have an attractive complexity and achieve quite promising performance.

*2) Synthetic Image Denoising:* For this application, three standard benchmark datasets, i.e., Set12, BSD68 [48], and Urban100 [41] are evaluated. We compare the two versions of our COLA-Net with some state-of-the-art denoising methods, including some well-known denoisers, *e.g.*, DnCNN [4] and FFDNet [5], and recent local attention-based method ADNet [13] as well as some competitive non-local denoisers such as BM3D [17], N3Net [32] and NLRN[20]. Additive white Gaussian noise (AWGN) with different noise levels (15, 25, 50, 70) is added to the clean images. The quantitative results (PSNR and SSIM ) are shown in Table II, which clearly shows that the basic COLA-Net achieves the best performance in all three datasets, and the enhanced version can further improve the performance. It is worth noting that the superiority of our COLA-Net is more obvious in the case of strong noise. Especially for non-local methods, high-intensity noisy signals will disturb the construction of long-range dependence; thus, in this case, computing self-similarity based on pixels [20] or within a limited area [32], [33] is not reliable. The visual comparisons of denoising results of different methods are shown in Fig. 7. One can observe that our COLA-Net produces higher visual quality than other methods. Specifically, the non-local denoising methods [17], [20] produce the restoration results with over-smooth artifacts, and the local attention method [13] output oversharp results. In comparison, our method generates more accurate textures, demonstrating the effectiveness of our proposed mixed feature attention mechanism.

To further prove the superiority of COLA-Net, we also compare our method with recent very deep non-local neural networks RNAN [14], with the quantitative results shown in Table III. We can see that compared with RNAN, the basic COLA-Net obtains better performance in most cases, and the enhanced COLA-Net gets better performance in all test sets and noise levels. Simultaneously, the number of trainable parameters in two versions of COLA-Net is much lower than RNAN.

*3) Compression Artifact Reduction:* For this application, we compare our COLA-Net with some competitive methods: SA-DCT [46], ARCNN [38], TNRD [47], DnCNN [4], and
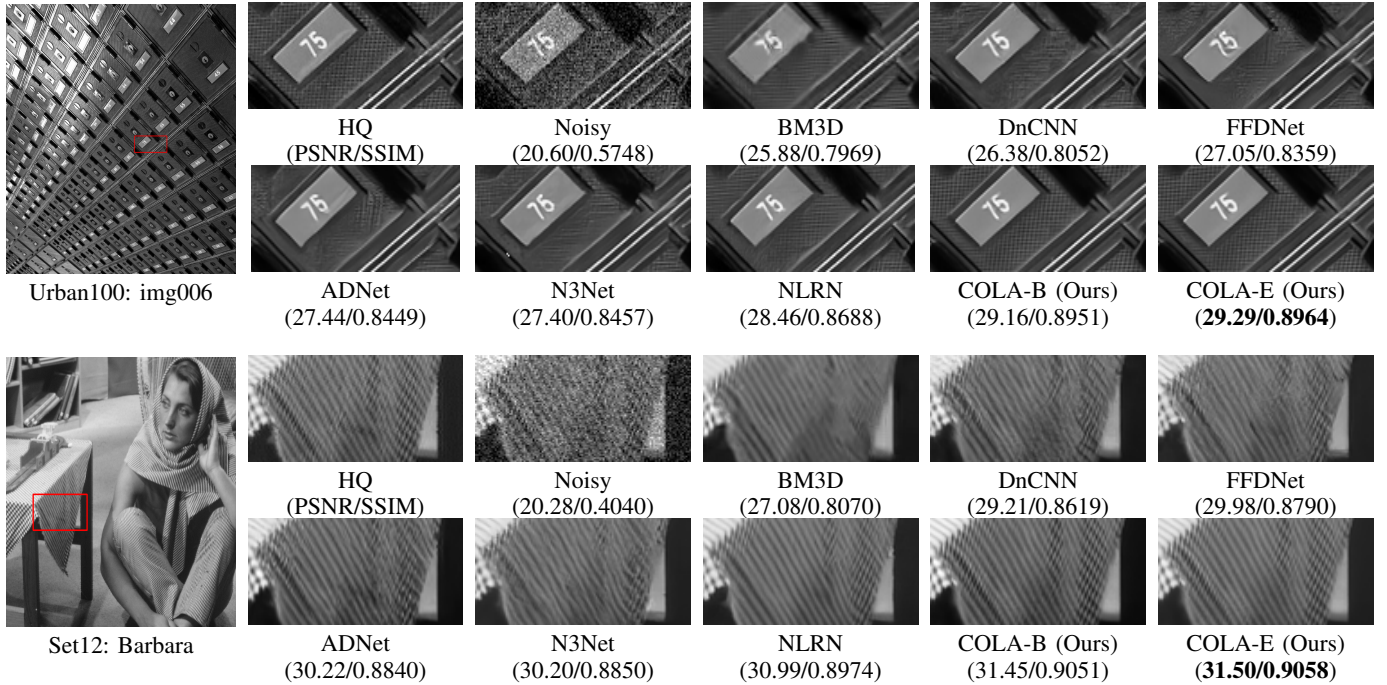
Fig. 7. Visual comparison for gray image denoising of various methods on two samples from Urban100 and Set12 with noise level $\sigma = 25$.

TABLE IV
QUANTITATIVE RESULTS (PSNR AND SSIM) OF COMPRESSION ARTIFACT REDUCTION. BEST RESULTS ARE **HIGHLIGHTED**.

| Dataset | $q$ | JPEG | SA-DCT [46] | ARCNN [38] | TNRD [47] | DnCNN [4] | RNAN [14] | COLA-B | COLA-E |
|---------|-----|------|-------------|------------|-----------|-----------|-----------|--------|--------|
| LIVE1 | 10 | 27.77/0.7905 | 28.86/0.8093 | 28.98/0.8076 | 29.15/0.8111 | 29.19 /0.8123 | <u>29.63/0.8239</u> | 29.60/0.8226 | **29.66/0.8234** |
|  | 20 | 30.07/0.8683 | 30.81/0.8781 | 31.29/0.8733 | 31.46/0.8769 | 31.59/0.8802 | <u>32.03/0.8877</u> | 31.98/0.8865 | **32.06/0.8880** |
|  | 30 | 31.41/0.9000 | 32.08/0.9078 | 32.69/0.9043 | 32.84/0.9059 | 32.98/0.9090 | <u>33.45/0.9149</u> | 33.40/0.9141 | **33.48/0.9152** |
|  | 40 | 32.35/0.9173 | 32.99/0.9240 | 33.63/0.9198 | -/- | 33.96/0.9247 | <u>34.47/0.9299</u> | 34.43/0.9293 | **34.49/0.9300** |
| Classic5 | 10 | 27.82/0.7800 | 28.88/0.8071 | 29.04/0.7929 | 29.28/0.7992 | 29.40/0.8026 | 29.96/0.8178 | <u>29.96/0.8180</u> | **30.03/0.8184** |
|  | 20 | 30.12/0.8541 | 30.92/0.8663 | 31.16/0.8517 | 31.47/0.8576 | 31.63/0.8610 | 32.11/0.8693 | <u>32.18/0.8695</u> | **32.28/0.8706** |
|  | 30 | 31.48/0.8844 | 32.14/0.8914 | 32.52/0.8806 | 32.74/0.8837 | 32.91/0.8861 | 33.38/0.8924 | <u>33.48/0.8929</u> | **33.54/0.8935** |
|  | 40 | 32.43/0.9011 | 33.00/0.9055 | 33.34/0.8953 | -/- | 33.77/0.9003 | 34.27/0.9061 | <u>34.33/0.9063</u> | **34.38/0.9066** |
| Parameters | | - | - | 0.12M | - | 0.56M | 8.96M | 1.10M | 1.88M |

RNAN [14]. The compressed images are generated by Matlab standard JPEG encoder with quality factor q = 10, 20, 30, 40. We evaluate the performance on LIVE1 [49] and Classic5 [46] test sets, and the quantitative results are presented in Table IV. One can see that our proposed COLA-E achieves the best performance in all test sets under the evaluation of both PSNR and SSIM [45], and the COLA-B presents an attractive performance maintaining fewer parameters. Visual comparison of compression artifacts reduction is shown in Fig. 8. We can find that our proposed model has better visual quality than other methods. Compared with the very deep non-local method [14], our image restoration results have fewer over-smooth artifacts, *e.g.,* in the region of shutters and wicker chair. The superiority of our method benefits from the fact that our proposed COLA-Net can balance the contribution of local and non-local attention operations according to the characteristic of specific image content.

*4) **Real Image Denoising**:* To further prove the merits of our proposed COLA-Net, we apply it to the more challenging task of real image denoising. Different from synthetic image denoising, in this case, images are corrupted by realistic noise,

which can be well explained by a Poisson-Gaussian distribution and can be further approximated with a heteroscedastic Gaussian distribution defined as below:

$$\begin{cases} n(L) \sim \mathcal{N}(0, \sigma^2(L)) \\ \sigma^2(L) = L \cdot \sigma_s^2 + \sigma_c^2, \end{cases} \qquad (12)$$

where $L$ is the irradiance intensity of raw pixels, $\sigma_s^2$ and $\sigma_c^2$ refer to the spatially variant noise variance and stationary noise variance, respectively. Since it is generally expensive to acquire hundreds of noisy images with realistic noise and their corresponding high-quality labels, existing real noise datasets are relatively small which can not be used to train neural networks adequately. Thus, this approximate noise model is widely used to perform data argumentation in real image denoising tasks.

During training, the heteroscedastic Gaussian distribution (Eq. 12) is used to generate the synthetic training samples. The commonly used DND [51] dataset is used for evaluation, which contains 50 images corrupted by realistic noise, and their high-quality labels are not available. We get the quantitative results from the benchmark website:
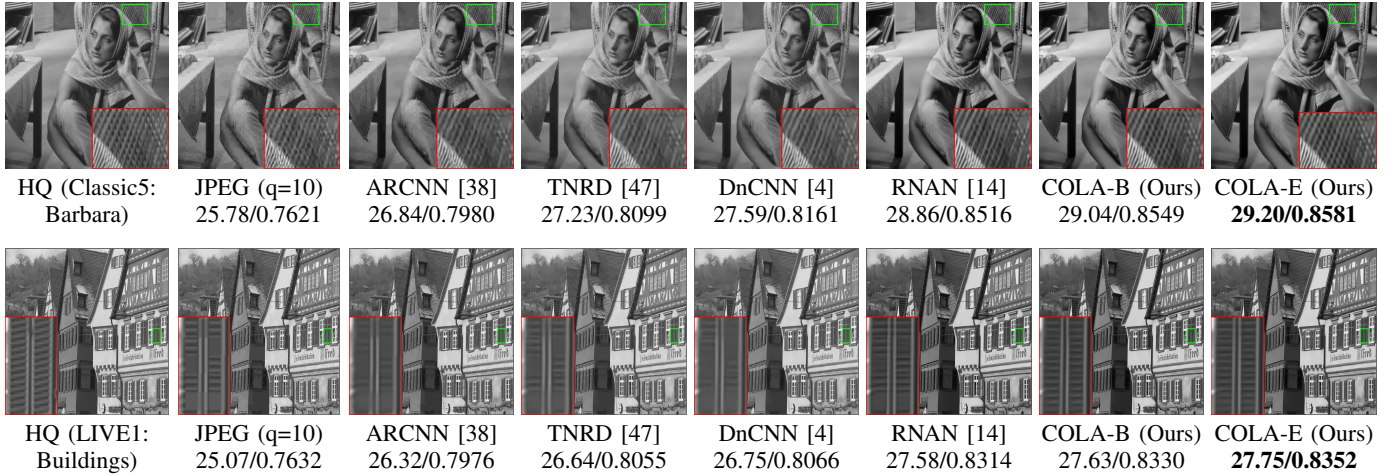
Fig. 8. Visual comparison of image compression artifact reduction application of various methods with JPEG quality q = 10.
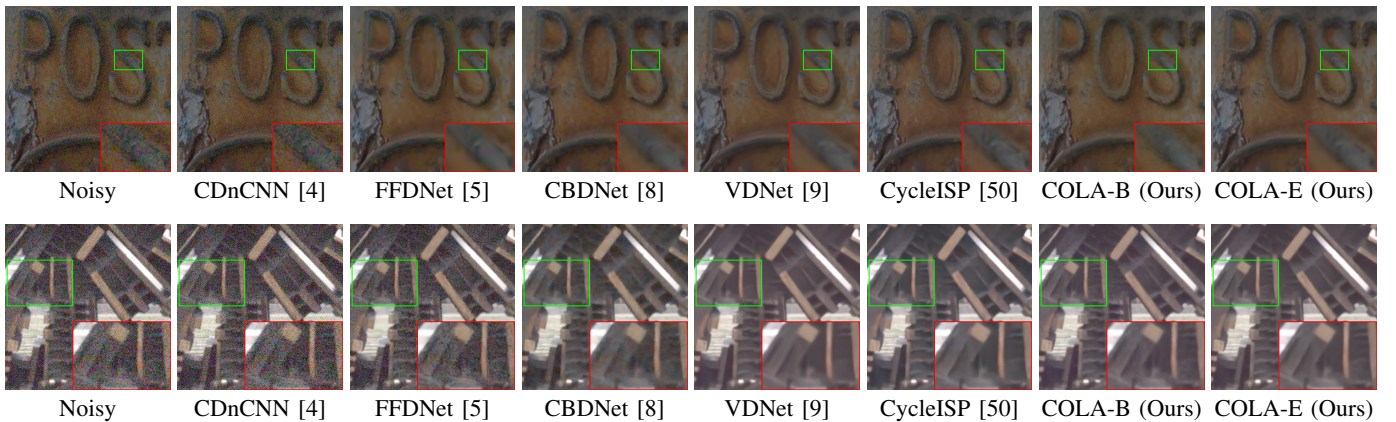


Fig. 9. Visual comparison of real image denoising application of various methods. These noisy images are corrupted by realistic noise from DND [51] dataset.

TABLE V
THE QUANTITATIVE RESULTS ON THE DND BENCHMARK.

| Algorithm | Parameters | Blind/Non-blind | sRGB PSNR | sRGB SSIM |
|-----------|-----------|-----------------|-----------|-----------|
| BM3D [17] | - | Non-blind | 34.51 | 0.851 |
| CDnCNN [4] | 0.67M | Blind | 32.43 | 0.790 |
| CFFDNet [5] | 0.85M | Non-blind | 37.61 | 0.914 |
| TWSC [52] | - | Blind | 37.94 | 0.940 |
| CBDNet [8] | 4.36M | Blind | 38.06 | 0.942 |
| VDNet [9] | 7.82M | Blind | 39.38 | 0.952 |
| CycleISP [50] | 2.60M | Blind | 39.56 | 0.956 |
| AINDNet [11] | 13.76M | Blind | 39.77 | **0.959** |
| MIRNet [23] | 31.79M | Blind | **39.88** | 0.956 |
| COLA-B | 1.18M | Blind | 39.07 | 0.949 |
| COLA-E | 1.92M | Blind | 39.64 | 0.954 |

**https://noise.visinf.tu-darmstadt.de/**, and also release the official evaluation result of our method on this website. The objective comparison is shown in Table V. It is clear to see that the basic COLA-Net (COLA-B) achieves attractive performance with much fewer parameters, and the enhanced COLA-Net (COLA-E) outperforms some very recent competitive methods (*e.g.,* CBDNet [8], VDNet [9] and CycleISP [50]) while maintaining an attractive computational complexity. Even compared with existing top-performing methods [11], [23], our method can achieve comparable performance with

much fewer parameters (only about $\frac{1}{8}$ of [11] and $\frac{1}{17}$ of [23]). Fig. 9 further provides a visual comparison of different methods on two samples from DND dataset. Compared with other methods, COLA-Net is able to produce higher visual quality with much sharper image contents and fewer noisy signals, especially in the boundary areas, which benefits from the proposed collaborative attention mechanism.

### C. Ablation Study

In this subsection, we show the ablation study in Table VI to investigate the effect of different components in COLA-Net. Note that Case 0 denotes our basic COLA-Net with the default setting. In the ablation study, we compare our method with several baseline models: ($Case1$) We replace the dual-branch fusion modules with non-local neural networks [21]. ($Case2$) We replace our patch-wise non-local attention submodule with non-local neural networks [21]. ($Case3$) We remove all patch-wise non-local attention submodules from our COLA-Net. ($Case4$) We remove all local attention submodules. ($Case5$) We remove all dual-branch fusion modules. The detailed analysis of the ablation study is presented below.

- **Non-local attention.** Compared with Case 0, Case 2 removes the non-local attention submodule. The obvious performance decrease in Case 2 indicates the positive

TABLE VI
ABLATION STUDY OF DIFFERENT COMPONENTS IN COLA-NET. PSNR VALUES ARE EVALUATED ON URBAN100 ($\sigma = 25$).

| Case Index | 0 (default) | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| Patch-wise Non-local Attention | ✓ | × | × | ✓ | × | × | ✓ | ✓ | ✓ |
| Local Attention | ✓ | × | ✓ | × | ✓ | × | ✓ | ✓ | ✓ |
| Pixel-wise Non-local Attention [21] | × | ✓ | × | × | ✓ | × | × | × | × |
| Number of CAB | 4 | 4 | 4 | 4 | 4 | 4 | 3 | 2 | 1 |
| PSNR (dB) | 31.17 | 30.70 | 30.18 | 31.04 | 30.75 | 29.79 | 31.12 | 31.01 | 30.68 |
| Parameters | 1.10M | 0.98M | 1.00M | 1.06M | 1.02M | 0.88M | 0.97M | 0.70M | 0.42M |

effect of our proposed non-local attention operation. In Case 3, we only apply non-local attention submodule, and in Case 5, we remove both local and non-local attention submodules. The performance comparison between Case 3 and Case 5 demonstrates the necessity of the non-local attention submodule.

- **Local Attention.** To verify the effect of local attention submodules, we also provide two groups of comparison. From Case 0 and Case 3, one can see that local attention submodule brings 0.13 dB PSNR gains together with patch-wise non-local attention submodule. From Case 2 and Case 5, the local attention submodule contributes 0.39 dB PSNR gains.

- **Patch-wise Non-local Attention.** To verify the effect of our proposed patch-wise non-local attention model, we make a comparison between Case 0 and Case 4, as well as between Case 3 and case 1. The 0.42 dB PSNR gains by Case 0 over Case 4 and 0.34dB gains by Case 3 over Case 1 clearly show the superiority of our proposed patch-wise non-local attention operation over existing pixel-wise non-local attention operation [21].

- **Collaborative Attention.** To verify the necessity of collaborative attention, we compare it with only local attention and only non-local attention separately. The 0.99 dB gains by Case 0 over Case 2 and the 0.13 dB gains by Case 0 over Case 3 fully demonstrate the superiority of our proposed collaborative attention block.

- **Number of CAB.** To study the effect of the number of CAB, we make a comparison among Cases 0, 6, 7, and 8, which clearly shows that the performance increases with the number of CAB. By making a trade-off between performance and complexity, the default number of CAB is set to be 4. It is also worth emphasizing that our proposed COLA-Net with only one CAB, that is, Case 8, still achieves better performance than ADNet [13] and N3Net [32], despite having fewer parameters (425K vs. 560K and 720K) and fewer layers (12 vs. 17 and 24).

## V. CONCLUSION

In this paper, a novel **COL**laborative **A**ttention **Net**work, dubbed **COLA-Net** is proposed by incorporating both local and non-local attention mechanisms into deep network for image restoration tasks. COLA-Net is the first attempt to adaptively combine local and non-local operations and to enable restoring complex textures and repetitive details distinguishingly. An effective and robust patch-wise non-local attention model is also developed to establish a more reliable long-range

dependence during image restoration. Extensive experiments on three typical image restoration tasks, i.e., synthetic image denoising, real image denoising and compression artifact reduction, show that the proposed COLA-Net achieves state-of-the-art results, while maintaining an attractive computational complexity.

## REFERENCES

[1] J. Zhang, D. Zhao, R. Xiong, S. Ma, and W. Gao, "Image restoration using joint statistical modeling in a space-transform domain," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 24, no. 6, pp. 915–928, 2014.

[2] J. Zhang, C. Zhao, D. Zhao, and W. Gao, "Image compressive sensing recovery using adaptively learned sparsifying basis via l0 minimization," *Signal Processing*, vol. 103, pp. 114–126, 2014.

[3] C. Zhao, J. Zhang, R. Wang, and W. Gao, "CREAM: CNN-regularized ADMM framework for compressive-sensed image reconstruction," *IEEE Access*, vol. 6, pp. 76 838–76 853, 2018.

[4] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a Gaussian denoiser: Residual learning of deep CNN for image denoising," *IEEE Transactions on Image Processing*, vol. 26, no. 7, pp. 3142–3155, 2017.

[5] K. Zhang, W. Zuo, and L. Zhang, "FFDNet: toward a fast and flexible solution for CNN-based image denoising," *IEEE Transactions on Image Processing*, vol. 27, no. 9, pp. 4608–4622, 2018.

[6] K. Zhang, W. Zuo, S. Gu, and L. Zhang, "Learning deep CNN denoiser prior for image restoration," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3929–3938.

[7] Y. Tai, J. Yang, X. Liu, and C. Xu, "MemNet: A persistent memory network for image restoration," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 4539–4547.

[8] S. Guo, Z. Yan, K. Zhang, W. Zuo, and L. Zhang, "Toward convolutional blind denoising of real photographs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 1712–1722.

[9] Z. Yue, H. Yong, Q. Zhao, D. Meng, and L. Zhang, "Variational denoising network: Toward blind noise modeling and removal," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 1688–1699.

[10] T. Brooks, B. Mildenhall, T. Xue, J. Chen, D. Sharlet, and J. T. Barron, "Unprocessing images for learned raw denoising," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 11 036–11 045.

[11] Y. Kim, J. W. Soh, G. Y. Park, and N. I. Cho, "Transfer learning from synthetic to real-noise denoising with adaptive instance normalization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 3482–3492.

[12] T. Wang, M. Sun, and K. Hu, "Dilated deep residual network for image denoising," in *Proceedings of the 2017 IEEE 29th International Conference on Tools with Artificial Intelligence*, 2017, pp. 1272–1279.

[13] C. Tian, Y. Xu, Z. Li, W. Zuo, L. Fei, and H. Liu, "Attention-guided CNN for image denoising," *Neural Networks*, vol. 124, pp. 117–129, 2020.

[14] Y. Zhang, K. Li, K. Li, B. Zhong, and Y. Fu, "Residual non-local attention networks for image restoration," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2019, pp. 1–18.

[15] X. Mao, C. Shen, and Y.-B. Yang, "Image restoration using very deep convolutional encoder-decoder networks with symmetric skip connections," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2016, pp. 2802–2810.

[16] J. Cai, Z. Meng, and C. Man Ho, "Residual channel attention generative adversarial network for image super-resolution and noise reduction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2020, pp. 454–455.

[17] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-D transform-domain collaborative filtering," *IEEE Transactions on Image Processing*, vol. 16, no. 8, pp. 2080–2095, 2007.

[18] J. Zhang, D. Zhao, and W. Gao, "Group-based sparse representation for image restoration," *IEEE Transactions on Image Processing*, vol. 23, no. 8, pp. 3336–3351, 2014.

[19] A. Buades, B. Coll, and J.-M. Morel, "A non-local algorithm for image denoising," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2005, pp. 60–65.

[20] D. Liu, B. Wen, Y. Fan, C. C. Loy, and T. S. Huang, "Non-local recurrent network for image restoration," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 1673–1682.

[21] X. Wang, R. Girshick, A. Gupta, and K. He, "Non-local neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7794–7803.

[22] S. Anwar and N. Barnes, "Real image denoising with feature attention," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019, pp. 3155–3164.

[23] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, M.-H. Yang, and L. Shao, "Learning enriched features for real image restoration and enhancement," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 492–511.

[24] C. Mou and J. Zhang, "Synergic feature attention for image restoration," in *Proceedings of the IEEE International Conference on Acoustics, Speech, & Signal Processing (ICASSP)*, 2021, pp. 1–5.

[25] H. Liu, R. Xiong, J. Zhang, and W. Gao, "Image denoising via adaptive soft-thresholding based on non-local samples," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 484–492.

[26] C. Zhao, J. Zhang, S. Ma, X. Fan, Y. Zhang, and W. Gao, "Reducing image compression artifacts by structural sparse representation and quantization constraint prior," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 27, no. 10, pp. 2057–2071, 2016.

[27] J. Zhang, R. Xiong, C. Zhao, Y. Zhang, S. Ma, and W. Gao, "CONCOLOR: Constrained non-convex low-rank model for image deblocking," *IEEE Transactions on Image Processing (TIP)*, vol. 25, no. 3, pp. 1246–1259, 2016.

[28] S. Lefkimmiatis, "Non-local color image denoising with convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 3587–3596.

[29] P. Qiao, Y. Dou, W. Feng, R. Li, and Y. Chen, "Learning non-local image diffusion for image denoising," in *Proceedings of the 25th ACM International Conference on Multimedia*, 2017, pp. 1847–1855.

[30] S. Lefkimmiatis, "Universal denoising networks: a novel CNN architecture for image denoising," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 3204–3213.

[31] B. Lecouat, J. Ponce, and J. Mairal, "Fully trainable and interpretable non-local sparse models for image restoration," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020, pp. 238–254.

[32] T. Plötz and S. Roth, "Neural nearest neighbors networks," in *Proceedings of the Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 1087–1098.

[33] Z. Xia and A. Chakrabarti, "Identifying recurring patterns with deep neural networks for natural image denoising," in *Proceedings of the IEEE Winter Conference on Applications of Computer Vision (WACV)*, 2020, pp. 2426–2434.

[34] S. Woo, J. Park, J.-Y. Lee, and I. So Kweon, "CBAM: Convolutional block attention module," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 3–19.

[35] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 7132–7141.

[36] X. Li, W. Wang, X. Hu, and J. Yang, "Selective kernel networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 510–519.

[37] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 286–301.

[38] C. Dong, Y. Deng, C. Change Loy, and X. Tang, "Compression artifacts reduction by a deep convolutional network," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015, pp. 576–584.

[39] Z. He, Y. Cao, L. Du, B. Xu, J. Yang, Y. Cao, S. Tang, and Y. Zhuang, "MRFN: Multi-receptive-field network for fast and accurate single image super-resolution," *IEEE Transactions on Multimedia*, vol. 22, no. 4, pp. 1042–1054, 2019.

[40] C. Tian, Y. Xu, W. Zuo, B. Zhang, L. Fei, and C.-W. Lin, "Coarse-to-fine CNN for image super-resolution," *IEEE Transactions on Multimedia*, vol. 30, pp. 976–985, 2020.

[41] J.-B. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 5197–5206.

[42] R. Timofte, E. Agustsson, L. Van Gool, M.-H. Yang, and L. Zhang, "Ntire 2017 challenge on single image super-resolution: Methods and results," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2017, pp. 114–125.

[43] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2001, pp. 416–423.

[44] J. Anaya and A. Barbu, "RENOIR-a dataset for real low-light noise image reduction," *Journal of Visual Communication of Image Representation*, vol. 51, pp. 144–154, 2018.

[45] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.

[46] A. Foi, V. Katkovnik, and K. Egiazarian, "Pointwise shape-adaptive DCT for high-quality denoising and deblocking of grayscale and color images," *IEEE Transactions on Image Processing*, vol. 16, no. 5, pp. 1395–1411, 2007.

[47] Y. Chen and T. Pock, "Trainable nonlinear reaction diffusion: A flexible framework for fast and effective image restoration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 6, pp. 1256–1272, 2016.

[48] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *Proceedings of the Eighth IEEE International Conference on Computer Vision (ICCV)*, 2001, pp. 416–423.

[49] H. Sheikh, "Live image quality assessment database release 2," *http://live. ece. utexas. edu/research/quality*, 2005.

[50] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, M.-H. Yang, and L. Shao, "Cycleisp: Real image restoration via improved data synthesis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 2696–2705.

[51] T. Plotz and S. Roth, "Benchmarking denoising algorithms with real photographs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 1586–1595.

[52] J. Xu, L. Zhang, and D. Zhang, "A trilateral weighted sparse coding scheme for real-world image denoising," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 20–36.