

R2RNet: Low-light Image Enhancement via Real-low to Real-normal Network

Jiang Hai, Zhu Xuan, Songchen Han, Ren Yang, Yutong Hao, Fengzhu Zou, and Fang Lin

Abstract—Images captured in weak illumination conditions could seriously degrade the image quality. Solving a series of degradation of low-light images can effectively improve the visual quality of images and the performance of high-level visual tasks. In this study, a novel Retinex-based Real-low to Real-normal Network (R2RNet) is proposed for low-light image enhancement, which includes three subnets: a Decom-Net, a Denoise-Net, and a Relight-Net. These three subnets are used for decomposing, denoising, contrast enhancement and detail preservation, respectively. Our R2RNet not only uses the spatial information of the image to improve the contrast but also uses the frequency information to preserve the details. Therefore, our model achieved more robust results for all degraded images. Unlike most previous methods that were trained on synthetic images, we collected the first Large-Scale Real-World paired low/normal-light images dataset (LSRW dataset) to satisfy the training requirements and make our model have better generalization performance in real-world scenes. Extensive experiments on publicly available datasets demonstrated that our method outperforms the existing state-of-the-art methods both quantitatively and visually. In addition, our results showed that the performance of the high-level visual task (*i.e.* face detection) can be effectively improved by using the enhanced results obtained by our method in low-light conditions. Our codes and the LSRW dataset are available at: <https://github.com/abcdef2000/R2RNet>.

Index Terms—Retinex, Low-light image enhancement, Image processing, Real-world low/normal-light image pairs.

I. INTRODUCTION

INSUFFICIENT illumination in the image capturing seriously affects the image quality from many aspects, such as low contrast and low visibility. Removing these degradations and transforming a low-light image into a high-quality sharp image is helpful to improve the performance of high-level visual tasks, such as image recognition [1], object detection [2], semantic segmentation [3], *etc.*, and can also improve the performance of intelligent systems in some practical applications, such as autonomous driving, visual navigation [4], *etc.* Low-light image enhancement, therefore, is highly desired.

Over the past few decades, there have been a large number of methods employed to enhance degraded images captured under insufficient illumination conditions. These methods have made great progress in improving image contrast and can obtain enhanced images with better visual quality. In addition

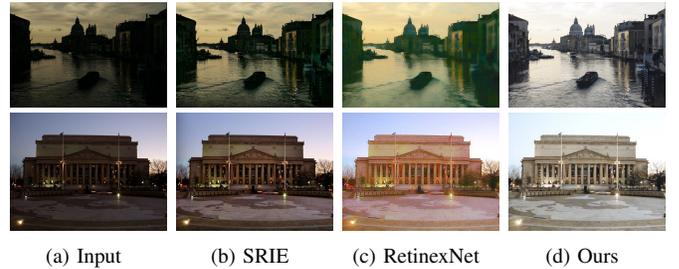


Fig. 1. Examples of low-light image enhanced results. The proposed method can not only improve the contrast of the image but also suppress the noise and artifacts in the dark regions. SRIE generates under-enhancement results and RetinexNet generates results of blur and color distortion.

to contrast, another special degradation of low-light images is noise. Many methods utilized additional denoising methods as pre-processing or post-processing. However, using denoising methods as pre-processing will cause blurring, while applying denoising as post-processing will result in noise amplification [5]. Recently, some methods [6] have designed effective models to perform denoising and contrast enhancement simultaneously and obtain satisfactory results.

It is noteworthy that many previous methods focused on using the spatial domain information of the image for enhancement, and image processing in frequency domain is also one of the important methods in the image enhancement field. High-frequency information usually represents the image details (*e.g.* contour and edge) or noise, so we proposed a novel Real-low to Real-normal Network for low-light image enhancement, dubbed R2RNet, which utilizes both spatial and frequency information to obtain high visual quality enhanced images. The design of our network is built on the Retinex theory [7], which includes three subnets: a Decom-Net, a Denoise-Net, and a Relight-Net. The Decom-Net aims to decompose the input low-light image into an illumination map and a reflectance map under the guidance of the Retinex theory. The Denoise-Net takes the decomposition results as input and uses the illumination map as a constraint to suppress the noise in the reflectance to obtain decomposition results with better visual quality. The illumination map obtained by Decom-Net and the reflectance map obtained by Denoise-Net are sent to Relight-Net to improve the contrast and brightness of the image. In Decom-Net and Denoise-Net, we only utilized the spatial information of low-light images, because the purpose of Decom-Net is to decompose the input image into an illumination map and a reflectance map without any further processing. According to the Retinex theory, the reflectance map contains

This work is partially supported by Key R&D project of Sichuan Province, China (No.22ZDYF3720). (*Corresponding author: Songchen Han.*)

The authors are with the School of Aeronautics and Astronautics, Sichuan University, Chengdu 610065, China (e-mail: j1269998232@163.com, j_zhxxx@163.com, renyang20212021@163.com, 15522808379@163.com, Y15286544560@163.com, photoatoms@163.com, hansongchen@scu.edu.cn).

the inherent attributes of the image, so if the frequency information is used to suppress the noise in Denoise-Net, the details of the reflectance map may be suppressed at the same time. Therefore, instead of using the frequency information in Decom-Net and Denoise-Net, we used the spatial information of the image to improve the image contrast and extracted the frequency information of the image based on the fast Fourier transform to better preserve the image details in Relight-Net. By a well-designed network, our method could appropriately enhance the image contrast, preserve more image details, and suppress noise. Moreover, the performance of high-level visual tasks can be effectively improved by using the enhanced results obtained by our method in low-light conditions.

Another difficulty in the low-light image enhancement task is that the learning-based model requires a lot of data for training, and the ability of the model is usually closely related to the quality of the training data. However, it is very difficult to collect sufficient real-world data, especially for paired images. Most learning-based enhancement methods use synthetic datasets for training, which limits their generalization ability in real-world scenarios. As far as we know, the existing real-world low-light image datasets only have the LOL dataset [8] and the SID dataset [9], but the number of images contained in these two datasets cannot satisfy the training requirement of deep neural networks. Therefore, we collected the first Large-Scale Real-World paired low/normal-light images dataset, named LSRW dataset, for our network training.

The remainder of this paper is organized as follows. Section II briefly reviews the relevant works of low-light enhancement methods, image denoising methods, and low-light image datasets. Section III presents the architecture of the proposed R2RNet and loss function settings. Section IV presents the experimental results, and Section V provides some concluding remarks.

II. RELATED WORK

A. Low-light Image Enhancement methods

In the past few decades, there have been extensive methods to enhance the contrast of weakly illuminated images. Traditional methods are mainly based on histogram equalization and the Retinex theory. Histogram equalization is a simple but effective image enhancement technology, it takes effect by changing the histogram of the image to enhance the contrast, such as brightness preserving bi-HE [10]. The Retinex theory assumes that the color image observed by human beings can be decomposed into an illumination map and a reflectance map, in which the reflectance map is the inherent attribute of the image and cannot be changed. The purpose of enhancing contrast can be achieved by changing the dynamic range of pixels in the illumination map. MSRCR [11] utilized a multi-scale Gaussian filter to restore color based on the Retinex theory. SRIE [12] proposed a weighted variational model to estimate reflectance and illumination at the same time. MF [13] tried to improve the local contrast of the illumination map and maintain naturalness. BIMEF [14] used a dual-exposure algorithm for image enhancement. In addition to

the illumination map and reflectance map, Mading *et al.* [15] added a noise map to form a robust Retinex model for further enhancement and denoising. LIME [16] first estimated the illumination by a prior hypothesis, obtained the estimated illumination through a weighted vibration model, subsequently used BM3D [17] as post-processing. Recently, Liu *et al.* [18] proposed a Retinex-inspired Unrolling with Architecture Search (RUAS) and designed a cooperative reference-free learning strategy to discover low-light prior architectures from a compact search space. Deep learning has been widely used in the field of computer vision and achieves excellent results. Many excellent methods, such as CNN, GAN, *etc.*, have made remarkable achievements in a variety of low-level visual tasks, including image de-hazing [19], [20], image super-resolution [21], [22], *etc.* Many researchers also build learning-based models based on the Retinex theory. MSR-Net [23] utilized different Gaussian convolution kernels to learn the mapping of low/normal-light images. RetinexNet [8] combined the Retinex theory with DeepCNN to estimate and adjust the illumination map to achieve image contrast enhancement and uses BM3D for post-processing to achieve denoising. Zhang *et al.* [24] also designed an effective network based on the Retinex theory to enhance low-light images. Lim *et al.* [25] proposed a deep-stacked Laplacian restorer (DSLRL) to recover the global illumination and local details from the original input. Furthermore, some methods that are not based on the Retinex theory are also proposed. Dong *et al.* [21] proposed an algorithm that improves the contrast in dark regions and improves the visual quality by using a de-hazing method. Ying *et al.* [26] used the camera response model for weakly illuminated image enhancement. Lore *et al.* [27] proposed a stacked sparse denoising autoencoder for image contrast enhancement and denoising. Ziaur *et al.* [28] used a bright channel prior to obtain an initial transmission map and adopted L1-norm regularization to refine scene transmission. Guo *et al.* [2] proposed a lightweight network named Zero-DCE, to transform the image enhancement problem into a curve estimation problem. Jiang *et al.* [1] proposed a network based on GAN for low-light image enhancement and used unpaired images for training for the first time.

The key to Retinex-based methods is the estimation of the illumination map and the reflectance map. Due to the limited decomposition ability, traditional methods often cause over/under-enhancement results. Learning-based methods can get better decomposition results and can properly improve the contrast. It is noteworthy that most learning-based methods only focus on using the spatial information of weakly illuminated images to obtain high-quality normal-light images, and combining spatial and frequency domain information to perform low-light image enhancement can obtain more satisfactory enhanced results. Therefore, our R2RNet uses both spatial and frequency information of the image for enhancement. Spatial information is used for contrast enhancement, and frequency information is used to restore more image details.

B. Denoising methods

Enhancing weakly illuminated images needs noise suppression in addition to contrast enhancement. The traditional image denoising methods relied on hand-crafted features and used discrete cosine transform or wavelet transform to modify transform coefficients. NLM [29] and BM3D used self-similar patches to achieve outstanding results in image fidelity and visual quality. Image denoising methods based on supervised learning, such as DnCNN-B [30], FFDNet [31], and CBDNet [32] utilized the Gaussian mixture model to perform denoising. Mei *et al.* [33] made full use of shallow pixel-level features and self-similarity to achieve a balance between pixel features and semantic features to preserve more details. Kim *et al.* [34] proposed CBAM to focus on learning the difference between noisy images and clear images. Chen *et al.* [35] used GAN to model the noise information extracted from the real noise image and combined the noise blocks generated by the generator with the original clear image to synthesize a new noise image. ADGAN [36] proposed a feature pyramid attention network to improve the ability of network feature extraction when modeling noise.

These methods can achieve impressive denoising results. However, directly using these methods as pre-processing or post-processing of low-light image enhancement methods will result in blurring or noise amplification. To avoid this, our method can perform contrast enhancement and denoising simultaneously.

C. Low-light Image Datasets

Another difficulty in the low-light image enhancement task is that learning-based models usually require a lot of data, but it is difficult to collect sufficient low-light images. Due to the lack of real-world paired images, most methods use synthetic images based on normal-light images. Lore *et al.* [27] applied gamma correction to each channel to synthesize low-light images. Lv *et al.* [37] used the same image synthesis strategy as LLNet. Lv *et al.* [38] and Wang *et al.* [39] combined linear transformation and gamma transformation to obtain paired images. Wang *et al.* [40] obtained synthetic images by using the camera response function and modeling the noise distribution in the low-light image.

As far as we know, the existing real-world low-light image datasets only have LOL dataset and SID dataset, both of them capture pairs of low/normal-light images by fixing the camera position and changing the ISO and exposure time. LOL dataset contains 500 low/normal-light image pairs. SID dataset contains 5094 short-exposure images and 424 long-exposure images; multiple short-exposure images correspond to one long-exposure image. However, the number of images contained in the above two datasets cannot support the training of DeepCNN, and the SID dataset is mainly suitable for extremely weakly illuminated image enhancement, which is not the same as what we focus on. In order to meet the training requirement of our network, we use a Nikon D7500 camera and a HUAWEI P40 Pro mobile phone to collect real-world paired images to form our LSRW dataset.

TABLE I
THE LSRW DATASET CONTAINS 5650 LOW/NORMAL-LIGHT IMAGE PAIRS CAPTURED BY A NIKON D7500 CAMERA AND A HUAWEI P40 PRO MOBILE PHONE IN REAL SCENES.

	Exposure Time(s)	ISO	Image Pairs
Nikon D7500	[1/200, 1/20]	[50, 100]	3170
Huawei P40 Pro	[1/400, 1/15]	[50, 100]	2480

III. LSRW DATASET

One of the difficulties in the task of low-light image enhancement is the lack of paired low/normal-light images captured in real scenes. The existing real-world paired images datasets are only the LOL dataset and the SID dataset, and SID is mainly suitable for extremely low-light image enhancement, which is not consistent with our concern. To satisfy the training requirements of DeepCNNs and provide support for follow-up researches, we propose the first large-scale real-world paired image dataset, named LSRW dataset. The LSRW dataset contains 5650 paired images captured by a Nikon D7500 camera and a HUAWEI P40 Pro mobile phone. We collected 3170 paired images using the Nikon camera and 2480 paired images using the Huawei mobile phone.

The low-light images can be obtained by reducing ISO and using a shorter exposure time to reduce the amount of light input, while the normal-light images can be obtained by using a larger ISO and a longer exposure time. We chose to collect the LSRW dataset for both indoor and outdoor scenarios. When obtaining low-light images in indoor scenes, the exposure time will be increased to avoid capturing extremely dark images. Similarly, when obtaining normal-light images in outdoor scenes, the exposure time will be reduced to avoid capturing over-exposure images. The ISO value of the low-light condition is fixed to 50, and the normal-light condition is fixed to 100. We can obtain paired low/normal-light images by changing the exposure time. Note that if there are moving objects or camera/phone shaking when reducing the exposure time, the low-light image will be blurred. Therefore, in order to avoid camera/mobile phone shaking, we use a tripod to fix the position of the camera/mobile phone and adjust the ISO and exposure time through long-range control. At the same time, the scenes we select are static without any moving objects, which can ensure that the captured low-light images will not be blurred. The ISO value of the low-light condition is fixed to 50, and the normal-light condition is fixed to 100. When using Nikon to obtain low-light images, the exposure time is limited to 1/200 to 1/80, while the exposure time of normal-light images is limited to 1/80 to 1/20. When using Huawei to obtain low-light images, the exposure time is limited to 1/400 to 1/100, while the exposure time of normal-light images is limited to 1/100 to 1/15. We selected 5600 paired images from the LSRW dataset for training and the remaining 50 pairs for evaluation. Table I summarizes the LSRW dataset. Fig.2 shows several image pairs in the LSRW dataset, including indoor and outdoor scenes.

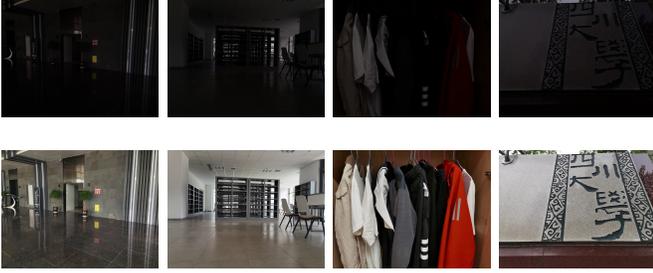


Fig. 2. Several examples for paired low/normal-light images in the LSRW dataset.

IV. METHOD

In this section, we will introduce the details of our R2RNet, including the network architecture and the loss function.

A. Network Architecture

We propose a novel deep convolutional neural network, dubbed R2RNet, which consists of three subnets: a Decom-Net, a Denoise-Net, and a Relight-Net. The Decom-Net decomposes the input weakly illuminated image into an illumination map and a reflectance map based on the Retinex theory. The Denoise-Net takes the decomposed results as input and uses the illumination map as a constraint to suppress the noise in the reflectance map. Subsequently, the illumination map obtained by Decom-Net and the reflectance map obtained by Denoise-Net are sent to the Relight-Net to obtain a normal-light image with better visual quality. Therefore, our method can improve the contrast, retain more details, and suppress the noise simultaneously. The network architecture of R2RNet is illustrated in Fig.3. The detailed descriptions are provided below.

Decom-Net: The key of Retinex-based methods is to obtain the high-quality illumination map and reflectance map, the quality of decomposition results will also affect the subsequent enhancement and denoising process. Therefore, it is important to design an efficient network to decompose the weakly illuminated image. Residual network [41] has been widely used in many computer vision tasks and achieved excellent results. Benefiting from the skip connection structure, the residual network can make the deep neural network easier to optimize during the training stage, and not cause gradient disappearance or explosion. Inspired by this, we use multiple Residual Modules (RM) in DecomNet to get better decomposition results. Each RM contains 5 convolutional layers with the kernel size of $\{1, 3, 3, 3, 1\}$, the number of kernels is $\{64, 128, 256, 128, 64\}$, respectively. And we add a convolutional layer of $64 \times 1 \times 1$ at the shortcut connection. There is also a convolutional layer of $64 \times 3 \times 3$ before and after each RM.

The Decom-Net takes in paired low/normal-light images (S_{low} and S_{normal}) each time and learns the decomposition for both low-light and its corresponding normal-light image under the guidance that the low-light image and normal-light image share the same reflectance map. During training, there is no need to provide the ground truth of the reflectance and

illumination. Only requisite knowledge including the consistency of reflectance and the smoothness of the illumination map is embedded into the network as loss functions. Note that the illumination map and reflectance map of the normal-light image neither participate in the follow-up training, but only provide a reference for decomposition.

Denoise-Net: Most traditional methods and previous learning-based methods based on the Retinex theory do not take the noise into account after obtaining the decomposition results, it will cause the final enhancement result to be affected by the noise in the reflectance map. Recently, researchers have designed effective models that can suppress noise while enhancing the contrast of low-light images. Inspired by that, we also designed a Denoise-Net to suppress the noise in the reflectance map. Similar to most learning-based methods, our Denoise-Net only uses the spatial information of the image, because eliminating noise by suppressing high-frequency signals in the reflectance map may result in the loss of inherent details.

U-Net [42] has achieved excellent results in a large number of computer vision tasks due to its excellent structural design. In the field of low-light image enhancement, a large number of networks have adopted the U-Net as the main architecture or a part of it. Chen *et al.* [9] directly uses U-Net to enhance the image without any modification to the network and achieve excellent results. Inspired by the residual network, Res-UNet [43] substitutes a module with residual connections for each sub-module of U-Net. However, U-Net and Res-UNet use multiple max-pooling layers in the feature extraction stage, and the max-pooling layer will lead to the loss of feature information, this is what we do not want. Inspired by [44], we replace the max-pooling layers with stride convolutional layers, which will slightly increase the network parameters, but improve the performance. Both U-Net and Res-UNet belong to "shallow-wide" architecture, Li *et al.* [45] demonstrate that "deep-narrow" architecture is more efficient, so we replace each sub-module of UNet with RM to build "deep-narrow" Res-UNet, which is named DN-ResUnet in this paper. The RM used in Denoise-Net is slightly different from that in Decom-Net, the number of convolutions is maintained at 128 without increasing, except for the last 1×1 convolutional layer. Moreover, we use dilated convolution in the first two layers of the network to extract more feature information. As shown in Fig.5, the illumination map obtained by our Denoise-Net retains the original image details while suppressing noise.

Relight-Net: After getting the decomposition results, it is necessary to improve the contrast of the illumination map to obtain a high visual quality result, which is the purpose of Relight-Net design. Inspired by the effectiveness of combing spatial and frequency information to restore high-quality sharp images in other image restoration tasks [46], our Relight-Net consists of two modules: Contrast Enhancement Module (CEM) and Detail Reconstruction Module (DRM). The CEM uses spatial information for contrast enhancement, its architecture is similar to Denoise-Net, we also utilize multi-scale fusion, concatenate the output of each deconvolutional layer in the expansive path to reduce the loss of feature information. The DRM extracts frequency information based on the Fourier

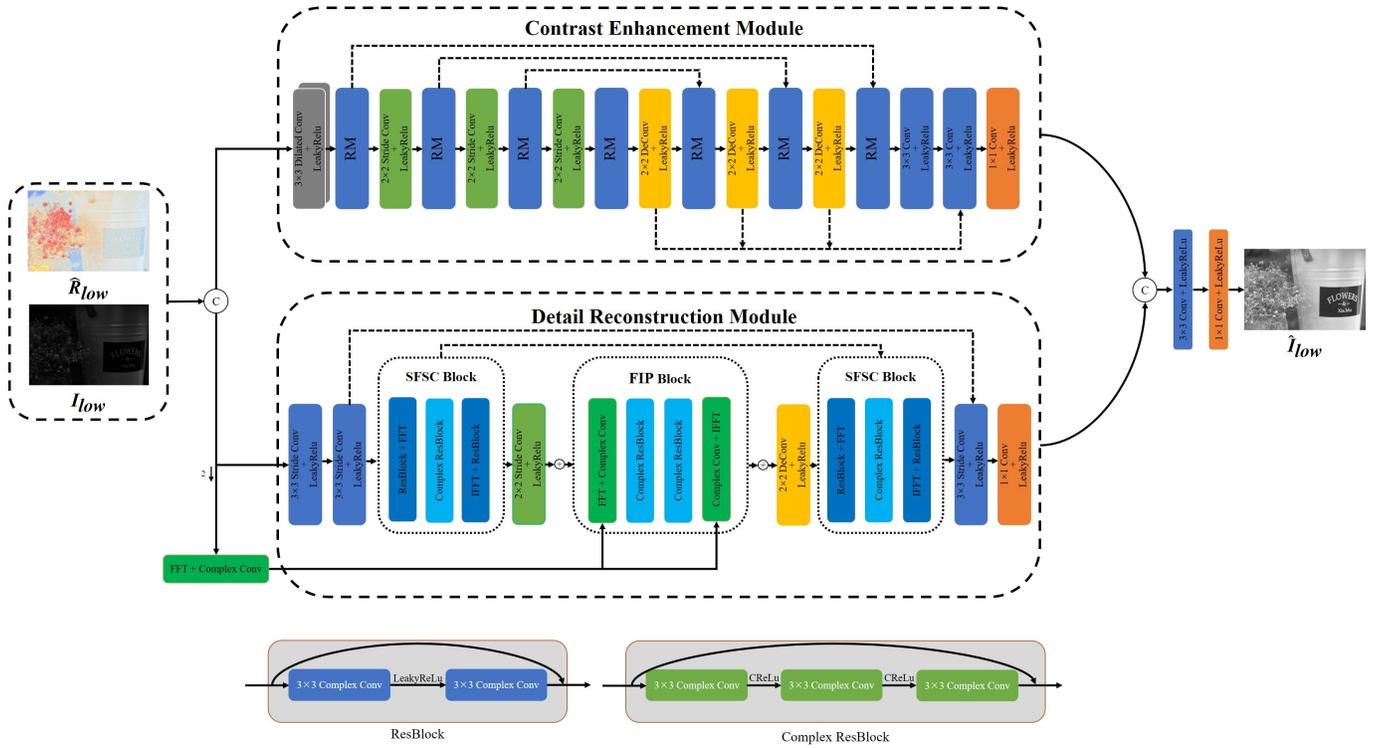


Fig. 4. The proposed Relight-Net architecture. The Relight-Net consists of two modules: Contrast Enhancement Module (CEM) and Detail Reconstruction Module (DRM). CEM uses spatial information for contrast enhancement and DRM uses frequency information to preserve image details.

The decomposition results obtained by our method are illustrated in Fig.5. The reflectance map obtained by Denoise-Net retains the original image details while suppressing the noise, and the Relight-Net properly improves the contrast of the illumination map and retains more details.

B. Loss Function

In the training phase, the three subnets are trained separately, so the whole loss function consists of three parts: the decomposition loss \mathcal{L}_{Dc} the denoise loss \mathcal{L}_{Dn} , and the relight loss \mathcal{L}_{Re} . Each of them consists of two parts: content loss and perceptual loss [48].

Decomposition loss: Our decomposition loss contains two parts: content loss \mathcal{L}_{Dc-con} , and perceptual loss \mathcal{L}_{Dc-per} . We use the L1 loss as content loss,

$$\begin{aligned} \mathcal{L}_{Dc-con} = & \sum_{i=1}^N |R_{low} \circ I_{low} - S_{low}| + \\ & \sum_{i=1}^N |R_{nor} \circ I_{nor} - S_{nor}| + \\ & \lambda_1 \sum_{i=1}^N |R_{nor} \circ I_{low} - S_{low}| + \\ & \lambda_2 \sum_{i=1}^N |R_{low} \circ I_{nor} - S_{nor}| \end{aligned} \quad (2)$$

and we calculate the perceptual loss based on features extracted from a VGG-16 pre-trained model, and in contrast to

previous methods, we adopt features before rather than after the activation layer,

$$\begin{aligned} \mathcal{L}_{Dc-per} = & \frac{1}{C_j H_j W_j} \|\phi_j(R_{low} \circ I_{low}) - \phi_j(S_{low})\|_2^2 + \\ & \frac{1}{C_j H_j W_j} \|\phi_j(R_{nor} \circ I_{nor}) - \phi_j(S_{nor})\|_2^2 \end{aligned} \quad (3)$$

The decomposition loss is formulated as:

$$\mathcal{L}_{Dc} = \mathcal{L}_{Dc-con} + \lambda_3 \mathcal{L}_{Dc-per} \quad (4)$$

Denoise loss: Similar to the decomposition loss, denoise loss contains two parts: content loss \mathcal{L}_{Dn-con} , and perceptual loss \mathcal{L}_{Dn-per} . We also adopt L1 loss as content loss,

$$\mathcal{L}_{Dn-con} = \sum_{i=1}^N |\hat{R}_{low} - R_{nor}| \quad (5)$$

and use features before the activation layer extracted from a VGG-16 pre-trained model to calculate the perceptual loss.

$$\mathcal{L}_{Dn-per} = \frac{1}{C_j H_j W_j} \|\phi_j(\hat{R}_{low}) - \phi_j(R_{nor})\|_2^2 \quad (6)$$

The denoise loss is formulated as:

$$\mathcal{L}_{Dn} = \mathcal{L}_{Dn-con} + \lambda_4 \mathcal{L}_{Dn-per} \quad (7)$$

Relight Loss: Relight loss contains content loss, perceptual loss, and detail preserve loss. We use the same strategy as decomposition loss and denoise loss to build content loss and perceptual loss.

$$\mathcal{L}_{Re-con} = \sum_{i=1}^N |\hat{S}_{low} - S_{nor}| \quad (8)$$

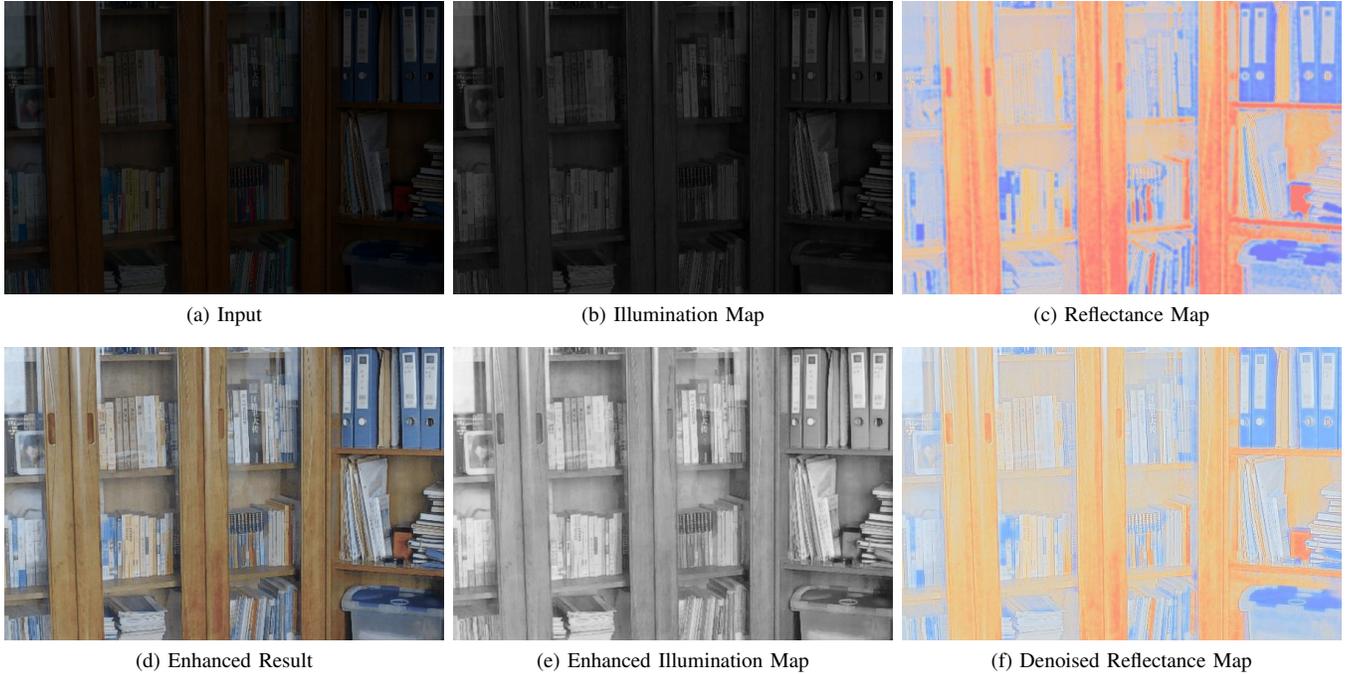


Fig. 5. The decomposition results obtained by our method. The illumination map and reflectance map are obtained by the Decom-Net, the enhanced illumination map is obtained by Relight-Net, and the denoised reflectance map is obtained by Denoise-Net.

$$\mathcal{L}_{Re-per} = \frac{1}{C_j H_j W_j} \|\phi_j(\hat{S}_{low}) - \phi_j(S_{nor})\|_2^2 \quad (9)$$

Moreover, since we used the frequency information in the Relight-Net, we proposed a novel frequency loss to help Relight-Net recover more details. The enhanced image and the sharp image are converted into frequency domain by fast Fourier transform, and the Wasserstein distance is used to minimize the difference between the real part and imaginary part of enhanced result and the ground truth in frequency domain. The frequency loss is formulated as:

$$\mathcal{L}_{Re-fre} = \frac{1}{N^2} \sum_{i=real}^{imag} inf_{\gamma \sim \Pi(\hat{S}_{low}^i, S_{nor}^i)} \mathbb{E}_{(x,y) \sim [|\hat{S}_{low}^i - S_{nor}^i|]} \quad (10)$$

The relight loss is formulated as:

$$\mathcal{L}_{Re} = \mathcal{L}_{Ren-con} + \lambda_5 \mathcal{L}_{Re-per} + \lambda_6 \mathcal{L}_{Re-fre} \quad (11)$$

V. EXPERIMENTS

A. Implementation Details

Our implementation is done with PyTorch. The proposed network can be quickly converged after being trained for 20 epochs on a 1080Ti GPU with our LSRW dataset. We use the Adam [49] optimizer with $lr = 10^{-3}$, $\beta_1 = 0.9$, and $\beta_2 = 0.999$. The batch size and patch size are set to 4 and 96, respectively. We also use the learning rate decay strategy, which reduces the learning rate to 10^{-4} after 10 epochs. λ_1 and λ_2 in Eq.2 are set to 0.01, λ_3 in Eq.4, λ_4 in Eq.7, and λ_5 in Eq.11 are set to 0.1, λ_6 in Eq.11 is set to 0.01. For more implementation details of our network, please refer to the code we are going to release.

B. Comparison with State-of-the-Arts on the Real Datasets

We compare the proposed method with the existing state-of-art methods (MF, Dong, NPE, SRIE, BIMEF, MSRCR, LIME, RetinexNet, DSLR, MBLLN, EnlightenGAN, and Zero-DCE) on six publicly available datasets, including LOL, LIME, DICM [50], NPE [51], MEF [52], and VV¹. For fair comparison, we use the released codes of these methods without any modification, and use our LSRW dataset to train supervised learning-based methods, including Retinexnet, MBLLN, and DSLR. Since Zero-DCE and EnlightenGAN use unpaired data for training, we use their published pre-trained model for comparison. The LOL dataset captured 500 pairs of real low/normal-light images by changing the camera’s exposure time and ISO. This is the only existing real low/normal-light image dataset for low-light image enhancement (the SID dataset is used for extremely low-light image enhancement). The results are shown in Table II. It can be seen that our method outperforms the existing state-of-the-art methods on the LOL dataset for both PSNR and SSIM, the proposed R2RNet achieves the best performance with an average PSNR score of 20.207dB and SSIM score of 0.816, which exceed the second-best approach (MBLLN) by 1.347dB (20.207-18.860) on PSNR and 0.062 (0.816-0.754) on SSIM. The visual comparison is illustrated in Fig.6. It can be seen that some traditional methods (such as SRIE, NPE) will cause under-enhancement results, while other methods based on the Retinex theory (such as LIME, RetinexNet) will blur the details or amplify the noise. The enhancement results generated by our method can not only improve the local and global contrast, have clearer details, but also suppress the

¹<https://sites.google.com/site/vonikakis/datasets>



Fig. 6. Visual Comparison with state-of-the-art low-light image enhancement methods on the LOL dataset. Please zoom in for a better review. EG denotes EnlightenGAN.

noise well, which demonstrates that our method can enhance the image contrast and suppress noise simultaneously. Please zoom in to compare more details.

LIME, DICM, NPE, MEF, VV have commonly been used as benchmark datasets for low-light image enhancement methods evaluation, which only contain low-light images so that PSNR and SSIM cannot be used for quantitative evaluation. Therefore, we use non-reference image quality evaluation NIQE to evaluate the performance of our method. The results are shown in Table III. Some visual comparison is illustrated in Fig.7. Please zoom in to compare more details.

C. User Study

We conduct a user study to compare the performance of our method and other methods. We collect 20 additional low-light images in the real-world scenes for user study and invite 10 participants to evaluate the enhanced results of the real low-light images obtained using five different methods (NPE, LIME, EnlightenGAN, MBLLEN, and Our method). The participants should consider the contrast, artifacts, noise,

details, and color of the enhanced results and rate them according to the performance of the enhanced images (from 1 to 5, 1 means the best, 5 means the worst). Fig.8 shows the distribution of scores and our method gets the best result, which demonstrates that the enhanced images obtained by our method are more visually satisfactory.

D. Ablation Study

In this section, we quantitatively evaluate the effectiveness of different components and the loss function setting in our model based on the LOL dataset. The results are shown in Table IV.

1). Effectiveness of CEM and DRM: We evaluate the effectiveness of the Contrast Enhancement Module (CEM) and Detail Reconstruction Module (DRM) in the Relight-Net by removing CEM and DRM respectively to build our Relight-Net. Removing CEM or DRM will significantly reduce the performance of our model. As shown in Table IV, the experimental result demonstrated that combining spatial and

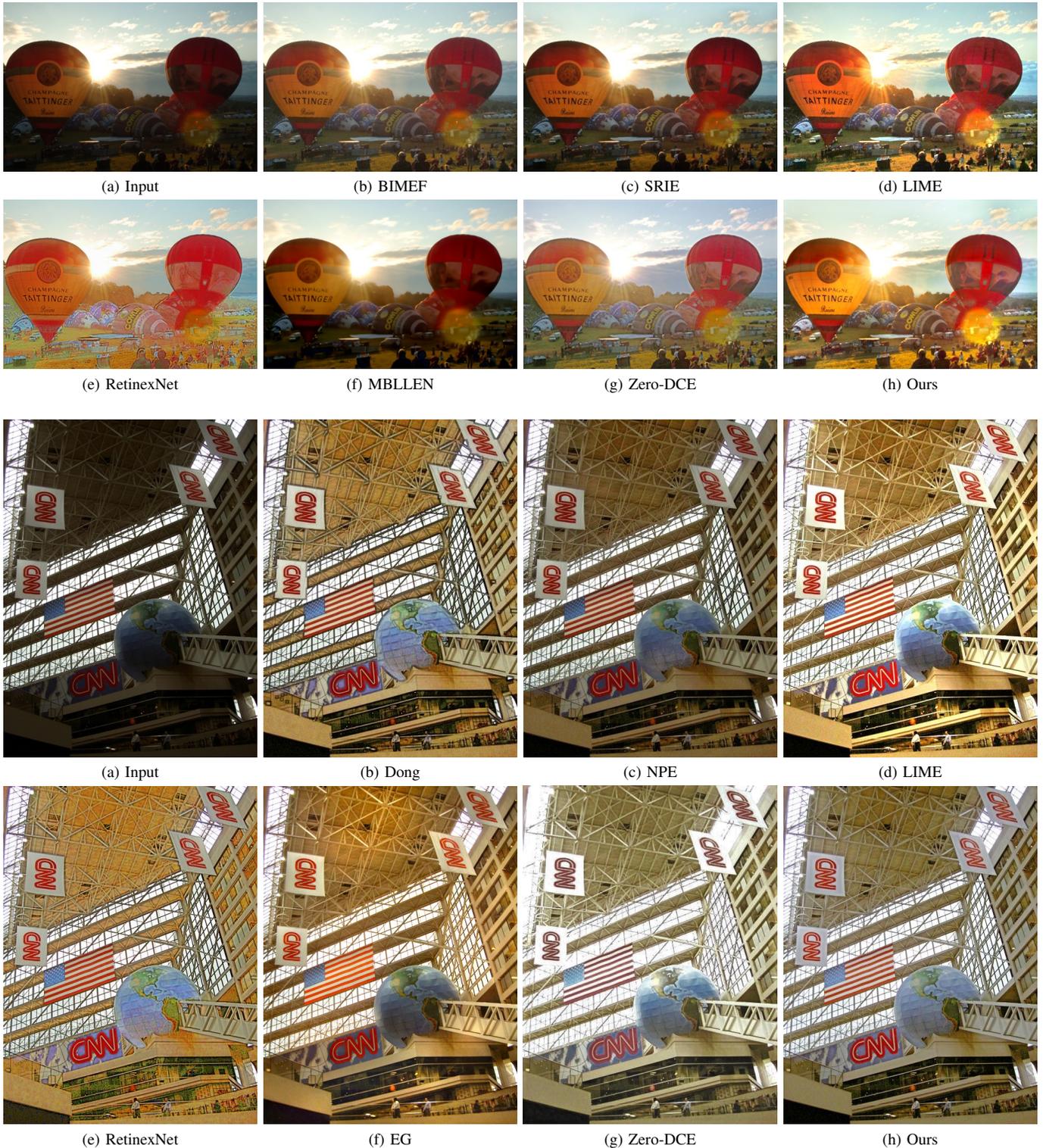


Fig. 7. Visual Comparison with state-of-the-art low-light image enhancement methods on the MEF dataset (row 1-2) and the DICM dataset (row 3-4). Please zoom in for a better review. EG denotes EnlightenGAN.

frequency information can obtain better performance than using one alone.

2). Effectiveness of deep-narrow architecture: We evaluate the effectiveness of DN-ResUnet and compare it with the corresponding "shallow-wide" ResUnet. We replace DN-

ResUnet architecture in DenoiseNet and RelightNet with ResUnet. Our proposed DN-ResUnet exceeds ResUnet with 0.971dB (=20.207-19.236) on PSNR and 0.011 (=0.816-0.805) on SSIM. The results demonstrate that our default architecture will result in better performance.

TABLE II

QUANTITATIVE EVALUATION OF LOW-LIGHT IMAGE ENHANCEMENT METHODS ON THE LOL DATASET. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD. NOTE THAT RETINEXNET, DSLR, AND MBLLEN ARE TRAINED ON OUR LSRW DATASET.

Methods	PSNR	SSIM	FSIM	MAE	GMSD
MF	16.907	0.605	0.906	0.115	0.090
Dong	15.905	0.537	0.875	0.142	0.119
NPE	17.354	0.546	0.875	0.093	0.115
SRIE	11.557	0.531	0.887	0.220	0.103
BIMEF	13.769	0.640	0.907	0.103	0.085
MSRCR	13.964	0.514	0.827	0.046	0.151
LIME	17.267	0.513	0.850	0.097	0.123
RetinexNet	16.013	0.661	0.851	0.071	0.146
DSLR	15.036	0.667	0.883	0.196	0.144
MBLLEN	18.860	0.754	0.904	0.032	0.103
EnlightenGAN	17.239	0.678	0.911	0.087	0.085
Zero-DCE	14.584	0.610	0.911	0.161	0.087
Ours	20.207	0.816	0.933	0.036	0.076

TABLE III

NIQE SCORES ON THE MEF, LIME, NPE, VV, DICM DATASET, RESPECTIVELY. THE BEST RESULTS ARE HIGHLIGHTED IN BOLD.

Methods	MEF	LIME	NPE	VV	DICM	AVG
Dong	4.695	4.052	4.334	3.284	3.263	3.926
NPE	4.256	3.905	3.403	3.031	2.845	3.488
LIME	4.447	4.155	3.796	2.750	3.001	3.630
RetinexNet	4.408	4.361	3.943	3.816	4.209	4.147
MBLLEN	3.654	4.073	5.000	4.294	3.442	4.063
EnlightenGAN	3.573	3.719	4.113	2.581	3.570	3.443
Zero-DCE	4.024	3.912	3.667	3.217	2.835	3.531
Ours	3.029	3.176	3.355	3.093	3.503	3.431

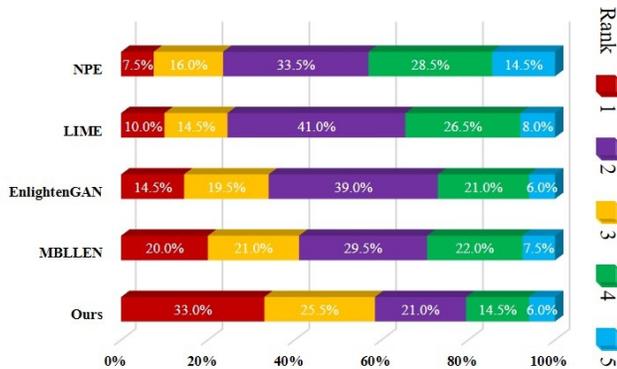


Fig. 8. Score distribution of user study. Our method has a better score distribution.

3). Loss function setting: In order to explore the effectiveness of the loss function setting, we conduct experiments by converting the content loss into MSE loss, removing the perceptual loss and removing the frequency loss, respectively. Using L1 loss exceeds MSE loss with 0.676dB (=20.207-19.531) on PSNR and 0.012 (=0.816-0.804) on SSIM. Removing the perceptual loss and the frequency loss will lead to performance degradation. After removing the perceptual loss,

PSNR decreased by 0.868db (=20.207-19.339) and SSIM decreased by 0.043(=0.816-0.773). After removing the frequency loss, PSNR decreased by 0.451db (=20.207-19.756) and SSIM decreased by 0.012 (=0.816-0.804). The experimental results verify the rationality of our loss function setting.

TABLE IV

ABLATION STUDY. THIS TABLE REPORTS THE PERFORMANCE UNDER EACH CONDITION BASED ON THE LOL DATASET. IN THIS TABLE, "w/o" MEANS WITHOUT.

Conditions	PSNR	SSIM
1. default	20.207	0.816
2. w/o CEM	17.965	0.784
3. w/o DRM	17.483	0.772
4. DN-ResUnet → ResUnet	19.236	0.804
5. L1 → MSE	19.531	0.804
6. w/o perceptual loss	19.339	0.773
7. w/o frequency loss	19.756	0.805

E. Pre-Processing for Improving Face Detection

Image enhancement as pre-processing for improving subsequent high-level vision tasks has recently received increasing attention [53], [54]. We investigate the impact of light enhancement on the DARK FACE dataset², which was specifically built for the task of face detection in low-light conditions. The DARK FACE dataset consists of 6,100 real-world low-light images captured during the nighttime, including 6000 images in the training/validation set, 100 images in the testing set. Because there are no corresponding labels in the test set, we randomly select 100 images from the training set for evaluation, applying our R2RNet as a pre-processing step, followed by two state-of-art pre-trained face detection methods: RetinaFace [55] and DSFD [56]. Using R2RNet as pre-processing improves the average precision (AP) from 17.12% (DSFD+Low-light image) and 15.28% (RetinaFace+Low-light image) to 33.98% (DSFD+R2RNet) and 25.97% (RetinaFace+R2RNet) after enhancement, which demonstrates that R2RNet can improve the performance of high-level vision tasks, in addition to producing visually pleasing results. We also conduct experiments using EnlightenGAN and MBLLEN. EnlightenGAN improves the AP to 32.75% and 23.44%, and MBLLEN improves the AP to 31.69% and 24.67%. Examples of face detection results are illustrated in Fig.9.

VI. CONCLUSION

In this research, we proposed a novel Real-low to Real-normal Network for low-light image enhancement based on the Retinex theory, the proposed network includes three sub-networks: a Decom-Net, a Denoise-Net, and a Relight-Net. The enhanced results obtained by our method have better visual quality. Unlike previous methods, we collected the first large-scale real-world paired low/normal-light images dataset known as LSRW dataset used for network training. The results

²<https://flyywh.github.io/CVPRW2019LowLight>



Fig. 9. Examples of face detection results. We use EnlightenGAN, MBLLEN, and our R2RNet as pre-processing step, followed by DSFD and RetinaFace for detection. EG denotes EnlightenGAN.

on the publicly available datasets showed that our method can properly improve the image contrast and suppress noise, and achieve the highest PSNR and SSIM scores, which outperform state-of-the-art methods by a large margin. We also showed that our R2RNet can effectively improve the performance of face detection methods under low-light conditions.

Overall, the main contributions of this work are threefold:

1). We proposed a novel Real-low to Real-normal Network (R2RNet) to transform the weakly illuminated image into the normal-light image. The proposed network consists of three subnets: a Decom-Net, a Denoise-Net, and a Relight-Net. The purpose of Decom-Net is to decompose the input image into an illumination map and a reflectance map based on the Retinex theory. The Denoise-Net is designed to suppress noise in the reflectance map. The Relight-Net uses spatial information of low-light images to improve contrast and uses frequency information for detail reconstruction. Additionally, a novel frequency loss function was used to help Relight-Net recover more image details.

2). Different from previous methods using synthetic image datasets, we collected the first large-scale real-world paired low/normal-light image dataset (LSRW dataset), which contains 5650 pairs of low/normal-light images to satisfy the training requirement of deep neural networks and make our model has better generalization performance in real-world scenes.

3). The experimental results on the publicly available datasets demonstrated that our method outperforms state-of-the-art methods by a large margin. The enhanced results generated by our method are excellent in contrast, brightness, detail preservation, and noise suppression. And we also showed that our method can effectively improve the performance of face detection under insufficient illumination conditions.

In the future, we will explore a more effective model and apply the model to other enhancement tasks (such as low-light video enhancement, extremely low-light image enhancement,

etc.).

ACKNOWLEDGMENT

We would like to thank Mingrui Wu, Zhiyun Jiang, and Wenxuan Liu for helping us collect the LSRW dataset.

REFERENCES

- [1] C. Guo, C. Y. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, "Zero-reference deep curve estimation for low-light image enhancement," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, pp. 1780–1789.
- [2] Y. Jiang, X. Gong, D. Liu, Y. Cheng, C. Fang, X. Shen, J. Yang, P. Zhou, and Z. Wang, "Enlightengan: Deep light enhancement without paired supervision," *IEEE Trans. Image Process.*, vol. 30, pp. 2340–2349, 2021.
- [3] X. Ke, W. Lin, G. Chen, Q. Chen, X. Qi, and J. Ma, "Edllie-net: Enhanced deep convolutional networks for low-light image enhancement," in *IEEE 5th Int. Conf. Image. Vis. Comput.(ICIVC)*, 2020, pp. 59–68.
- [4] H. Jiang, Y. Hao, F. Zou, F. Lin, and S. Han, "A visual navigation system for uav under diverse illumination conditions," *Appl. Arti. Intell.*, pp. 1–21, 2021.
- [5] X. Ren, M. Li, W.-H. Cheng, and J. Liu, "Joint enhancement and denoising method via sequential decomposition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*. IEEE, 2018, pp. 1–5.
- [6] K. Xu, X. Yang, B. Yin, and R. W. Lau, "Learning to restore low-light images via decomposition-and-enhancement," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2020, p. 2281–2290.
- [7] E. H. Land, "The retinex theory of color vision," *Sci. Amer.*, vol. 237, no. 6, pp. 108–128, 1977.
- [8] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*. Newcastle, U.K., 2018, pp. 1–12.
- [9] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 3291–3300.
- [10] Y. T. Kim, "Contrast enhancement using brightness preserving bi-histogram equalization," *IEEE Trans. Consum. Electron.*, vol. 43, no. 1, pp. 1–8, 1997.
- [11] D. J. Jobson, Z. Rahman, and G. A. Woodell, "A multiscale retinex for bridging the gap between color images and the human observation of scenes," *IEEE Trans. Image Process.*, vol. 6, no. 7, pp. 965–976, 1997.
- [12] X. Fu, D. Zeng, Y. Huang, X.-P. Zhang, and X. Ding, "A weighted variational model for simultaneous reflectance and illumination estimation," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 2782–2790.

- [13] X. Fu, D. Zeng, Y. Huang, Y. L. X. Ding, and J. Paisley, "A fusion-based enhancing method for weakly illuminated images," *Signal Process.*, vol. 129, pp. 82–96, 2016.
- [14] Z. Ying, G. Li, and W. Gao, "A bio-inspired multi-exposure fusion framework for low-light image enhancement," *arXiv preprint arXiv:1711.00591*, 2017. [Online]. Available: <https://arxiv.org/abs/1711.00591>
- [15] L. Mading, J. Liu, W. Yang, X. Sun, and Z. Guo, "Structure-revealing low-light image enhancement via robust retinex model," *IEEE Trans. Image Process.*, vol. 27, no. 6, pp. 2828–2841, 2018.
- [16] X. Guo, Y. Li, and H. Ling, "Lime: Low-light image enhancement via illumination map estimation," *IEEE Trans. Image Process.*, vol. 26, no. 2, pp. 982–993, 2016.
- [17] K. Dabov, A. Foi, V. Katkovnik, and K. Egiazarian, "Image denoising by sparse 3-d transform-domain collaborative filtering," *IEEE Trans. Image Process.*, vol. 16, no. 8, pp. 2080–2095, 2007.
- [18] R. Liu, L. Ma, J. Zhang, X. Fan, and Z. Luo, "Retinex-inspired unrolling with cooperative prior architecture search for low-light image enhancement," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 10 561–10 570.
- [19] K. Zhang, W. Zuo, Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [20] J. Salmon, Z. Harmany, C. A. Deledalle, and R. Willett, "Poisson noise reduction with non-local pca," *J. math. imag. vis.*, vol. 48, no. 2, pp. 279–294, 2014.
- [21] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 38, no. 2, pp. 295–307, 2015.
- [22] M. Haris, G. Shakhnarovich, and N. Ukita, "Deep back-projection networks for super-resolution," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 1664–1673.
- [23] L. Shen, Z. Yue, F. Feng, Q. Chen, S. Liu, and J. Ma, "Msr-net: Low-light image enhancement using deep convolutional network," *arXiv preprint arXiv:1711.02488*, 2017. [Online]. Available: <https://arxiv.org/abs/1711.02488>
- [24] Y. Zhang, J. Zhang, and X. Guo, "Kindling the darkness: A practical low-light image enhancer," in *Proc. 27th. ACM Int. Conf. on Multimedia.*, 2019, pp. 1632–1640.
- [25] S. Lim and W. Kim, "Dslr: Deep stacked laplacian restorer for low-light image enhancement," *IEEE Trans. Multimedia.*, 2020.
- [26] Z. Ying, G. Li, Y. Ren, R. Wang, and W. Wang, "A new low-light image enhancement algorithm using camera response model," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, 2017, pp. 3015–3022.
- [27] K. G. Lore, A. Akintayo, and S. Sarkar, "Llnet: A deep autoencoder approach to natural low-light image enhancement," *Pattern Recognit.*, vol. 61, pp. 650–662, 2017.
- [28] Z. Rahman, Y. Pu, M. Aamir, S. Wali, and Y. Guan, "Efficient image enhancement model for correcting uneven illumination images," *IEEE Access*, vol. 8, pp. 109 038–109 053, 2020.
- [29] B. Goossens, H. Luong, A. Pizurica, and W. Philips, "An improved non-local denoising algorithm," in *2008 Int. Workshop. Loc Non-Loc. Approx. Image. Process. (NLNA)*, 2008, pp. 143–156.
- [30] K. Isogawa, T. Ida, T. Shiodera, and T. Takeguchi, "Deep shrinkage convolutional neural network for adaptive noise reduction," *IEEE Signal Process. Lett.*, vol. 25, no. 2, pp. 224–228, 2017.
- [31] K. Zhang, W. Zuo, and L. Zhang, "Ffdnet: Toward a fast and flexible solution for cnn-based image denoising," *IEEE Trans. Image Process.*, vol. 27, no. 9, pp. 4608–4622, 2018.
- [32] J. Chen, J. Chen, H. Chao, and M. Yang, "Image blind denoising with generative adversarial network based noise modeling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2018, pp. 3155–3164.
- [33] Y. Mei, Y. Fan, Y. Zhang, J. Yu, Y. Zhou, D. Liu, Y. Fu, T. S. Huang, and H. Shi, "Pyramid attention networks for image restoration," *arXiv preprint arXiv:2004.13824*, 2020. [Online]. Available: <https://arxiv.org/abs/2004.13824>
- [34] D. W. Kim, C. J. Ryun, and W. J. S., "Grdn: Grouped residual dense network for real image denoising and gan-based real-world noise modeling," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, 2019, pp. 1–9.
- [35] K. Lin, T. H. Li, S. Liu, and G. Li, "Real photographs denoising with noise domain adaptation and attentive generative adversarial network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW)*, 2019, pp. 1–5.
- [36] K. Zhang, Z. W. Y. Chen, D. Meng, and L. Zhang, "Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising," *IEEE Trans. Image Process.*, vol. 26, no. 7, pp. 3142–3155, 2017.
- [37] F. Lv, F. Lu, J. Wu, and C. Lim, "Mblen: Low-light image/video enhancement using cnns," in *Proc. Brit. Mach. Vis. Conf. (BMVC)*, 2018, pp. 1–13.
- [38] F. Lv and F. Lu, "Attention-guided low-light image enhancement," *arXiv preprint arXiv:1908.00682*, 2019. [Online]. Available: <https://arxiv.org/abs/1908.00682>
- [39] L. W. Wang, Z. S. Liu, W. C. Siu, and D. P. Lun, "Lightening network for low-light image enhancement," *IEEE Trans. Image Process.*, vol. 29, pp. 7984–7996, 2020.
- [40] Y. Wang, Y. Cao, Z. J. Zha, J. Zhang, Z. Xiong, W. Zhang, and F. Wu, "Progressive retinex: Mutually reinforced illumination-noise perception network for low-light image enhancement," in *Proc. 27th. ACM Int. Conf. on Multimedia.*, 2019, pp. 2015–2023.
- [41] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2016, pp. 770–778.
- [42] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Image Comput. Assist. Intervent.* Springer, 2015, pp. 234–241.
- [43] X. Xiao, S. Lian, Z. Luo, and S. Li, "Weighted res-unet for high-quality retina vessel segmentation," in *Proc. 9th Int. Conf. Inf. Tech. Med. Edu.(ITME)*, 2018, pp. 327–331.
- [44] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for simplicity: The all convolutional net," *arXiv preprint arXiv:1412.6806*, 2014. [Online]. Available: <https://arxiv.org/abs/1412.6806>
- [45] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z. Jiang, E. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," *arXiv preprint arXiv:2101.11986*, 2021. [Online]. Available: <https://arxiv.org/abs/2101.11986>
- [46] W. B. Z. M. J. Y. Zhang, C. Liang, Z. Lu, and Y. Wu, "Sdwnet: A straight dilated network with wavelet transformation for image deblurring," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2021, pp. 1895–1904.
- [47] T. Chihab, B. Olexa, Z. Ying, S. Dmitriy, S. Sandeep, J. Santos, M. Soroush, R. Negar, Y. Bengio, and C. Pal, "Deep complex networks," in *Int. Conf. on Learn. Repr.*, 2018.
- [48] J. Johnson, A. Alahi, and F. F. Li, "Perceptual losses for real-time style transfer and super-resolution," in *Proc. Eur. Conf. Comput. Vis. (ECCV)*. Springer, 2016, pp. 694–711.
- [49] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014. [Online]. Available: <https://arxiv.org/abs/1412.6980>
- [50] C. Lee, C. Lee, and C. S. Kim, "Contrast enhancement based on layered difference representation," in *Proc. 19th. Int. Conf. Imag. Process.*, 2012, pp. 965–968.
- [51] S. Wang, J. Zheng, H. M. Hu, and B. Li, "Naturalness preserved enhancement algorithm for non-uniform illumination images," *IEEE Trans. Image Process.*, vol. 22, no. 9, pp. 3538–3548, 2013.
- [52] K. Ma, K. Zeng, and Z. Wang, "Perceptual quality assessment for multi-exposure image fusion," *IEEE Trans. Image Process.*, vol. 24, no. 11, pp. 3345–3356, 2015.
- [53] B. Li, X. Peng, Z. Wang, J. Xu, and D. Feng, "Aod-net: All-in-one dehazing network," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2017, pp. 4770–4778.
- [54] Y. Liu, G. Zhao, B. Gong, Y. Li, R. Raj, N. Goel, S. Kesav, S. Gottimukkala, Z. Wang, and W. Ren, "Improved techniques for learning to dehaze and beyond: A collective study," *arXiv preprint arXiv:1807.00202*, 2018. [Online]. Available: <https://arxiv.org/abs/1807.00202>
- [55] J. Deng, J. Guo, Y. Zhou, J. Yu, I. Kotsia, and S. Zafeiriou, "Retinaface: Single-stage dense face localisation in the wild," *arXiv preprint arXiv:1905.00641*, 2019. [Online]. Available: <https://arxiv.org/abs/1905.00641>
- [56] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang, "Dsfid: dual shot face detector," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR)*, 2019, pp. 5060–5069.