

# Unsupervised Low-light Image Enhancement with Decoupled Networks

Wei Xiong, Ding Liu, Xiaohui Shen, Chen Fang, and Jiebo Luo

**Abstract**—In this paper, we tackle the problem of enhancing real-world low-light images with significant noise in an unsupervised fashion. Conventional unsupervised learning-based approaches usually tackle the low-light image enhancement problem using an image-to-image translation model. They focus primarily on illumination or contrast enhancement but fail to suppress the noise that ubiquitously exists in images taken under real-world low-light conditions. To address this issue, we explicitly decouple this task into two sub-tasks: illumination enhancement and noise suppression. We propose to learn a two-stage GAN-based framework to enhance the real-world low-light images in a fully unsupervised fashion. To facilitate the unsupervised training of our model, we construct samples with pseudo labels. Furthermore, we propose an adaptive content loss to suppress real image noise in different regions based on illumination intensity. In addition to conventional benchmark datasets, a new unpaired low-light image enhancement dataset is built and used to thoroughly evaluate the performance of our model. Extensive experiments show that our proposed method outperforms the state-of-the-art unsupervised image enhancement methods in terms of both illumination enhancement and noise reduction.

**Index Terms**—Low-light Image Enhancement, Generative Adversarial Networks, Image Denoising.

## I. INTRODUCTION

Real-world low-light image enhancement [1]–[4] is challenging since images captured under low-light conditions usually exhibit low illumination and contain heavy noise. Enhancing these images requires adjusting illumination, contrast, color, as well as suppressing the noise while preserving the details simultaneously. Traditional methods for this task primarily focus on adjusting contrast via a fixed tone-mapping [5], resulting in limited performance on challenging cases. Recently, learning-based methods have been utilized to learn content-aware illumination enhancement from data with deep neural networks [6]–[11]. Despite promising performance, many of them heavily rely on low-light and normal-light image pairs, which are expensive or even impossible to obtain in real-world scenarios. One alternative way to cheaply generate such training pairs is to synthesize a low-light image from its counterpart captured under a normal light condition [12]. However, due to the significant signal distribution gap between synthesized dark images and ones taken under real-world low-light conditions, models trained on synthesized image pairs usually fail to generalize well in realistic scenarios [13], [14].

Wei Xiong and Jiebo Luo are with University of Rochester. (Email: wei.xiong@rochester.edu; jluo@cs.rochester.edu)

Ding Liu and Xiaohui Shen are with ByteDance Inc. (Email: liuding@bytedance.com; shenxiaohui@bytedance.com)

Chen Fang is with Tencent Inc. (Email: fangchen1988@gmail.com)

Recently, several unsupervised deep learning-based methods have been developed for image enhancement to eliminate the reliance on paired data. As a general-purpose method, unsupervised image-to-image translation or transformation models [15]–[18] such as CycleGAN [17], and UNIT [16] can be applied to image enhancement. These methods adopt generative adversarial networks (GANs) [19] to encourage the distribution of the generated images to be close to that of the target images without paired supervision. Recently, the GAN-based models have been specially designed to address the task of illumination enhancement [9], [11], [12], [20], [21]. These unsupervised learning approaches can generate images with better illumination and color in several cases. However, they have common limitations in the real-world low-light image enhancement task in two aspects (as we will show in our experiments later): 1) the contrast and illumination of enhanced images can be unsatisfactory and usually suffer from color distortion and inconsistency. The bright regions of a dark image can be over-exposed. Moreover, continuous regions may exhibit a sharp color or brightness inconsistency due to the unstable training of the unsupervised models. 2) when applied to low-light images with heavy noise, models with a single image-to-image mapping network primarily address illumination enhancement but usually fail on noise suppression. *This suggests that a single network may be inadequate to model both the complex illumination patterns and real-world noise patterns.* Although a few prior works have been developed to remove noise after the illumination enhancement as a post-processing step, they either use traditional methods such as BM3D [22] that requires a given noise level as an input, or use learning-based denoising methods that are originally designed for synthesized noise such as additive white Gaussian noise [23]–[25].

To address these issues, we propose to explicitly decouple the whole unsupervised enhancement task into two sub-tasks: 1) illumination enhancement and 2) noise suppression. We propose a two-stage framework to handle each sub-task separately. Specifically, in Stage I, a Retinex-based deep network is trained under a GAN framework in an unsupervised manner to enhance the illumination of low-light images while preserving the contextual details. Instead of predicting the final enhanced image directly, we predict the illumination map of the target image first then use the predicted illumination to generate the enhanced image. It has been suggested by recent low-light image enhancement works [7] that the illumination maps for natural images usually have simple forms. Therefore, the Retinex-based illumination modeling can benefit the generalization ability of the generator. To tackle the inconsistency

problem in color and brightness, we introduce a pyramid module to enlarge the receptive field of the generator.

In Stage II, we propose a guided unsupervised denoising model based on GANs. Our model is adaptively guided by the illumination conditions of the original low-light image and the enhanced image from Stage I. Inspired by the success of pseudo labeling methods in semi-supervised learning [26], [27], in this work, we propose to construct *pseudo triples*, i.e., a pseudo low-light image, a pseudo image after illumination enhancement, and a real noise-free normal-light image with the same content at each training iteration, to facilitate the unpaired training for image denoising. We find this technique is crucial to the success of unsupervised noise modeling. We design an adaptive content loss for the *pseudo triples* to preserve the illumination and color of the input image.

It is noteworthy that a recent work, GCBD [28], also adopts a GAN-based denoising framework that can be learned without paired data. However, our Stage II model significantly differs from this work in the following aspects. First, GCBD limits noise estimation in smooth regions only before learning the GAN-based denoiser, while we jointly perform noise modeling and noise removal in a single GAN framework. Second, GCBD does not consider the influence of illumination. In contrast, our denoising model is explicitly guided by the illumination conditions of the input images to handle illumination-correlated noise adaptively. Third, we propose an adaptive content loss using *pseudo triples* to remove real image noise based on the illumination intensity, while GCBD may produce severe color distorted results without such a constraint, as demonstrated by our experiments later.

We evaluate the performance of our proposed approach over the LOw-Light (LOL) dataset [29] and an unpaired enhancement dataset from [21]. To further demonstrate the effectiveness of our method, we contribute an Unpaired Real-world Low-light image enhancement dataset (URL) for evaluation. Our dataset is composed of 1) low-light images captured under real-world low-light conditions with varying levels of noise, and 2) normal-light images collected from existing data galleries, which consist of diverse scenes ranging from outdoor scenes to indoor pictures. We compare our method with the state-of-the-art unsupervised learning-based enhancement methods on these datasets. Extensive experiments show that our method outperforms other methods in terms of both illumination enhancement and noise suppression.

Our primary contributions are:

- We propose a decoupled framework for unsupervised low-light image enhancement;
- We propose an illumination guided unsupervised denoising model. To facilitate unsupervised training of our denoising model, we construct *pseudo triples* and propose an adaptive content loss to denoise regions guided by both the original lighting condition and enhanced illumination;
- We build an unpaired low-light image enhancement dataset containing varying noise and good diversity as an important complement to the existing low-light enhancement datasets.

## II. RELATED WORK

### A. Image Enhancement.

Traditional image enhancement methods are primarily built upon histogram equalization (HE) [30] or Retinex theory [31], [32]. HE-based methods aim to adjust the histogram of pixel intensities to obtain an image with better contrast [30]. Retinex-based methods assume that an image is the composition of illumination map and reflectance, and thus low-light images can be restored by estimating the illumination map and reflectance map [32]. Recently, learning-based methods have been proposed to learn the illumination enhancement from data [33]–[35]. Later deep neural networks have been used and achieved promising results [6], [12], [29], [36]. A more recent work [10] proposes the DeepLPF model to enhance images by learning different types of filters. [8] *et al.* propose a frequency-based decomposition-and-enhancement model for low-light image enhancement. However, the learning of these models heavily rely on image and reference pairs.

Due to the difficulty of acquiring paired data in real-world scenarios, several weakly supervised and unsupervised enhancement approaches have been proposed, such as WESPE [37], Deep Photo Enhancer [20], and EnlightenGAN [21]. Ignatov *et al.* propose a transitive GAN-based enhancement model that can be learned without paired data [37]. Chen *et al.* propose an unpaired model based on 2-way GANs to enhance images [20]. Jiang *et al.* further propose EnlightenGAN which is specifically designed for illumination enhancement of low-light images [21]. There are also general-purpose unsupervised image translation models that can be used for image enhancement, such as CycleGAN [17], UNIT [16] and ADN [38]. A primary limitation of existing enhancement approaches is that they mainly perform illumination enhancement and fail to address noise suppression, such as CycleGAN [17] and EnlightenGAN [21]. In contrast, we address both illumination enhancement and noise suppression simultaneously.

It is worth noting that a recent work named Zero-DCE [9] also tackles the image enhancement problem in an unsupervised fashion. They learn a curve estimation model with deep networks and use it for pixel-wise dynamic range adjustment of the input image. However, their work primarily focuses on contrast enhancement without paying much attention to noise removal.

### B. Real-world Image Denoising.

There have been a number of works for image denoising, including conventional methods such as BM3D [22] and Non-local means [39], and deep learning-based models such as DnCNN [40], Residual Dense Networks [41] and Non-local Recurrent Networks [42]. However, most of the models are limited to synthetic noise removal and are difficult to generalize to real-world noise removal. Recently, real-world blind denoising models have been proposed to learn a blind denoiser from real-world paired data [13], [43]–[47]. Xu *et al.* [43] design a multi-channel weighted nuclear norm minimization model to use channel redundancy. Guo *et al.* propose CBDNet [13] to directly learn a blind denoiser from real-world paired

data. Kim et al. leverage a GAN based deep network for real-world noise modeling [44]. Other approaches [45]–[47] also show promising results.

Most learning-based denoising models need to be trained with paired data, which is expensive to obtain for real-world noise removal tasks. Recently, several unsupervised denoising methods have been devised, including self-supervised learning approaches, such as Noise2Noise [23] and Noise2Void [24], as well as unpaired training approaches [28], [48]. Our Stage II model is also an unsupervised denoising approach. Unlike previous methods, we capture the real-world noise pattern with the explicit guidance of image illumination conditions and denoise the images with *Pseudo Labeling* technique and an adaptive content loss, which prove to be crucial to the success of real-world image denoising.

### III. OUR APPROACH

As shown in Fig. 1 (a), our approach for real-world low-light image enhancement consists of two stages. In Stage I, we perform illumination enhancement on the real-world low-light images while preserving contextual details. In Stage II, we propose an unsupervised learning-based denoising model to suppress the noise in the output image from Stage I and enhance the contextual details.

#### A. Stage I: Illumination Enhancement

Given a noisy low-light image  $I_l$ , our goal in this stage is to learn a model  $G_e$  to generate an enhanced image  $I_e$  with proper illumination, color, as well as realistic content details. Conventional learning-based models such as CycleGAN and EnlightenGAN usually adopt a U-net [49] like architecture to predict the enhanced image from the low-light input image directly. However, under the unsupervised learning scheme, directly applying such an architecture is easy to produce results with unstable illumination, such as color distortion or inconsistency [21]. Recent work [7] on low-light image enhancement suggests that illumination maps for natural images usually have relatively simple forms, thus using Retinex-based illumination modeling can facilitate the learning process, leading to an enhancer with better generalization ability. Inspired by this, we adopt a Retinex-based model to enhance the low-light images. Based on the Retinex theory, an image  $I$  can be modeled as  $I = S \circ R$ , where  $S$  is the illumination,  $\circ$  denotes element-wise multiplication, and  $R$  is the reflectance. Similar to [7], [29], we regard the reflectance as a well-exposed image  $I_e$ , then we have  $I_l = S \circ I_e$ . In reverse, the enhanced image  $I_e$  can be recovered from the low-light image  $I_l$  given the predicted illumination map  $S$ . As shown in Fig. 1, we use the generator  $G_e$  to estimate an illumination map  $S = G_e(I_l)$  from low-light image  $I_l$ . Then we obtain the enhanced image

$$I_e = I_l / S, \quad (1)$$

where  $/$  is element-wise division.

**Model Architecture.** As shown in Fig. 1, the input low-light image  $I_l$  is passed to the encoder, then a pyramid module, and then decoded by the decoder into an illumination map  $S$  with

RGB channels. The pyramid module, inspired by PSPNet [50], is customized to enlarge the receptive field of our network. In this module, we down-sample the feature maps into features at multiple resolutions, namely,  $1 \times 1$ ,  $2 \times 2$ ,  $4 \times 4$ ,  $16 \times 16$ . At each resolution, feature maps are followed by a convolution layer and a ReLU [51] layer. Then the transformed feature maps are upsampled and concatenated, then fused into the next convolution layer. The fused feature maps are then decoded by the decoder to generate the illumination map  $S$ . With the pyramid module, the receptive field of the model is enlarged, and the network can perceive the illumination information at different spatial levels, which is beneficial for reducing the color and contrast inconsistency.

**Loss Functions.** In Stage I, the generator  $G_e$  is trained with unpaired images. To achieve this goal, we adopt adversarial learning to encourage the distribution of the enhanced image  $I_e$  to be close to that of the normal-light images. Specifically, we use two discriminators to distinguish the generated image from the real normal-light image. The global discriminator  $D_g$  takes the whole image as the input, and outputs the realness of the image. The local discriminator  $D_l$  takes random patches extracted from the image and outputs the realness of each patch. The global discriminator encourages the global appearance of the enhanced image to be similar to a normal-light image, while the local discriminator ensures that the local context (shadow, local contrast, highlight, *etc.*) can be as realistic as the real normal-light images.

We adopt the LSGAN version of relativistic average GAN loss [52] for training the global discriminator. When updating the discriminator, we have:

$$L_D^g = \mathbb{E}_{x_r \in \mathbb{P}} [(D_g(x_r) - \mathbb{E}_{x_f \in \mathbb{Q}} D_g(x_f) - 1)^2] + \mathbb{E}_{x_f \in \mathbb{Q}} [(D_g(x_f) - \mathbb{E}_{x_r \in \mathbb{P}} D_g(x_r))^2] \quad (2)$$

When updating the generator, we have:

$$L_G^g = \mathbb{E}_{x_r \in \mathbb{P}} [(D_g(x_r) - \mathbb{E}_{x_f \in \mathbb{Q}} D_g(x_f))^2] + \mathbb{E}_{x_f \in \mathbb{Q}} [(D_g(x_f) - \mathbb{E}_{x_r \in \mathbb{P}} D_g(x_r) - 1)^2] \quad (3)$$

where  $\mathbb{P}$  and  $\mathbb{Q}$  are the real image (the normal-light images) distribution and the generated image distribution, respectively.  $x_r$  and  $x_f$  are samples from distribution  $\mathbb{P}$  and  $\mathbb{Q}$ , respectively.

We adopt the original LSGAN loss for training the local discriminator. When updating the discriminator, we have:

$$L_D^l = \mathbb{E}_{x_r \in \mathbb{P}} [(D_l(x_r) - 1)^2] + \mathbb{E}_{x_f \in \mathbb{Q}} [(D_l(x_f))^2] \quad (4)$$

When training the generator, we have:

$$L_G^l = \mathbb{E}_{x_f \in \mathbb{Q}} [(D_l(x_f) - 1)^2] \quad (5)$$

Similar to [21], we use a perceptual loss (computed on VGG features) between the output image and the input image to preserve content details from the input image. To prevent the output image from being as dark as the input image, we alleviate the influence of image brightness and force the network to focus on only the content preservation by using instance normalization on the VGG feature maps before performing the perceptual loss. Similar to the adversarial loss, we calculate perceptual loss on both the whole image and the

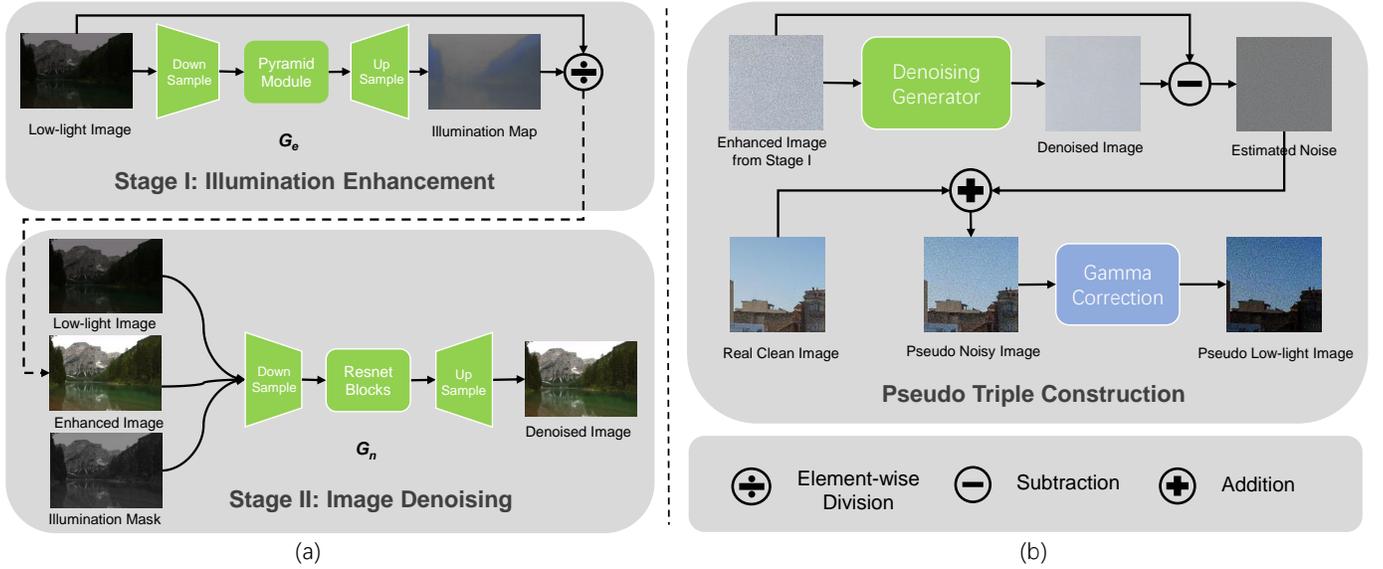


Fig. 1. (a): An overview of our proposed two-stage decoupled networks. In Stage I, given a low-light image, we learn a deep network to predict an illumination map, then use the Retinex theory to obtain the illumination-enhanced image. Then in Stage II, we input the original low-light image, the enhanced image from Stage I and the illumination mask to generate an image with reduced noise and better details. The illumination mask indicates how much the illumination is improved in Stage I. (b): An illustration of pseudo triple construction. We first estimate a noise image from our noisy training data. Then we randomly fetch a clean normal light image and add the estimated noise to the clean image to obtain a pseudo noisy image. Next, we perform Gamma correction to obtain a pseudo low-light counterpart of the noisy image. In this way, we obtain a pseudo triplet {pseudo low-light image, pseudo noisy image, real clean image} as additional training data to train our denoising model.

image patches. We formulate the global perceptual loss  $L_P^g$  and the local version  $L_P^l$  as:

$$L_P^g = \|\Phi(I_l) - \Phi(I_e)\|_2^2 / N, \quad (6)$$

$$L_P^l = \|\Phi(I_l^p) - \Phi(I_e^p)\|_2^2 / N, \quad (7)$$

where  $\Phi$  denotes to the VGG19 feature extractor,  $N$  is the number of elements in the image,  $I_l^p$  and  $I_e^p$  are random patches extracted from  $I_l$  and  $I_e$ , respectively.

By minimizing both the adversarial losses and the perceptual losses, we are able to learn a good illumination predictor and produce results without color distortion and with better contextual details.

**Differences between Prior Work.** The primary differences between our Stage I model and previous unsupervised low-light image enhancement method EnlightenGAN [21] lie in the theory we follow and the model architecture. First, EnlightenGAN uses an image-to-image translation model to map the input low-light image to the illumination-enhanced image directly. In contrast, our model is derived from the Retinex theory, which predicts the illumination map instead of the enhanced image. Such a model proves to be effective and has better generalization ability for image contrast enhancement. Second, we propose to use a pyramid module to connect the encoder and the decoder in our generator. Such a design enlarges the receptive field of our network and helps the model better perceive both the global and local illumination conditions.

### B. Stage II: Noise Suppression

As shown in Fig. 1, our noise suppression model  $G_n$  adopts the original low-light image  $I_l$ , the enhanced image (with

noise)  $I_e$  from Stage I, and an illumination mask  $M$  as inputs to generate the final clean image  $I_c$  with reduced noise as well as a good illumination condition, i.e.,  $I_c = G_n(I_l, I_e, M)$ .  $M$  serves as an indicator showing how much illumination is increased from low-light image  $I_l$  to enhanced image  $I_e$ . We have  $M = \max(\text{illu}(I_e) - \text{illu}(I_l), 0)$ , where  $\text{illu}(\cdot)$  means extracting the illumination of an image. In our work, we directly use the gray-scale version of an image as its illumination map. With  $M$ ,  $I_l$ , and  $I_e$  as inputs, our denoising model is explicitly guided by the illumination conditions.

**Model Architecture.** Our denoising generator  $G_n$  in Stage II adopts an encoder-decoder architecture, with several convolutional blocks followed by several Resnet blocks then decoded back to an image. We adopt a multi-scale discriminator [16] to predict the realness of images at multiple resolutions. The detailed architecture of the networks can be found in the appendix section.

**Loss Functions.** Since there is no ground-truth for the input low-light image during training, to learn a noise-free image from its noisy counterpart, we adopt an LSGAN-based adversarial loss to encourage the generated image to be as clean as the real-world clean normal-light images. Note that the discriminator needs not only to judge whether the illumination and color of the generated image are realistic enough but also needs to judge whether the generated image is clean without much noise. Training a single discriminator with clean images or synthesized images to do both tasks simultaneously is challenging. As our goal in this stage is noise suppression, when feeding the discriminator, we first perform an instance normalization on both the synthesized image and the normal-light clean image to reduce the influence of image illumination, color, and contrast.

During training, we randomly match the denoised image  $I_c$  to a normal-light clean image  $J_c$ . Our adversarial loss for training the discriminator  $D_n$  is:

$$L_D = \mathbb{E}_{J_c \in \mathbb{P}}[(D_n(J_c) - 1)^2] + \mathbb{E}_{I_c \in \mathbb{Q}}[(D_n(I_c))^2]. \quad (8)$$

The corresponding loss for updating the generator is  $L_G = \mathbb{E}_{I_c \in \mathbb{Q}}[(D_n(I_c) - 1)^2]$ , where  $\mathbb{P}$  and  $\mathbb{Q}$  are the normal-light clean image distribution and the distribution of generated images in Stage II, respectively.

Merely using the adversarial loss can cause color shifting problem, i.e., the color of the generated images can be easily distorted, since we only constrain the images after instance normalization to be similar to the normal-light images. As we have already obtained an image  $I_e$  with a satisfactory contrast and color from Stage I, in Stage II, we only need to preserve the contrast and color of  $I_e$ . Therefore, we use a color loss to constrain the generated image  $I_c$  to have the same color as  $I_e$ .

Specifically, we first down-sample the images with average pooling to  $I_c^\downarrow$  and  $I_e^\downarrow$  to suppress the noise in  $I_e$ , then perform the color matching. We have

$$L_{color} = \sum_p \angle((I_c^\downarrow)_p, (I_e^\downarrow)_p) / N_n,$$

where  $p$  is the location of a pixel in the down-sampled image,  $\angle(x, y)$  calculates the inner product between two 3-D vectors which are composed of RGB channels of a pixel location,  $N_n$  is the number of pixels in the downsized image.

**Pseudo Labeling: Constructing Pseudo Triples for Unsupervised Learning.** We propose a *Pseudo Labeling* technique to facilitate the unsupervised training of the denoising model. As shown in Fig. 1 (b), we first estimate the noise in image  $I_e$  as  $I_n = I_e - I_c$ . Then given the randomly matched normal-light clean real image  $J_c$ , we can simulate a pseudo noisy image  $J_e$  by adding the estimated noise to the clean image, i.e.,  $J_e = J_c + I_n$ . We also use Gamma Correction [53] to decrease the brightness of  $J_e$ , to obtain a corresponding pseudo low-light image  $J_l = (J_e)^\lambda$ , where  $\lambda$  is estimated as  $\lambda = \log \bar{I}_c / \log \bar{I}_l$ .  $\bar{I}_c$  and  $\bar{I}_l$  are the average pixel values over all pixel locations of image  $I_c$  and  $I_l$ , respectively. After these steps, we obtain a *pseudo triple*  $\mathcal{J} = \{J_l, J_e, J_c\}$ , where  $J_l$  is the constructed pseudo low-light image,  $J_e$  is the pseudo enhanced image (with noise) and  $J_c$  is the real clean image from our training set. Similarly, we construct the illumination mask for the pseudo triple as  $M_J = \max(\text{illu}(J_e) - \text{illu}(J_l), 0)$ . We can then predict a denoised image  $J_g = G_n(J_l, J_e, M_J)$  from the constructed fake images, and use  $J_c$  as the supervision to train  $G_n$ .

**Adaptive Content Loss with Pseudo Triples.** To train  $G_n$ , we adopt an adaptive content loss to constrain the generated image  $J_g$  to be perceptually close to the clean normal-light image  $J_c$ . This is achieved by using both perceptual loss and L1 reconstruction loss on the pixel space between  $J_g$  and  $J_c$ . As different regions may have different lighting conditions, regions with a significant brightness increase after the first stage may contain heavy noise, and regions without a large brightness increase may contain less noise. When imposing the

reconstruction constraint, we encourage the network to focus more on dark regions where noise is usually heavier. We then formulate the adaptive content loss for the pseudo triples as:

$$L_{con}^{adapt} = \sum_l \|M_J^{(l)} \circ (\Phi_l(J_g) - \Phi_l(J_c))\|_2^2 / N_l + \gamma_p \|M_J \circ (J_g - J_c)\|_1 / N, \quad (9)$$

where  $M_J^{(l)}$  is the downsized version of  $M_J$  that matches the spatial size of the VGG features at the  $l$ -th VGG layer.  $M_J$  serves as the weight mask for each pixel in the image.  $N$  is the number of elements in image  $J_g$ ,  $N_l$  is the number of elements in the feature maps of the  $l$ -th layer in VGG network.  $\Phi_l(I)$  is the feature in the  $l$ -th layer of VGG given the input image  $I$ .  $\gamma_p$  is the weight to balance the losses from the RGB image domain and the VGG feature domain. In our work, we choose the layers of “relu1\_2”, “relu2\_2”, “relu3\_2”, “relu4\_4”, “relu5\_4” to perform both low-level and high-level feature matching, and  $\gamma_p$  is set as 10. We do not use instance normalization on the VGG feature maps of the *pseudo triple*, since we need to preserve the color and contrast.

**Interpretation of Pseudo Labeling.** Note that at the early training stage, the estimated noise may contain clear structures of the objects in the input noisy image  $I_e$ . As a result, the constructed pseudo image  $J_e$  may also contain object structures from  $I_e$ . Since our network is trained to remove noise, it will be difficult for the network to remove the high-frequency object structure patterns. Then the generated image  $J_g$  may also contain object structures from  $I_e$ . By minimizing our proposed adaptive content loss, we are essentially encouraging the estimated noise to contain fewer object structures, i.e., encouraging the network to estimate the noise pattern more accurately.

**Content Preserving Loss.** To make sure that the denoised image  $I_c$  preserves contextual details of the enhanced image  $I_e$ , we also impose a perceptual loss and a reconstruction loss between images  $I_c$  and  $I_e$ . To reduce the influence of color and contrast, we perform instance normalization to the images before imposing the perceptual loss and reconstruction loss. The content loss is formulated as:

$$L_{con} = \sum_l \|\Phi_l(I_e) - \Phi_l(I_c)\|_2^2 / N_l + \gamma_c \|I_e - I_c\|_1 / N, \quad (10)$$

where the layers and operations used are the same as  $L_{con}^{adapt}$ .  $\gamma_c$  is a weight balance term similar to  $\gamma_p$  and set as 10 in our work. The total loss  $L$  for training  $G_n$  is a combination of all the losses.

$$L = L_G + \lambda_c L_{color} + \lambda_C^p L_{con}^{adapt} + \lambda_C^r L_{con}, \quad (11)$$

and we empirically find that setting  $\lambda_c, \lambda_C^p, \lambda_C^r$  as 10, 1, 1, respectively, yields the best result.

## IV. EXPERIMENTS

In this section, we first compare the performance of each model with respect to only illumination enhancement on an unpaired enhancement dataset from EnlightenGAN [21]

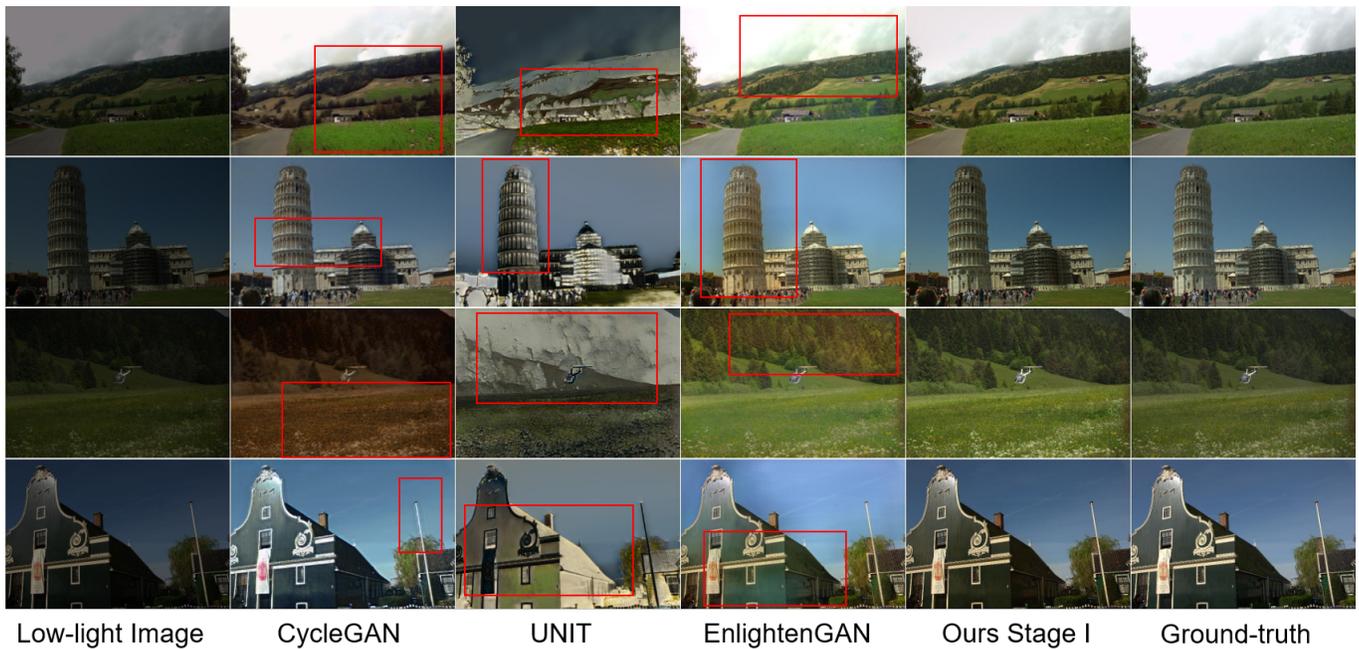


Fig. 2. Illumination enhancement results on the Unpaired Enhancement Dataset from [21]. Please pay attention to the regions in the red boxes, where severe color distortion and contrast inconsistency occur.

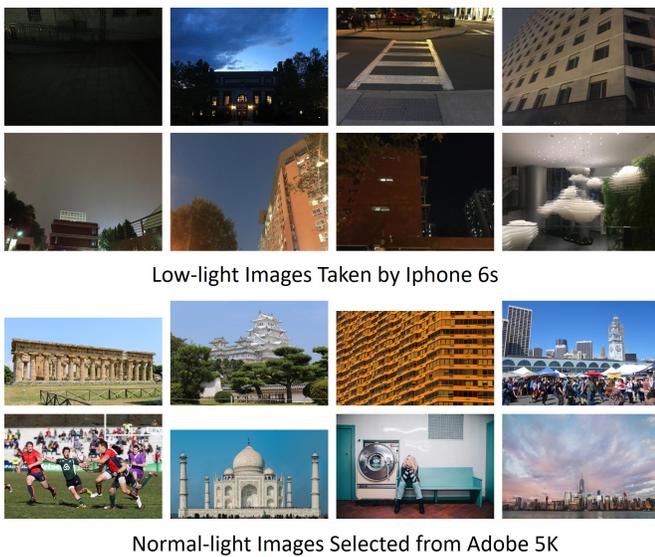


Fig. 3. Examples from our URL dataset. The top two rows display the low-light images collected with an iPhone 6s cell phone. We can see that the illumination conditions are quite diverse. Some images are very dark. The bottom two rows display the normal-light natural images collected from Adobe 5K.

which *does not contain noticeable noise*. Then we conduct experiments for both illumination enhancement and noise suppression on LOL dataset [29] and our collected unsupervised real-world low-light dataset (URL dataset), which *both contain low-light images with noticeable noise*.

A. Datasets

**Unpaired Enhancement Dataset.** Jiang et al. [21] collect an unpaired dataset for training contrast enhancement models.

The training set is composed of 914 low-light images which are dark yet *containing no significant noise*, and 1016 normal-light images from public datasets. We use this dataset to compare the performance of contrast enhancement of each model. The evaluation set is composed of 148 low-light/normal-light image pairs from public datasets. All the images from both the training and evaluation sets have been resized to  $400 \times 600$ .

**Low-Light (LOL) Dataset** [29]. LOL is composed of 500 low-light and normal-light image pairs and is split into 485 training pairs and 15 testing pairs. The low-light images contain noise produced during the photo capture process. Most of the images are indoor scenes. To adapt the dataset to our unsupervised setting, we adopt the training images as our low-light train set and adopt the normal-light images in the Unpaired Enhancement Dataset [21] as the normal-light train set. The testing images remain the same as the LOL dataset. All the images have a resolution of  $400 \times 600$ .

**URL Dataset.** There are a quite limited number of real-world low-light datasets publicly available. Among the public low-light datasets, some of them are composed of synthetic images while many other datasets such as ExDark [54] or Adobe FiveK [55] contain dark images without significant noise. These datasets do not meet the objective of our study. Therefore, we collect an Unsupervised Real-world Low-light dataset (URL dataset) composed of 414 high-resolution real-world low-light images taken by iPhone-6s and 3,837 normal-light images selected from Adobe FiveK. To collect the low-light images, we first take photos with an iPhone-6s from various scenes in different cities around the world. We remove images that are too dark (cannot be recovered since details are lost), blurry, or with high brightness. We also remove images that are very similar to other images in order to boost the diversity of the dataset. In the end, we are able to select 414

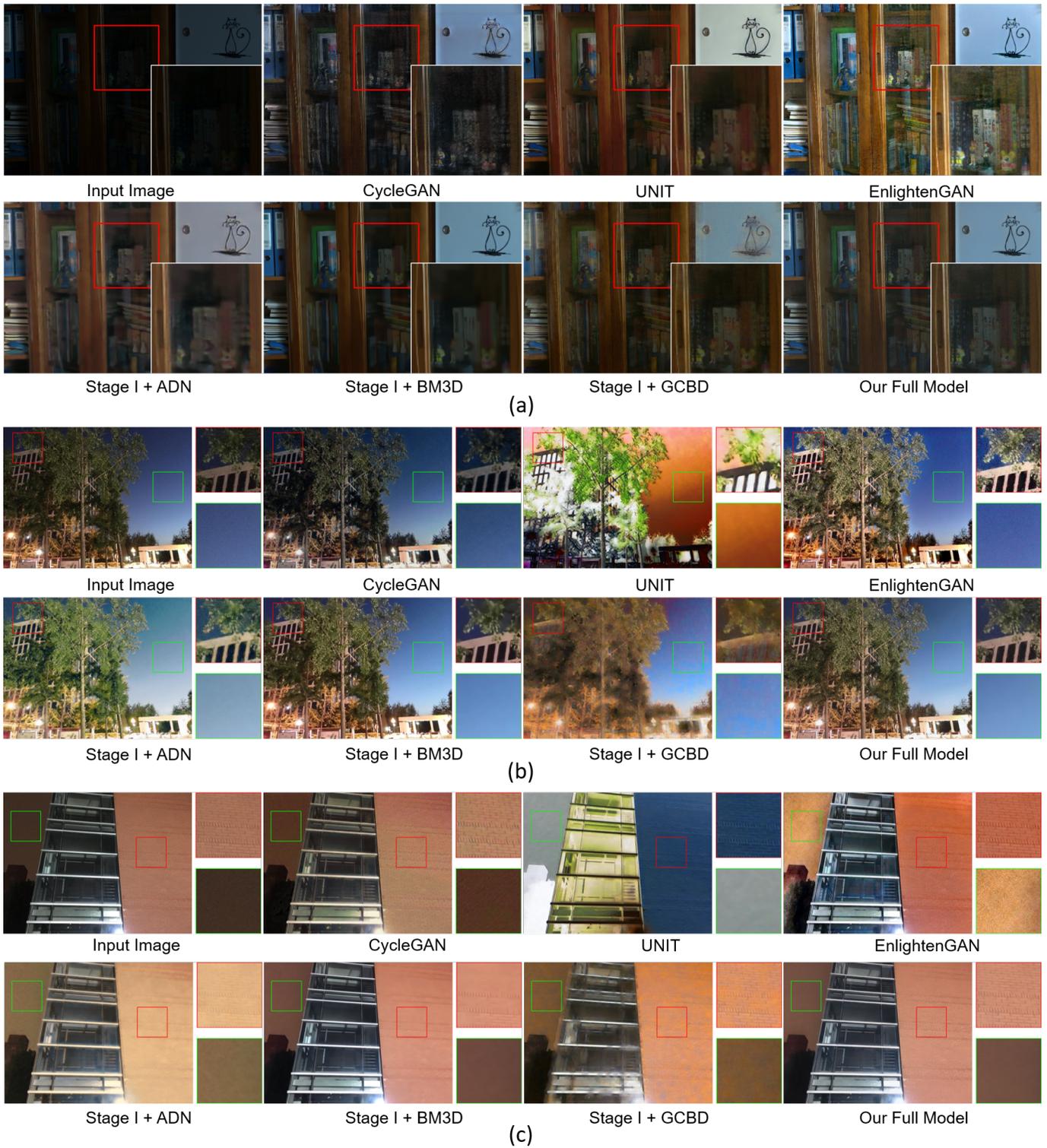


Fig. 4. Qualitative results (enhancement with denoising) on the LOL (a) and URL (b and c) datasets. *Please zoom in to the maximum for the very details.*

TABLE I  
QUANTITATIVE RESULTS FOR ILLUMINATION ENHANCEMENT ON THE UNPAIRED ENHANCEMENT DATASET.

Model	CycleGAN	UNIT	EnlightenGAN	Ours Stage I
PSNR	18.22	8.42	17.31	<b>19.78</b>
SSIM	0.7284	0.2549	0.8047	<b>0.8197</b>

low-light images from over 4,000 photos. Our URL dataset is quite diverse, containing various indoor and outdoor scenes under different light conditions. Consequently, the level of noise contained in each image or even different regions of the same image varies considerably across the dataset. We divide the low-light images into 328 training images and 86 testing images. This dataset thus complements the existing datasets in those two regards. Note that there is no corresponding ground-truth image for each testing image. Each low-light image is resized to  $1008 \times 756$ . Fig. 3 shows low-light image examples from our dataset.

### B. Implementation Details

**Implementation Details of Stage I** In Stage I, we first train the networks with a learning rate of 0.0001 using Adam optimizer for 100 epochs, then the learning rate is decreased linearly to zero within the next 100 epochs. We split the training data into train and validation sets, and select the best model on the validation set, then apply the model to the test images. We use a batch size of 32. We crop  $320 \times 320$  patches for training and randomly flip the patches for data augmentation. The whole model is trained on two 1080Ti GPUs.

**Implementation Details of Stage II** In Stage II, we train the networks with an initial learning rate of 0.0001 for 2,000 epochs. We train the model with randomly cropped image patches (with a size of  $128 \times 128$ ), as our model in this stage primarily needs to capture the noise pattern in the local regions. We randomly rotate and flip the patches for data augmentation. The whole model is trained on two 1080Ti GPUs.

### C. Experiments for Illumination Enhancement

We compare our Stage I model which does only illumination enhancement on the Unpaired Enhancement Dataset [21] with state-of-the-art models, including CycleGAN [17], UNIT [16] and EnlightenGAN [21]. As shown in Fig. 2, our model can generate normal-light images with reasonable contrast and color in both global and local regions. EnlightenGAN can produce visually pleasing images. However, they may still suffer from color distortion in several local regions, as indicated by the red boxes. The other unsupervised image translation methods can synthesis roughly good images. However, the color, contrast are not perfect.

We report the PSNR and SSIM of the generated images as a complement to the visual results on Unpaired Enhancement Dataset. Results in Table I show that our model performs significantly better than the existing models, which are consistent with the visual results.

TABLE II  
QUANTITATIVE RESULTS ON THE LOL DATASET.

Dataset Model	LOL		
	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
CycleGAN	14.75	0.6852	0.3808
UNIT	15.49	0.7280	0.3476
EnlightenGAN	18.36	0.7839	0.2915
Stage I + ADN	17.72	0.7776	0.3531
Stage I + BM3D	19.36	0.8154	0.2837
Stage I + GCBD	18.76	0.7753	0.3579
Our Full Model	<b>20.04</b>	<b>0.8216</b>	<b>0.2661</b>

TABLE III  
USER STUDY ON THE URL DATASET.

Model	POS
CycleGAN	3
UNIT	3
EnlightenGAN	28
Stage I + ADN	15
Stage I + BM3D	64
Stage I + GCBD	2
Our Full Model	<b>85</b>

### D. Experiments for Both Illumination Enhancement and Noise Suppression

**Experiment Settings.** We train the two stages of our model separately, as we observe that jointly training the two stages yields unstable results under our unsupervised learning setting. We evaluate our model on the real-world low-light image enhancement datasets: LOL and our URL datasets, and compare it with the state-of-the-art unsupervised image-to-image translation or contrast enhancement models, including CycleGAN [17], UNIT [16] and EnlightenGAN [21]. We also compare our model with the combination of illumination enhancement model and denoising model. Specifically, we compare our full model with our Stage I + BM3D [22], our Stage I + ADN [38], [56] and our Stage I + GCBD [28]. BM3D [22] is a robust image denoising method. The limitation is that it requires a known noise level as input. To use BM3D in our task, we first estimate a rough noise level for each testing image, then apply BM3D on the testing images.

**Quantitative Results.** On the LOL dataset, we report PSNR, SSIM results. We also report the perceptual score (LPIPS) of the enhanced images to better quantify the perceptual quality of images. From Table II we can see that our model is consistently better than the existing models, demonstrating the superiority of our decoupled networks.

**Qualitative Results.** Fig. 4 show the qualitative results on both the LOL dataset and our URL dataset. CycleGAN generates heavy artifacts and slightly suffers from color distortion. UNIT suffers from heavy color distortion on the URL dataset and cannot preserve details on the LOL dataset. EnlightenGAN can improve the illumination of the images, but there are still obvious noise and artifacts on the image. *The visual results indicate that it is challenging to handle image illumination enhancement and denoising simultaneously with a single model for real-world low-light image enhancement.*

However, simply cascading an illumination enhancement model with a denoising model still cannot produce satisfactory

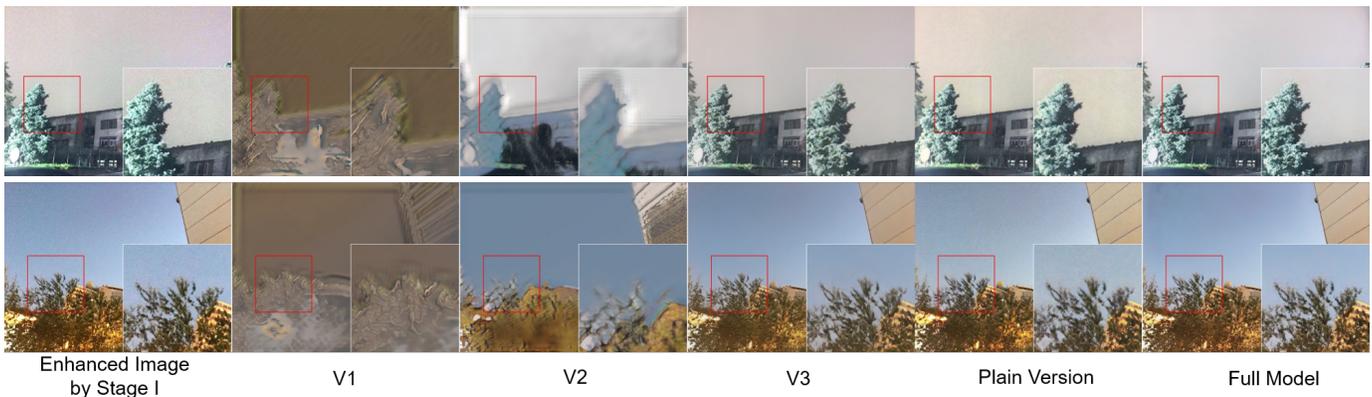


Fig. 5. Ablation Study on our URL dataset. We strongly encourage readers to zoom in to the maximum for the very details.

TABLE IV  
CONFIGURATION OF DIFFERENT VERSIONS OF OUR MODEL.

Model	$L_G$	$L_{color}$	$L_{con}^{adapt}$	$L_{con}$
Version 1	✓			
Version 2	✓	✓		
Version 3	✓	✓	✓	
Full Model	✓	✓	✓	✓
Plain	✓	✓	vanilla	✓

TABLE V  
NON-REF IMAGE QUALITY SCORES ON URL. (LOWER IS BETTER)

Model	V3	Plain	Full
Brisque	33.66	26.95	<b>22.13</b>
NIQE	4.12	3.44	<b>2.53</b>

results. From Fig. 4, we see that the results of ADN and GCBD still exhibit color distortion or contain many artifacts on both datasets. ADN even suffers from overexposure in bright regions on URL dataset. Using BM3D to post-process the results of our Stage I yields good color and illumination. However, from Fig. 4, we observe that the texture details of images on both the LOL and URL datasets are over-smoothed. A possible reason is that BM3D was proposed for synthetic noise removal. It may not generalize well to real-world denoising. Compared to existing methods, our full model performs noise removal adaptively and considers preserving the details and color when denoising. Therefore, It can suppress noise as well as preserving more details.

**User Study.** We perform a user study on the URL dataset. Specifically, we randomly select 20 low-light testing images from the URL dataset. For each input image, we show each user the final enhanced result of each method and ask the user to select the most visually pleasing result among all methods, considering illumination condition, color, texture realness, and noise level. In total, we obtain 200 preference opinions. Table III shows the preference opinion score (POS) of each method. Our method outperforms state-of-the-art methods. It is worth noting that although EnlightenGAN can generate images with good illumination, it fails to remove noise when noise is heavy, as shown in Fig. 4. As a result, its POS is relatively low. The other one-stage models CycleGAN and UNIT suffer from color distortion. The results further indicate that it is challenging for a single model to perform illumination enhancement and noise removal simultaneously.

### E. Ablation Study

In this section, we study how each component of our model contributes to the final performance. We primarily analyze the components in Stage II, which are the core contributions in this work. Specifically, as shown in Table IV, we compare the versions of our Stage II model with different losses imposed. Besides the versions regarding the loss functions, we also compare our model to a non-adaptive denoising model, which we call the Plain Model, *i.e.*, the generator only takes the enhanced image as input without illumination guidance. The four losses used in the Plain Model remain the same as our Full Model, except that we use the vanilla content loss instead of adaptive content loss.

Experiments are conducted on our URL dataset. From Fig. 5 we observe that the results from Stage I still contain heavy noise, further indicating that it is challenging to model both illumination and noise patterns with a single network. Merely using adversarial loss (Version 1) can help to smooth the image, but the color is shifting and cannot be well controlled. Using the color loss (Version 2) significantly helps to constrain the output image to have similar color and contrast as the input image. However, without learning with *pseudo triples* and adaptive content loss (*i.e.*, loss  $L_{con}^{adapt}$ ), the content of the output is not explicitly constrained and the images contain distorted contextual details. When imposing the learning with *pseudo triples* (Version 3), the network can produce realistic contents and perform noise suppression, *indicating that Pseudo Labeling plays a key role in stabilizing the training and improving the performance of noise modeling*. However, several local regions still lack details. Please pay attention to the trees in these results in Fig. 5. The texture details of the trees and leaves are smoothed in Version 3 and preserved in our Full Model, indicating that the content loss between the real input image and its output image further helps to preserve



Fig. 6. Failure cases of our method. In each row, from left to right we present the input low-light image, the enhanced image from our Stage I model, and the final denoised result from our Stage II model. Please zoom in to see the details and pay special attention to the red marked regions.

contextual details during denoising.

Comparing our Full Model in Stage II with the Plain Model, we observe that the Plain Model cannot effectively suppress the noise. There is still notable noise in the sky region and other regions, as shown in Fig. 5. A possible reason is that the plain model may not be able to perceive all the noise patterns under different illumination conditions precisely. In contrast, our model is explicitly guided by the illumination of the image. Therefore it can capture the noise pattern under various illumination conditions more effectively and produce better results.

Since the URL dataset does not contain ground truth images, we use non-reference image quality assessment methods BRISQUE [57] and NIQE [58] to quantify the quality of the testing images. Results are shown in Table V. We do not include the scores of V1 or V2, as these two versions can only produce heavily distorted images, as shown in Fig. 5. These quantitative results further demonstrate the importance of using illumination as guidance for real-world noise modeling.

#### F. Failure Cases

Fig. 6 shows two failure cases of our model on the URL dataset. In the top case, there is still noise at the upper-right corner of the denoised image. In the bottom case, there is still noise at the bottom-right corner of the denoised image, as indicated by the red circles. These results demonstrate our model’s limitation on dealing with image borders or corners. A possible reason may be that the noise pattern at the corners or the borders of the image may be different from that of the other regions, since we are dealing with the real-world noise that is spatially variant. We will continue working on it to improve our model in the future.

#### V. CONCLUSION

In this paper, we have presented the decoupled networks to address the real-world low-light image enhancement problem

in an unsupervised fashion. Our model in Stage I enhances a low-light image to generate an image with satisfactory illumination and color. Our model in Stage II further denoises the enhanced image to obtain a clean image, while preserving good contrast, color and contextual details. We conduct experiments on three real-world datasets. The results show that our model outperforms the state-of-the-art models in terms of both illumination enhancement and noise removal.

#### APPENDIX

##### ARCHITECTURE OF OUR MODEL

##### A. Generator of Stage I

The generator of Stage I is composed of an encoder, a pyramid module, and a decoder. The encoder is composed of several convolutional blocks and Max pooling layers. Each convolutional block is composed of a convolution layer, a Batch Normalization layer (except the first convolutional block), and a Leaky ReLU activation function with slope to be 0.2.

The decoder is composed of three Deconvolutional blocks. Each Deconvolutional block is composed of an upsampling layer, two convolution layers, a Batch Normalization layer and a Leaky ReLU activation layer (the activation of the last Deconvolutional block is Sigmoid instead of Leaky ReLU).

The overall structure of our generator in Stage I is:  $CCMCCMCCMPDDD$ , where  $C$  denotes to Convolutional block,  $M$  means Max pooling,  $P$  means pyramid module,  $D$  means deconvolutional block.

##### B. Discriminators of Stage I

In Stage I, we have a global discriminator and a local discriminator. The global discriminator is composed of five convolutional layers. Each convolution layer is followed by a Leaky ReLU activation function layer, except the last layer. The local discriminator has the same structure as the global discriminator.

### C. Generator of Stage II

The generator of Stage II is composed of an encoder, six Resnet blocks and a decoder. The encoder is composed of two convolutional blocks. The decoder is composed of three deconvolutional blocks, as introduced in the previous section. The detailed structure of our generate in Stage II is:  $CCRRRRRRRDDDD$ , where  $R$  means a Resnet block.

### D. Discriminator of Stage II

We use two-scale discriminators in Stage II. Each discriminator has the same structure as the discriminators used in Stage I. The only difference between discriminators at different scale is the input. The input to the first discriminator is the original full-sized image. The input to the second discriminator is a downsized version of the full-sized image. We use Average pooling to downsize the image.

## REFERENCES

- [1] S. Lee, "An efficient content-based image enhancement in the compressed domain using retinex theory," *IEEE transactions on circuits and systems for video technology*, vol. 17, no. 2, pp. 199–213, 2007.
- [2] Y. Ren, Z. Ying, T. H. Li, and G. Li, "Lecarm: Low-light image enhancement using the camera response model," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 29, no. 4, pp. 968–981, 2018.
- [3] Z. Zhao, B. Xiong, L. Wang, Q. Ou, L. Yu, and F. Kuang, "Retinexpdip: A unified deep framework for low-light image enhancement," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [4] J. Li, X. Feng, and Z. Hua, "Low-light image enhancement via progressive-recursive network," *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [5] R. Mantiuk, S. Daly, and L. Kerofsky, "Display adaptive tone mapping," in *ACM SIGGRAPH 2008 papers*, 2008, pp. 1–10.
- [6] M. Gharbi, J. Chen, J. T. Barron, S. W. Hasinoff, and F. Durand, "Deep bilateral learning for real-time image enhancement," *ACM Transactions on Graphics (TOG)*, vol. 36, no. 4, pp. 1–12, 2017.
- [7] R. Wang, Q. Zhang, C.-W. Fu, X. Shen, W.-S. Zheng, and J. Jia, "Underexposed photo enhancement using deep illumination estimation," in *CVPR*, 2019, pp. 6849–6857.
- [8] K. Xu, X. Yang, B. Yin, and R. W. Lau, "Learning to restore low-light images via decomposition-and-enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 2281–2290.
- [9] C. Guo, C. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, "Zero-reference deep curve estimation for low-light image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1780–1789.
- [10] S. Moran, P. Marza, S. McDonagh, S. Parisot, and G. Slabaugh, "Deepplf: Deep local parametric filters for image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12 826–12 835.
- [11] W. Yang, S. Wang, Y. Fang, Y. Wang, and J. Liu, "From fidelity to perceptual quality: A semi-supervised approach for low-light image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 3063–3072.
- [12] K. G. Lore, A. Akintayo, and S. Sarkar, "Llnet: A deep autoencoder approach to natural low-light image enhancement," *Pattern Recognition*, vol. 61, pp. 650–662, 2017.
- [13] S. Guo, Z. Yan, K. Zhang, W. Zuo, and L. Zhang, "Toward convolutional blind denoising of real photographs," in *CVPR*, 2019, pp. 1712–1722.
- [14] C. H. Martin and M. W. Mahoney, "Rethinking generalization requires revisiting old ideas: statistical mechanics approaches and complex learning behavior," *arXiv preprint arXiv:1710.09553*, 2017.
- [15] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.
- [16] M.-Y. Liu, T. Breuel, and J. Kautz, "Unsupervised image-to-image translation networks," in *Advances in neural information processing systems*, 2017, pp. 700–708.
- [17] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *ICCV*, 2017, pp. 2223–2232.
- [18] W. Xiong, Y. He, Y. Zhang, W. Luo, L. Ma, and J. Luo, "Fine-grained image-to-image transformation towards visual recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 5840–5849.
- [19] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [20] Y.-S. Chen, Y.-C. Wang, M.-H. Kao, and Y.-Y. Chuang, "Deep photo enhancer: Unpaired learning for image enhancement from photographs with gans," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 6306–6314.
- [21] 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 Transactions on Image Processing*, vol. 30, pp. 2340–2349, 2021.
- [22] 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.
- [23] J. Lehtinen, J. Munkberg, J. Hasselgren, S. Laine, T. Karras, M. Aittala, and T. Aila, "Noise2noise: Learning image restoration without clean data," in *ICML*, 2018, pp. 2965–2974.
- [24] A. Krull, T.-O. Buchholz, and F. Jug, "Noise2void-learning denoising from single noisy images," in *CVPR*, 2019.
- [25] J. Batson and L. Royer, "Noise2self: Blind denoising by self-supervision," *arXiv preprint arXiv:1901.11365*, 2019.
- [26] D.-H. Lee, "Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks," in *Workshop on challenges in representation learning, ICML*, vol. 3, 2013, p. 2.
- [27] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. A. Raffel, "Mixmatch: A holistic approach to semi-supervised learning," in *Advances in Neural Information Processing Systems*, 2019, pp. 5050–5060.
- [28] J. Chen, J. Chen, H. Chao, and M. Yang, "Image blind denoising with generative adversarial network based noise modeling," in *CVPR*, 2018, pp. 3155–3164.
- [29] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," in *BMVC*, 2018.
- [30] S. M. Pizer, E. P. Amburn, J. D. Austin, R. Cromartie, A. Geselowitz, T. Greer, B. ter Haar Romeny, J. B. Zimmerman, and K. Zuiderveld, "Adaptive histogram equalization and its variations," *Computer vision, graphics, and image processing*, vol. 39, no. 3, pp. 355–368, 1987.
- [31] E. H. Land, "The retinex theory of color vision," *Scientific american*, vol. 237, no. 6, pp. 108–129, 1977.
- [32] D. J. Jobson, Z.-u. Rahman, and G. A. Woodell, "A multiscale retinex for bridging the gap between color images and the human observation of scenes," *IEEE Transactions on Image processing*, vol. 6, no. 7, pp. 965–976, 1997.
- [33] X. Fu, D. Zeng, Y. Huang, X.-P. Zhang, and X. Ding, "A weighted variational model for simultaneous reflectance and illumination estimation," in *CVPR*, 2016, pp. 2782–2790.
- [34] X. Guo, Y. Li, and H. Ling, "Lime: Low-light image enhancement via illumination map estimation," *IEEE Transactions on Image Processing*, vol. 26, no. 2, pp. 982–993, 2016.
- [35] M. Li, J. Liu, W. Yang, X. Sun, and Z. Guo, "Structure-revealing low-light image enhancement via robust retinex model," *IEEE Transactions on Image Processing*, vol. 27, no. 6, pp. 2828–2841, 2018.
- [36] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," in *CVPR*, 2018, pp. 3291–3300.
- [37] A. Ignatov, N. Kobyshev, R. Timofte, K. Vanhoey, and L. Van Gool, "Wespe: weakly supervised photo enhancer for digital cameras," in *CVPR Workshops*, 2018, pp. 691–700.
- [38] H. Liao, W.-A. Lin, Y. Jianbo, S. K. Zhou, and J. Lou, "Artifact disentanglement network for unsupervised metal artifact reduction," in *International Conference on Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, 2019.
- [39] A. Buades, B. Coll, and J.-M. Morel, "A non-local algorithm for image denoising," in *CVPR*, vol. 2. IEEE, 2005, pp. 60–65.
- [40] 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.
- [41] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image restoration," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

- [42] D. Liu, B. Wen, Y. Fan, C. C. Loy, and T. S. Huang, "Non-local recurrent network for image restoration," in *Advances in Neural Information Processing Systems*, 2018, pp. 1673–1682.
- [43] J. Xu, L. Zhang, D. Zhang, and X. Feng, "Multi-channel weighted nuclear norm minimization for real color image denoising," in *ICCV*, 2017, pp. 1096–1104.
- [44] D.-W. Kim, J. Ryun Chung, and S.-W. Jung, "Grdn: Grouped residual dense network for real image denoising and gan-based real-world noise modeling," in *CVPR Workshops*, 2019, pp. 0–0.
- [45] M. Lebrun, M. Colom, and J.-M. Morel, "Multiscale image blind denoising," *IEEE Transactions on Image Processing*, vol. 24, no. 10, pp. 3149–3161, 2015.
- [46] 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.
- [47] J. Xu, L. Zhang, and D. Zhang, "A trilateral weighted sparse coding scheme for real-world image denoising," in *ECCV*, 2018, pp. 20–36.
- [48] H. Yan, V. Tan, W. Yang, and J. Feng, "Unsupervised image noise modeling with self-consistent gan," *arXiv preprint arXiv:1906.05762*, 2019.
- [49] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.
- [50] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia, "Pyramid scene parsing network," in *CVPR*, 2017.
- [51] B. Xu, N. Wang, T. Chen, and M. Li, "Empirical evaluation of rectified activations in convolutional network," *arXiv preprint arXiv:1505.00853*, 2015.
- [52] A. Jolicoeur-Martineau, "The relativistic discriminator: a key element missing from standard gan," *arXiv preprint arXiv:1807.00734*, 2018.
- [53] C. Poynton, *Digital video and HD: Algorithms and Interfaces*. Elsevier, 2012.
- [54] Y. P. Loh and C. S. Chan, "Getting to know low-light images with the exclusively dark dataset," *Computer Vision and Image Understanding*, vol. 178, pp. 30–42, 2019.
- [55] V. Bychkovsky, S. Paris, E. Chan, and F. Durand, "Learning photographic global tonal adjustment with a database of input / output image pairs," in *The Twenty-Fourth IEEE Conference on Computer Vision and Pattern Recognition*, 2011.
- [56] H. Liao, W. Lin, S. K. Zhou, and J. Luo, "Adn: Artifact disentanglement network for unsupervised metal artifact reduction," *IEEE Transactions on Medical Imaging*, 2019.
- [57] A. Mittal, A. K. Moorthy, and A. C. Bovik, "No-reference image quality assessment in the spatial domain," *IEEE Transactions on image processing*, vol. 21, no. 12, pp. 4695–4708, 2012.
- [58] A. Mittal, R. Soundararajan, and A. C. Bovik, "Making a "completely blind" image quality analyzer," *IEEE Signal processing letters*, vol. 20, no. 3, pp. 209–212, 2012.