# Learning an Adaptive Model for Extreme Low-light Raw Image Processing

*Qingxu Fu[1], Xiaoguang Di[1*], Yu Zhang[2]*

[1] *Control and Simulation Center, Harbin Institute of Technology, Harbin, 150080, PeopleâĂŹs Republic of China*
[2] *National Key Laboratory of Tunable Laser Technology, Harbin Institute of Technology, Harbin, 150080, PeopleâĂŹs Republic of China*
*\* E-mail: dixiaoguang@hit.edu.cn*

**Abstract:** Low-light images suffer from severe noise and low illumination. Current deep learning models that are trained with real-world images have excellent noise reduction, but a ratio parameter must be chosen manually to complete the enhancement pipeline. In this work, we propose an adaptive low-light raw image enhancement network to avoid parameter-handcrafting and to improve image quality. The proposed method can be divided into two sub-models: Brightness Prediction (BP) and Exposure Shifting (ES). The former is designed to control the brightness of the resulting image by estimating a guideline exposure time $t_1$. The latter learns to approximate an exposure-shifting operator $ES$, converting a low-light image with real exposure time $t_0$ to a noise-free image with guideline exposure time $t_1$. Additionally, structural similarity (SSIM) loss and Image Enhancement Vector (IEV) are introduced to promote image quality, and a new Campus Image Dataset (CID) is proposed to overcome the limitations of the existing datasets and to supervise the training of the proposed model. Using the proposed model, we can achieve high-quality low-light image enhancement from a single raw image. In quantitative tests, it is shown that the proposed method has the lowest Noise Level Estimation (NLE) score compared with the state-of-the-art low-light algorithms, suggesting a superior denoising performance. Furthermore, those tests illustrate that the proposed method is able to adaptively control the global image brightness according to the content of the image scene. Lastly, the potential application in video processing is briefly discussed.

## 1 Introduction

Images play an irreplaceable role in industry, military and entertainment. As the art of light, how to obtain better images in the low-light environment has been extensively studied in the literature. Some advanced photography equipment can help but at an expensive cost. On the other hand, there are few limitations to post-process low-light images. However, it is challenging due to problems such as noise, color distortion et al.
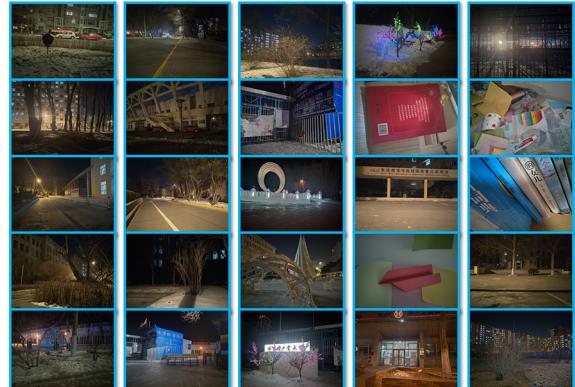
Exposure time, ISO (which measures the sensitivity of the image sensor) and aperture are known as three pillars of photography. With extended exposure time, the low-light problem can be easily solved, but it is not realistic because if the camera is not fixed or the scene contains moving objects. Higher ISO introduces more noise and is not always available for mobile devices. As a flexible solution, low-light image enhancement provides an alternative for imaging in the low-light environment.

In general, current low-light image enhancement method can be classified as follows:

Classic low-light image enhancement methods, with no application of convolutional neural networks. Histogram equalization (HE)[1] and gamma correction[2] provide simple solution for low-light image processing. Based on Retinex theory[3], Single-scale Retinex (SSR)[4], Multi-scale Retinex (MSR)[5], SRIE[6] and LIME[7] were proposed, enhancing low-light image by estimating the illumination map and reflectance map. Dehazing[8] methods can also be applied to this subject, regarding low-light images as inversed hazed image.

With the rapid development of the deep neural network, researchers have combined classic low-light image enhancement methods with Convolutional Neural Networks (CNNs). Chen et al. proposed Retinex-Net[9] to decompose low-light image into illumination and reflectance and use BM3D[10] for denoising.

Those methods assume low-light images are free of noise or noise is already removed by additional denoising process. Therefore, when applying those methods on real-world low-light images, an additional denoising process is required. However, most low-light images suffer from severe noise, traditional denoising methods like BM3D are not able to provide satisfactory results. To be rid of image noise,
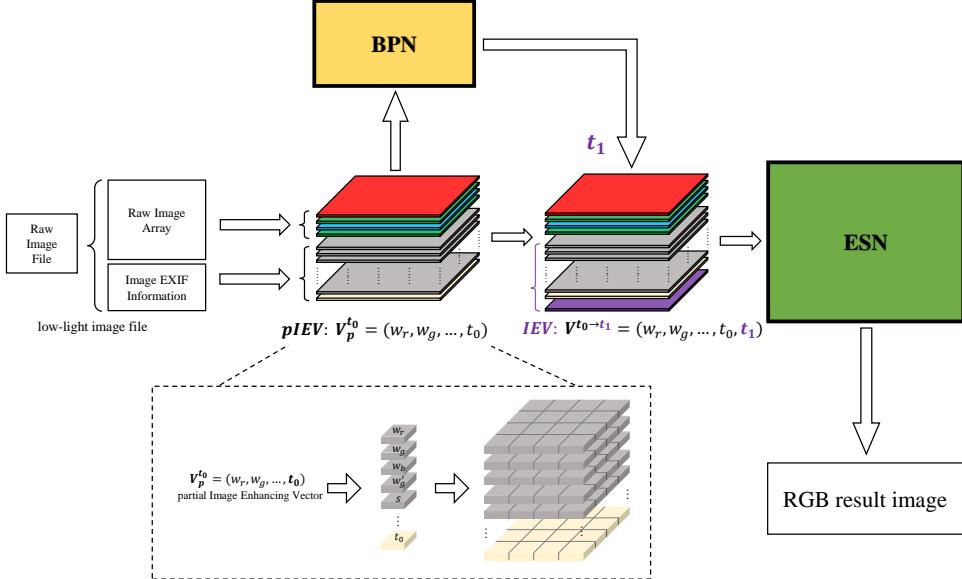


**Fig. 1**: Several scenes in CID datasets.

researchers also explored the method of deep learning[11–14] in denoising topics.

There are some low-light images which are barely visible before post processing, as they are captured in extremely dark environments or with very short exposure time. In this work, those images are named as extreme low-light images. While the classic methods are unable to tackle the severe noise and serious color distortion in those images, it is then found that with the help of deep learning and raw-format image datasets, we can still obtain satisfactory results. Chen et al. proposed an end-to-end solution for raw low-light image enhancement[15], dealing with the low-light problem and denoising in a single model. However, extra brightness amplification ratio is needed and requires manual parameter tuning for different scenes, gaining excellent denoising effects but losing the flexibility and adaptability of classic methods.

Currently, all datasets available for low-light enhancement models chose a longer-exposed image as the ground-truth image in training, which was heavily dependent on the dataset collector's experience to select exposure parameters such as exposure time, ISO and white

**Fig. 2**: The proposed framework. For a underexposed low-light raw image with guideline exposure time $t_1$ already known, the Exposure Shifting Network (ESN) is applied to complete the estimation of the RGB result image. In a more general case, Brightness Prediction Network (BPN) is utilized to give a proper estimation for the uncertain parameter $t_1$ with raw image array and available Exif metadata.

balance mode. Thus these parameters were usually not optimal, considering the content of the scene, illumination of the environment et al., which can lead to an undesired trained model.

It is hard for deep CNNs to learn how we determine the best exposure parameters for the ground-truth images because this can be very subjective and lacks specific standards. It can lead to more problems if there is more than one dataset collector. On the other hand, making a baseline for the ground-truth selection is equally difficult. Nevertheless, it is possible to predict how an image varies when exposure parameters change. We name this process Exposure Shifting (ES). With a learned ES model, it becomes possible to reversely study what good exposure parameters (e.g. exposure time) *should* be suitable for each low-light image, utilizing the characteristics of the back-propagation algorithm.

The contributions of our work are summarized as follows:

1) We presented Campus Image Dataset (CID), which contains a variety of scenes captured by multi-level exposure. CID provides powerful support for the training of our data-driven low-light enhancement model.

2) We proposed an adaptive two-stage low-light image enhancement model, which provides state-of-the-art low-light noise suppression as well as adaptive global brightness control effect.

In stage one, a Brightness Prediction Network (BPN) was introduced to estimate a proper exposure time based on the content of images as well as Exif (Exchangeable image file) metadata. BPN was trained to provide adaptive image brightness control during image enhancement and to get rid of the external parameters required in the Exposure Shifting stage.

In stage two, an Exposure Shifting Network (ESN) was proposed to estimate a longer-exposed version of the image with target exposure time estimated by BPN. Moreover, ESN also completes image denoising and ends up with the final resulting RGB image.

The rest of the article is organized as follows. In Section 2, previous research works are provided. In Section 3, we demonstrate the importance and advantages of our CID (Campus Image Dataset) dataset. In Section 4, our model and its training details are provided, as well as the design concept of the model. The proposed model is compared with other state-of-the-art methods in Section 5, from multiple perspectives including denoising, adaptive brightness control performance and computational cost. The possible extension in video processing is briefly discussed here. In Section 6, the advantages and the limitations of our model are concluded, and we also present the future expectations of this work.

## 2 Related Works

The method we proposed was inspired not only by the application of deep neural networks, but also by the classic methods based on the Histogram Equalization model, dehazing model and Retinex model.

### 2.1 Histogram Equalization Methods

Histogram Equalization (HE) method is used extensively to process digital images. The basic idea is to improve image visibility by changing the image histogram into a uniform distribution, which effectively increases image entropy for low-light images. In [16] it is argued that a color-invariant representation can be obtained by applying HE. However, HE tends to cause noise amplification, details disappearance and color distortion in low-light image enhancement.

On the account of the drawbacks of HE, a number of improvements have been put forward. Adaptive Histogram Equalization (AHE) [17] and its variation Contrast-limited Adaptive Histogram Equalization (CLAHE) [18] perform the histogram equalization considering not only the global histogram, but also the neighbors of each pixel. Hue-preserving color image enhancement [19] works on contrast-enhanceing and hue-preserving, [20] and [21] on preserving the original image luminance.

### 2.2 Dehazing and Retinex Model-Based Methods

A number of low-light enhancement methods are based on the dehazing and Retinex model. It is observed that low-light images and inversed haze images share many similarities. For this reason, Dong et al. [8] presented a method and achieved the low-light image enhancement in which the low-light images is inversed, dehazed and then inversed back again. This method was studied further by [22] and [23].

On the other hand, the Retinex model [3] revealed that a natural image is made up of the illumination map and the reflectance map. With Retinex decomposition and the assumption that the illumination map is smooth, researchers have proposed single-scale Retinex [24] and multi-scale Retinex [4]. Based on image fusion, adjustments can be applied to the illumination map to improve the performance [25]. [26] focused on preserving natural characteristics with a Bright-Pass filter. In [27], a variational Retinex model was proposed and the illumination estimation problem was formulated

as a Quadratic Programming optimization problem. In [28] illumination map is estimated considering local structure and [29] focused on revealing the structure details from the reflectance map. A recent study shows that the Retinex model can combine with the atmospheric scattering model[30]. More refined models were presented by [31, 32] based on variational Retinex theory.

However, noise modeling is not well established in those algorithms. In [8] it is assumed that noise is insignificant and can be removed before or after the Dehazing stage. And the Retinex based models rely on classic noise reduction algorithms (e.g. BM3D) to reduce noise found in the reflectance map. Therefore it can achieve good results in mild low-light images but tends to amplify noise when processing severe or extreme low-light images. Besides, the variational Retinex model is very time-consuming as it takes many iterations to solve the optimization objectives, thus it is hard to achieve real-time processing.

### 2.3  Data-driven Methods

Image-denoising and low-light enhancement can be realized with the data-driven approach, separately or integrally.

Noise reduction was the subject of extensive studies. Burger et al.[11] discovered that convolutional neural networks could compete with the state-of-the-art classic denoising algorithm BM3D. The works on date-driven denoising have been extended by [14, 33, 34]. Although these works were not targeted at low-light image processing, they provided some references for the research of other image processing studies.

The first application of the data-driven method is LLNET [12], which utilized a deep auto-encoder to learn from synthetic low-light image datasets and the denoising procedure is integrated into the low-light enhancement process. Also, it was demonstrated that Retinex based methods can be further improved by deep learning [9, 35], but those works still have undesirable denoising performance because synthetic datasets are used and they cannot reflect the characteristics of real low-light images.

To solve the drawbacks of synthetic datasets, a large multi-exposure image dataset was proposed in [36]. It was found by Chen et al.[15] that models utilizing raw-format images can provide much better results in low-light enhancement. In the SID model proposed in [15], it suggested that the state-of-the-art results can be achieved by using raw-format paired images and training in an end-to-end way. However, since the paired images have different exposure time, a ratio was introduced as an additional input to bridge the gap. Consequently, if the SID model is applied on images out of the dataset, manual adjustment is requires because there is no paired images to indicate the ratio, which results in many drawbacks in practical applications.
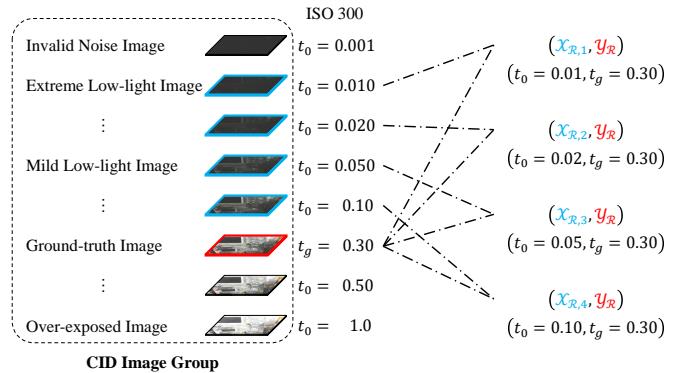
## 3  Dataset

For extreme low-light image enhancement, there are few optional datasets. The Google HDR+ dataset[37] and Darmstadt Noise Dataset[38] avoid the disadvantages of synthetic datasets, but they were mainly collected in environments with sufficient illumination. The RENOIR[39] and learning-to-See-In-the-Dark dataset (SID)[15] provides real low-light image pairs, which is captured by carefully selecting the ISO and exposure time for each scene. The Exclusively Dark Dataset (EDD) [40] is another low-light dataset with object-level annotations, which is targeting at object recognition. But these datasets cannot meet the requirements of our study. In our work, in order to adaptively enhance images with different noise levels, it is needed to consider how the exposure parameters have an impact on low-light images. More specifically, as the environment illumination grows lower, low-light images are gradually overwhelmed by noise and invalid pixels. Although it is not flexible to manipulate the environment illumination, the camera exposure time setting can be easily changed to simulate this illumination shift. Therefore, it is expected that there is a dataset containing, instead of just some image pairs, a number of multi-exposed image series, in

each of which there are low-light images of different levels and one reference image captured in the same scene.

Overall, there are three conditions to be met for a satisfactory dataset:

- Real scenes. Unlike synthetic images, photos captured in real scenes reflect much more diverse noise and distortion pattern
- Multiple exposure levels. Using a number of multi-exposed image series to train an adaptive model capable of enhancing low-light images of different levels.
- Raw-format images. Most data captured by the camera sensor is lost in format convertion, those lost infomation is essential for extreme low-light image enhancement.

**Fig. 3**: Demonstration of a CID image group.

Based on the above description, we proposed our Campus Image Dataset (CID)*. There are about 200 image groups in CID, each of which consists of 8 raw images with different exposure time, shot continuously in the same scene using a tripod (Fig.3). To be specific, for every group, an upper limit of exposure time was selected to make sure that the first image in the sequence was slightly overexposed, then a lower limit was chosen to capture an extreme low-light image as the last image in the sequence, and the gap between the limits was filled up by another 6 images. Finally, the ground-truth image for each group was manually selected, and invalid images that had nothing but noise due to extreme short exposure were tagged and excluded. A demonstration of scenes in CID is presented in Fig.1.
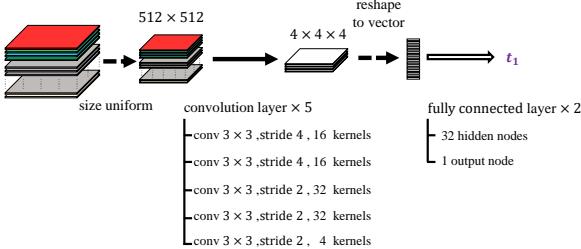
## 4  Method

### 4.1  Method Overview

Most of the existing methods are unable to suppress the severe noise present in extreme low-light images. On the other hand, recent studies e.g. [15] have concentrated on the extreme low-light enhancement with raw-format images and real-scene datasets in an end-to-end way, but unlike the classic methods, the brightness of the enhanced image cannot be automatically controlled in this framework, more specifically, an external ratio has to be manually adjusted when the input image has no paired image as the reference.

It is our aim to find a solution to enhance extreme low-light images, combining the advantages of adaptive brightness control and superior denoising performance. It is found that the problem of low-light enhancement can be considered as a special case of changing the time of exposure. With $t_0$ and $t_1$ representing exposure time, the low-light image $\mathcal{X}_{t_0}$ can be interpreted as a normal image $\mathcal{Y}_{t_1}$ corrupted by a hypothetical exposure-shifting operator $ES$, which only

---

**Fig. 4**: The components of BPN.

alters the exposure time but keeps the aperture and ISO unchanged:

$$BX_{t_0} = ES(\mathcal{Y}_{t_1}, t_0) \qquad (1)$$

In this process, more stochastic noise $\mathcal{N}$ is generated by operator $ES$ because $t_0 < t_1$, corresponding to the low PSNR (Peak Signal to Noise Ratio) value in low-light images. Contrariwise, the corrupted low-light image can also be restored by the same operator $ES$:

$$\hat{\mathcal{Y}}_{t_1} = ES(\mathcal{X}_{t_0}, t_1) \qquad (2)$$

The noise $\mathcal{N}$ is reduced because $t_1 > t_0$. Thus the normal image $\hat{\mathcal{Y}}_{t_1}$ is reconstructed.

The latter process, named as Exposure Shifting (ES), can enhance low-light images as well as suppress image noise. More importantly, it can be learned in a data-driven and end-to-end way, thus it is expected to have superior denoising performance. However, a proper guideline exposure time $t_1$ always has to be provided in the reconstruction $\hat{\mathcal{Y}}_{t_1}$. To make our model adaptive, a Brightness Prediction (BP) procedure is introduced to achieve guideline time estimation.

In brief, the extreme low-light enhancement process is split into two procedures. Firstly, for a raw-format low-light image (with real exposure time $t_0$), an optimal guideline exposure time $t_1$ is estimated via a sub-model called Brightness Prediction Network (BPN). Secondly, the final RGB-format enhanced image is obtained by another sub-model, namely Exposure Shifting Network (ESN), using estimated guideline time $t_1$ to control the image brightness.

In addition, besides the image itself, some extra data are involved in both BPN and ESN, such as white balance index, ISO, $t_0$ and $t_1$. The Image Enhancing Vector (IEV) is put forward to encode the involved extra data.

## 4.2 Image Enhancing Vector

Raw image array is the original data captured by the camera. Other key parameters, such as camera white balance, ISO, exposure time, time and date, are stored together with the raw image array as Exif metadata. In HDR imaging[41], exposure time and ISO information have played an essential role in calculating the response curve and LDR to HDR conversion. Although it is not necessary to explicitly train the network to learn the camera response curve, feeding these additional parameters into the network as input avoids underfitting caused by insufficient features.

Image Enhancing Vector (IEV) is proposed to introduce extra features into the convolution neural networks. In this paper, IEV consists of ISO $u_s$, white balance indices for 4 raw image channels $(w_r, w_g, w_b, w_{g2})$, exposure time of the low-light image $t_0$ and guideline exposure time $t_1$. Moreover, it is possible to append other Exif metadata into IEV to improve network performance when necessary,e.g. aperture and focal length.

The IEV is used as network input in both BPN (for $t_1$ estimation) and ESN (for exposure shifting), but a difference exists. As it is shown in Fig.2, in ES procedure the IEV can be written as:

$$\mathcal{V}^{t_0 \to t_1} = \mathcal{V}(t_0, t_1) = (w_r, w_g, w_b, w_{g2}, u_s, ..., t_0, t_1) \qquad (3)$$

Where $(w_r, w_g, w_b, w_{g2})$ are the white balace indices of 4 raw image channels and $u_s$ is the ISO value (controlling the camera sensitivity).

In BP procedure the $t_1$ in IEV is removed. It is renamed as partial IEV (pIEV) to indicate the difference:

$$\mathcal{V}_p^{t_0} = \mathcal{V}_p(t_0) = (w_r, w_g, w_b, w_{g2}, u_s, ..., t_0) \qquad (4)$$

IEV and pIEV are fed into the convolutional neural network as additional channels of the raw image array. For example, the first extra channel provided by IEV is a uniform image filled up with pixel value $w_r$.

## 4.3 Exposure Shifting

Exposure Shifting (ES) is the second procedure in our model but is trained first. The U-NET[42] is selected as the structure of the Exposure Shifting Network (ESN), Compared with the fully convolutional network, the U-NET architecture reduced the usage of GPU memory, thus it can easily process images with high resolution.

Recall that the basic unit of our Campus Image Dataset (CID) is an image group, which consists of 8 raw-format images with different exposure time but captured in the same scene. Let $(\mathcal{X}_R, \mathcal{Y}_R)$ be the paired raw image extracted from one group, in which $\mathcal{Y}_R$ is the pre-selected ground-truth image and $\mathcal{X}_R$ is a randomly-selected low-light image. More specifically, $\mathcal{X}_R$ is randomly chosen within this image group, with $\mathcal{Y}_R$ and over-exposed images excluded (Fig.3). Moreover, the exposure time of $\mathcal{X}_R$ is denoted as $t_0$. The corresponding RGB-format image of $\mathcal{Y}_R$ is denoted as $\mathcal{Y}$ and its exposure time as $t_g$.

It is expected for the ESN to learn the exposure-shifting operator $ES$ from a large number of paired images. Specifically, the ESN is required to estimate an image $\hat{\mathcal{Y}}$ that resembles the $\mathcal{Y}$ as closely as possible. In this way ESN successfully changes the raw-format low-light image $\mathcal{X}_R$ into an RGB-format longer-exposed version of itself. More importantly, the serious noise of $\mathcal{X}_R$ can also be removed. This procedure can be written as:

$$\hat{\mathcal{Y}}_{t_g}^{\Theta_1} = F_{ES}\left(\mathcal{X}_R, \mathcal{V}^{t_0 \to t_g} \middle| \Theta_1\right)$$
$$\mathcal{V}^{t_0 \to t_g} = \mathcal{V}(t_0, t_1 = t_g) \qquad (5)$$

Where ESN is denoted as $F_{ES}$, its parameters as $\Theta_1$. Note that since the BP procedure has not yet involved, the guideline exposure time $t_1$ is replace by $t_g$ in IEV $\mathcal{V}^{t_0 \to t_g}$. The enhanced result image is represented by $\hat{\mathcal{Y}}_{t_g}^{\Theta_1}$.
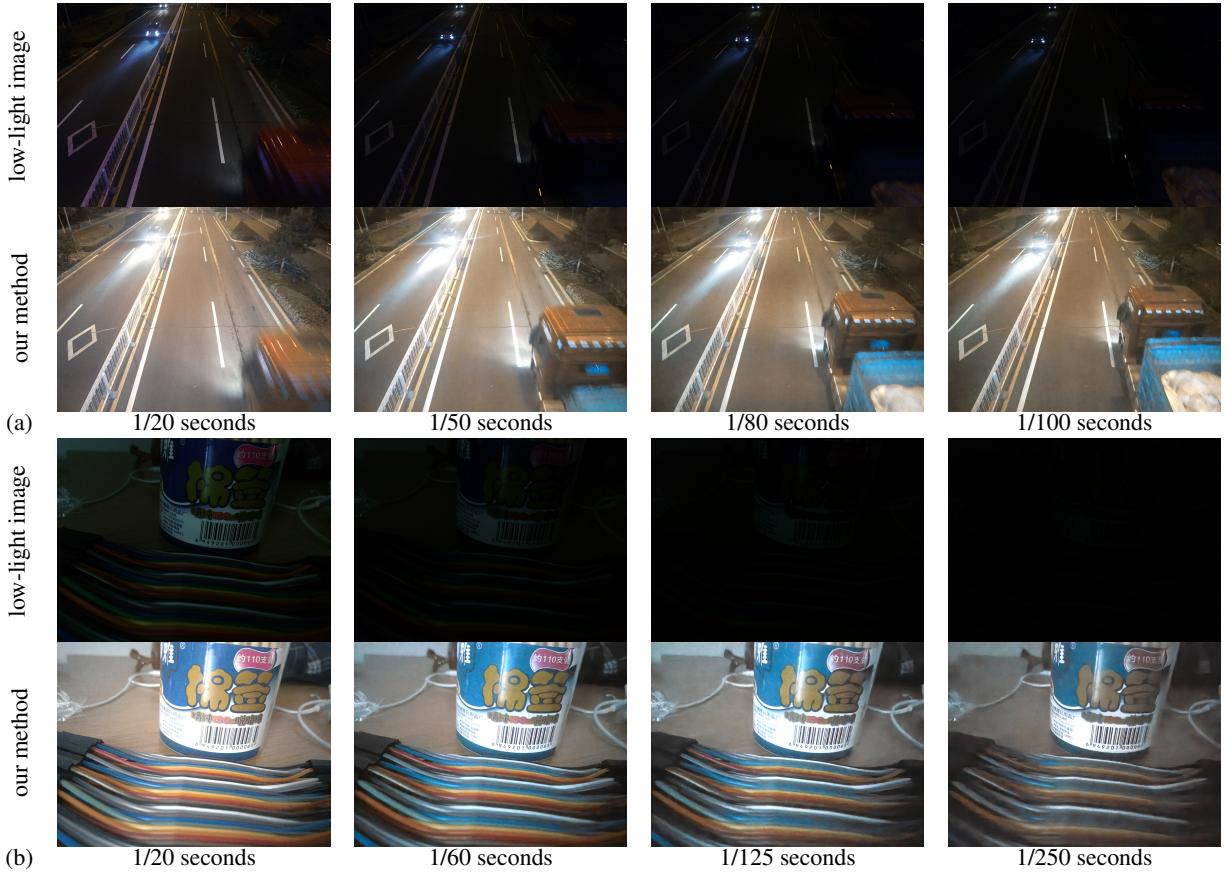
Two metrics to guide the training process of the ESN model are Mean Square Error (MAE) and multi-scale Structural Similarity (SSIM)[43]. The former is considered as a simple but effective indicator, evaluating the general similarity of the estimated image and the groud-truth image. And the later is a perceptual metric that is more sensitive to visible structures. Both of them are full-reference metrics. Let $L_{MAE}$ be the MAE loss and $L_{SSIM}$ be the SSIM loss. The width and length of the image are denoted as $m$ and $n$ respectively. With subscript $t_g$ omitted, the MAE loss $L_{MAE}$ and SSIM loss $L_{SSIM}$ are defined as:

$$L_{MAE}\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right) = MAE\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right) = \frac{1}{mn}\sum_{i=1}^{m}\sum_{j=1}^{n}\left|\hat{\mathcal{Y}}_{ij}^{\Theta_1} - \mathcal{Y}_{ij}\right|$$

$$L_{SSIM}\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right) = 1 - SSIM\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right) \qquad (6)$$

The final MAE loss is the average of the channel MAE loss, for simplicity, the image channel is not presented in the equation. The calculation of the multi-scale structural similarity (SSIM) can be refered in [43]. With subscript $t_g$ omitted, the loss function is defined as the linear combination of $L_{MAE}$ and $L_{SSIM}$:

$$L_{ES} = L_{ES}\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right)$$
$$= (1 - \alpha) L_{MAE}\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right) + \alpha L_{SSIM}\left(\hat{\mathcal{Y}}^{\Theta_1}, \mathcal{Y}\right) \qquad (7)$$

Where $\alpha$ is a constant and $0 \leqslant \alpha < 1$. The influence of $\alpha$ will be discussed in the *Experiments* section.

**Fig. 5**: The results of using our method to enhance multiple images in a same scene but captured with diverse exposure time. For the second image group, The exposure time $t_0$ of the last image in the sequence is only 8% of that of the first, yet the network compensates for this difference, keeping the brightness of each image in the scene approximately uniform.

The parameters of the ESN network $\Theta_1$ are learned by minimizing $L_{ES}$ over $K$ pairs of images in the training set:

$$\Theta_1^* = \underset{\Theta_1}{\arg\min} \sum_{k=1}^{K} L_{ES}\left(\hat{\mathcal{Y}}_k^{\Theta_1}, \mathcal{Y}_k\right) \qquad (8)$$

$\Theta_1^*$ represents the optimized ESN parameters.

The obtained sub-model $F_{ES}$ is an approximation to the exposure-shifting operator $ES$. However, for a low-light image out of the training set, a ground-truth counterpart does not exist, thus the guideline time $t_1$ in IEV $\mathcal{V}^{t_0 \to t_1}$ is absent, which leads to the Brightness Prediction (BP) sub-model and guideline time estimation.

### 4.4 Brightness Prediction

The Brightness Prediction (BP) procedure provides the estimation of guideline exposure time $t_1$:

$$t_1^{\Theta_2} = F_{BP}\left(\mathcal{X}_R, \mathcal{V}_p^{t_0} \middle| \Theta_2\right) \qquad (9)$$

Where $F_{BP}$ and $\Theta_2$ is BPN and its parameters respectively. The pIEV is represented by $\mathcal{V}_p^{t_0} = \mathcal{V}_p(t_0)$. It is the first enhancement procedure in our model but trained after the ESN. As it is illustrated in Fig.4, The Brightness Prediction Network (BPN) is a relatively small network, which consists of multiple convolution layers and then ends up with several fully connected layers. Because of the existence of fully connected layers, the width and length of the input image $\mathcal{X}_R$ should be uniformly changed into $512 \times 512$.

With the estimated $t_1^{\Theta_2}$ obtained in this BP procedure, the $t_g$ item in IEV in the ES procedure can be replaced by $t_1^{\Theta_2}$. Different

from $\hat{\mathcal{Y}}_{t_g}^{\Theta_1}$ derived in Equation (5), the new final result image can be formulated as:

$$\hat{\mathcal{Y}}^{\Theta_2|\Theta_1^*} = F_{ES}\left(\mathcal{X}_R, \mathcal{V}^{t_0 \to t_1^{\Theta_2}} \middle| \Theta_1^*\right) \qquad (10)$$
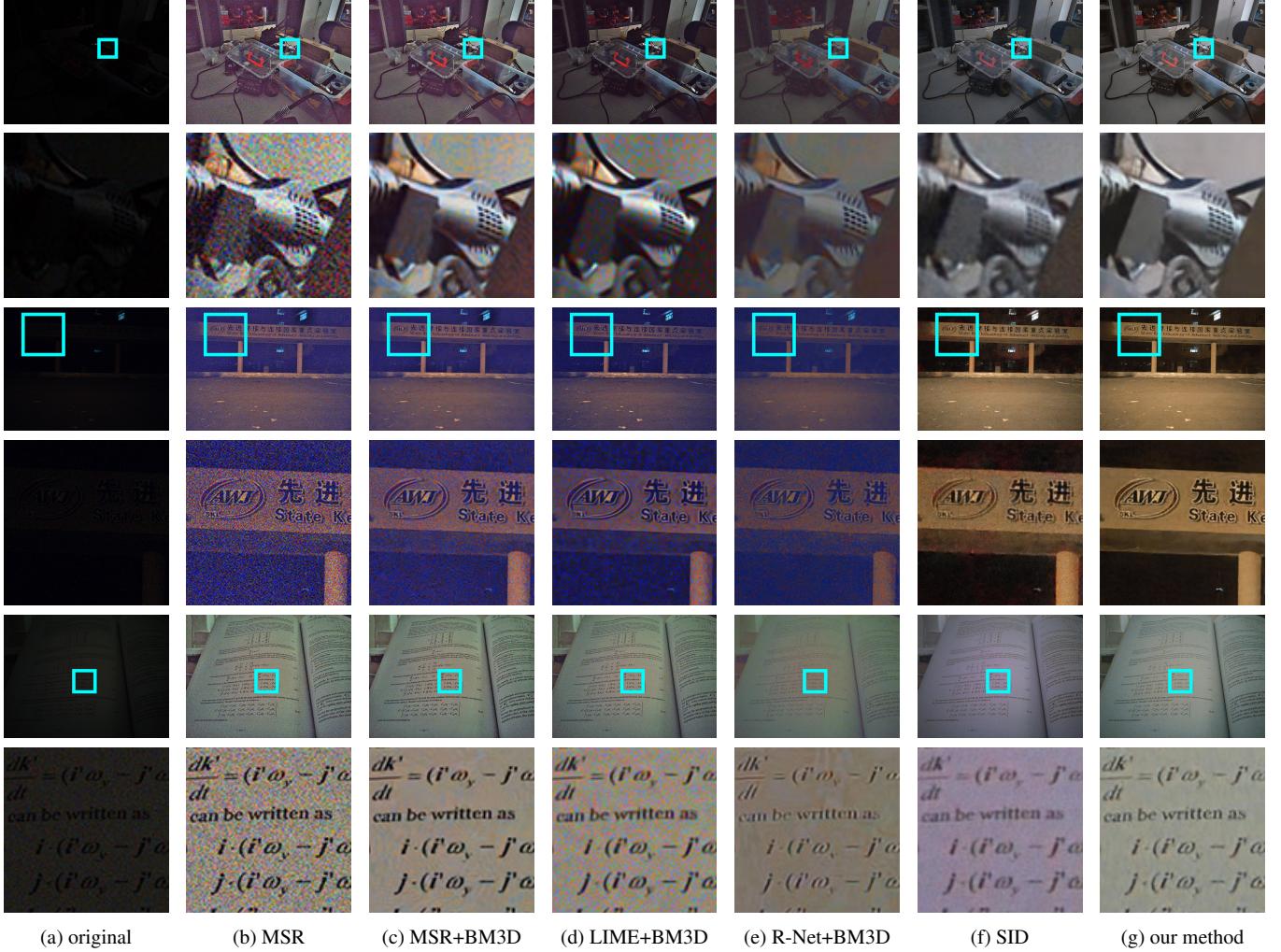
Then an approach to train the BPN is required. To begin with, the purpose of BPN is to control the brightness of the final result image, which is achieved by adjusting the guiding exposure time $t_1$ in the ES procedure. Since there are image pairs $(\mathcal{X}_R, \mathcal{Y}_R)$ and their exposure time $(t_0, t_g)$ available, intuitively it seems that the BPN can simply be learned from $(\mathcal{X}_R, t_g)$ in an end-to-end way, with $\mathcal{X}_R$ as the input and $t_g$ as the label. However it cannot be achieved, the reason for which will be discussed later in the next subsection.

The proposed BPN training approach and the loss function $L_{BP}$ are as follows. We define the brightness of a single pixel $Br$ as the average pixel value of all color channels. Given a certain dark scene for image capturing and with ISO, aperture settings and the scene all fixed, then the brightness of a pixel is only determined by exposure time. Let $\mathcal{R}$ be the pixel irradiance in this scene, when the exposure time $t$ in camera settings increases, the pixel Exposure $\mathcal{E}(\mathcal{R}, t)$ increases as well [44]:

$$\mathcal{E}(\mathcal{R}, t) = \mathcal{R}t \qquad (11)$$

Thus the pixel becomes brighter. It is feasible to design loss function $L_{BP}$ based on the pixel brightness.

It should be noted that the relationship between pixel brightness $Br(\mathcal{E})$ and Exposure $\mathcal{E} = \mathcal{R}t$ is not linear, instead, it can be typically described by an S-shaped curve because the pixel brightness is limited within $[0, 1]$ (or $[0, 255]$ for 8-bit image storage) and the camera sensor is less sensitive to the change of $t$ when the pixel value is close to 0 or 1. In other words, as the exposure time $t$

| (a) original | (b) MSR | (c) MSR+BM3D | (d) LIME+BM3D | (e) R-Net+BM3D | (f) SID | (g) our method |

**Fig. 6**: The results using different methods on CID test images. From (a) to (g), they are the original low-light images and the results of multiscale Retinex (MSR), MSR with BM3D, illumination map estimation based (LIME), LIME with BM3D, deep Retinex decomposition (R-net), R-net with BM3D, learning-to-see-in-the-dark (SID) model and ours respectively.

changes, some pixels are not sensitive and there are little changes in the pixel brightness compared with some other pixels. For example, when photographing a street at night, pixels revealing the lamp bulbs are always saturated unless the $t$ is extremely small, and pixels representing the dark sky are always close to zero unless it is affected by noise. If those pixels are involved in BPN training, they can cause gradient vanishing in the back propagation algorithm because $\frac{\partial Br(\mathcal{E})}{\partial t} = \frac{\partial Br(\mathcal{E})}{\partial \mathcal{E}} \cdot \mathcal{R}$ and as it is revealed in the S-shaped culve, $\frac{\partial Br(\mathcal{E})}{\partial \mathcal{E}}$ is close to 0. On the other hand, when pixel brightness $Br$ is around 0.5 the $\frac{\partial Br}{\partial t}$ reaches its maximum. The pixels meeting this condition can be collected into an Area of Interest (AoI), which identifies pixels that are sensitive to the process of Exposure Shifting.

First, in the ground-truth image $\mathcal{Y}$, an AoI weight map $\mathcal{W}$ can be obtained by filtering out pixels with brightness near 0.5. More specifically, the ground-truth image $\mathcal{Y}$ is converted into a single channel grayscale image $\mathcal{Y}_G$, then a Gaussian culve with mean $\mu_w = 0.5$ and variance $\sigma_w^2 = 0.01$ is applied on each pixel of $\mathcal{Y}_G$:

$$\mathcal{W}_G = \exp\left(-\frac{(\mathcal{Y}_G - \mu_w)^2}{2\sigma_w^2}\right) \quad (12)$$

Then $\mathcal{W}_G$ is normalized, the result is the AoI weight map $\mathcal{W}$. The value at position $(i, j)$ is denoted as subscripts, namely $\mathcal{W}_{G,ij}$ and $\mathcal{W}_{ij}$:

$$\mathcal{W}_{ij} = \frac{\mathcal{W}_{G,ij}}{\sum_{i=1}^{m}\sum_{j=1}^{n}\mathcal{W}_{G,ij}}, \forall 1 \leqslant i \leqslant m, 1 \leqslant j \leqslant n \quad (13)$$

Combining the Equation 12 and 13, they can be simplified by the softmax function:

$$\mathcal{W} = softmax\left(-\frac{(\mathcal{Y}_G - \mu_w)^2}{2\sigma_w^2}\right) \quad (14)$$

Lastly, the loss function $L_{BP}$ is designed to check if the final result image $\hat{\mathcal{Y}}^{\Theta_2|\Theta_1^*}$ has the same AoI. Having the same AoI means that in this AoI area $\hat{\mathcal{Y}}^{\Theta_2|\Theta_1^*}$ and $\mathcal{Y}$ share similar brightness. Moreover, recall that pixels outside AoI are not sensitive to the change of exposure time $t$ or the ES process, thus it can be derived that the overall image brightness of $\hat{\mathcal{Y}}^{\Theta_2|\Theta_1^*}$ resembles that of $\mathcal{Y}$. In this way, BPN completed the adaptive Brightness Prediction objective.

Then the enhanced image $\hat{\mathcal{Y}}^{\Theta_2|\Theta_1^*}$ is converted into a grayscale image $\hat{\mathcal{Y}}_G^{\Theta_2|\Theta_1^*}$. When the AoI of $\hat{\mathcal{Y}}^{\Theta_2|\Theta_1^*}$ and $\mathcal{Y}$ overlaps exactly,

the following loss function reaches its minimum:

$$L_{BP} = L_{BP}\left(\hat{\mathcal{Y}}_G^{\Theta_2|\Theta_1^*}, \mathcal{W}\Big|\Theta_2; \Theta_1^*\right)$$

$$= -\frac{1}{mn}\sum_{i=1}^{m}\sum_{j=1}^{n}\exp\left(-\frac{\left(\hat{\mathcal{Y}}_{G,ij}^{\Theta_2|\Theta_1^*} - \mu_w\right)^2}{2\sigma_v^2}\right)\mathcal{W}_{ij} \quad (15)$$

Where $\sigma_v^2 = 0.04$ is another variance constant. The image value at position $(i, j)$ is represented by $i, j$ subscripts. Finally, BPN is trained by optimizing $\Theta_2$ over $K$ pairs of images in the training set:

$$\Theta_2^* = \underset{\Theta_2}{\text{argmin}}\sum_{k=1}^{K}L_{BP}\left(\hat{\mathcal{Y}}_{G,k}^{\Theta_2|\Theta_1^*}, \mathcal{W}_k\Big|\Theta_2; \Theta_1^*\right) \quad (16)$$

### 4.5 Method Summary and Discussion

In Algorithm 1, the training process of our model is summarized, together with the method to evaluate and apply this model.

---

**Algorithm 1** Model training and evaluation

---

1: Initialize ESN, BPN parameters $\Theta_1$ and $\Theta_2$
2: Prepare CID dataset
3: **function** TRAIN(Dataset)
4:     **for** $EpochTrainESN = 1 \rightarrow E_1$ **do**
5:         **for** Image Pair $k = 1 \rightarrow K$ **do**
6:             Read raw image pairs $(\mathcal{X}_{R,k}, \mathcal{Y}_{R,k})$ and $\mathcal{V}_k^{t_0 \rightarrow t_g}$
7:             Convert to RGB format. $\mathcal{Y}_k \leftarrow \mathcal{Y}_{R,k}$
8:             Compute $\hat{\mathcal{Y}}_k^{\Theta_1}$ via Equation 5
9:             $\Theta_1 \leftarrow \text{argmin}_{\Theta_1}\sum_{k=1}^{K}L_{ES}\left(\hat{\mathcal{Y}}_k^{\Theta_1}, \mathcal{Y}_k\right)$
10:         **end for**
11:     **end for**
12:     $\Theta_1^* \leftarrow \Theta_1$
13:     **for** $EpochTrainBPN = 1 \rightarrow E_2$ **do**
14:         **for** Image Pair $k = 1 \rightarrow K$ **do**
15:             Read raw image pair $(\mathcal{X}_{R,k}, \mathcal{Y}_{R,k})$ and $\mathcal{V}_{p,k}^{t_0}$
16:             Compute $t_{1,k}^{\Theta_2}$ and $\hat{\mathcal{Y}}_k^{\Theta_2|\Theta_1^*}$ via Equation 9 and 10
17:             Convert to grayscale image. $\mathcal{Y}_{G,k} \leftarrow \mathcal{Y}_{R,k}$
18:             Convert to grayscale image. $\hat{\mathcal{Y}}_{G,k}^{\Theta_2|\Theta_1^*} \leftarrow \hat{\mathcal{Y}}_k^{\Theta_2|\Theta_1^*}$
19:             Compute $\mathcal{W}_k$ via Equation 14
20:             $\Theta_2 \leftarrow \text{argmin}_{\Theta_2}\sum_{k=1}^{K}L_{BP}\left(\hat{\mathcal{Y}}_{G,k}^{\Theta_2|\Theta_1^*}, \mathcal{W}_k\Big|\Theta_2; \Theta_1^*\right)$
21:         **end for**
22:     **end for**
23:     $\Theta_2^* \leftarrow \Theta_2$
24:     **return** $\Theta_1^*, \Theta_2^*$
25: **end function**
26:
27: **function** EVALUATE(Image)
28:     Read raw-format image $\mathcal{X}_R$
29:     Read all elements in $\mathcal{V}_p^{t_0}$ from raw-image file
30:     BPN: Compute $t_1^{\Theta_2^*} \leftarrow F_{BP}\left(\mathcal{X}_R, \mathcal{V}_p^{t_0}\Big|\Theta_2^*\right)$
31:     ESN: Compute $\hat{\mathcal{Y}}^{\Theta_2^*|\Theta_1^*} \leftarrow F_{ES}\left(\mathcal{X}_R, \mathcal{V}^{t_0 \rightarrow t_1^{\Theta_2^*}}\Big|\Theta_1^*\right)$
32:     **return** $\hat{\mathcal{Y}}^{\Theta_2^*|\Theta_1^*}$
33: **end function**

---

As mentioned earlier in the last section, training BPN with input $\mathcal{X}_R$ and label $t_g$ is straightforward but not feasible. First let $F_{BP}^*$ be the hypothetical optimal BPN model, then $t_1^*(\mathcal{X}_R) = F_{BP}^*(\mathcal{X}_R)$ is the optimal estimation of guideline exposure time. For a pair of image and time label $(\mathcal{X}_R, t_g)$, the $t_g$ can be considered as one sampling from a Gaussian distribution: $t_g \sim \mathcal{N}\left(t_1^*, \sigma_c^2\right)$. The variance $\sigma_c^2$ in this process is too large, the reason are:

**Table 1** Comparison results with non-reference image quality metrics. MSR, LIME, R-Net methods are combined with BM3D to make a fair comparison. No-reference quality metrics, including Statistical Noise Level Estimation (SNLE), NLE(Noise Level Estimation), NV(Noise Variance) and image entropy, are applied to evaluate model denoising and brightness control performance.

| Methods | SNLE | NLE | NV | Entropy |
|---------|------|-----|-----|---------|
| Original | - | - | - | 4.0689 |
| MSR+BM3D | **0.0014** | 1.3792 | 0.8998 | 7.0068 |
| LIME+BM3D | 0.0086 | 4.2437 | 10.1789 | **7.1084** |
| R-Net+BM3D | 0.0017 | 0.2467 | **0.4491** | 6.1503 |
| Ours | **0.0015** | 0.1061 | 0.4554 | **6.6112** |

- In every CID image group, the ground-truth image is chosen manually and is dominated by subjective human judgments.
- The ground-truth image is selected from only 8 candidate images, but the optimal $t_1^*$ exists somewhere in the time continuum $(0, \inf)$.

Furthermore, the number of samples is very limited since each image group can only provide one $t_g$ sample. As a result, this straightforward approach cannot be applied in the training of BPN due to the high variance of the label $t_g$ and the limited number of samples.
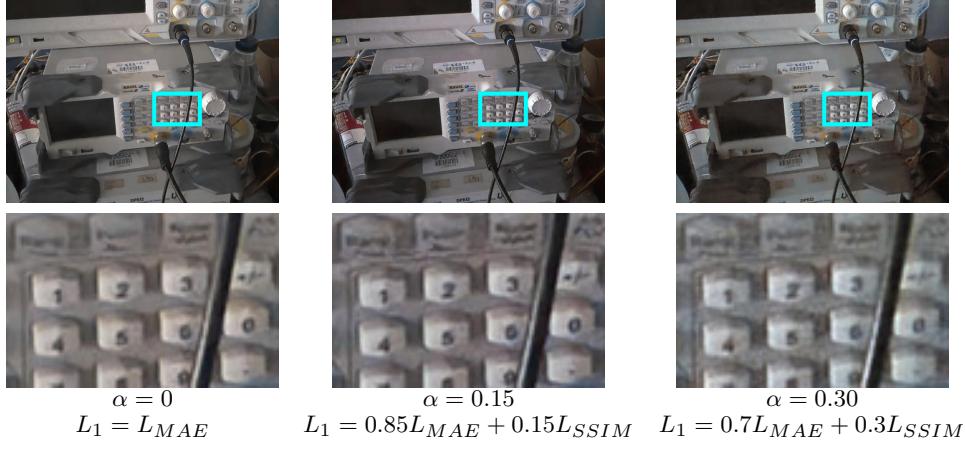
## 5 Experiments

### 5.1 Training

Limited by the GPU memory, we did not adopt the Batch Normalization [45] technique. To accelerate model training, before the training process started, each channel of the input image, as well as each element of IEV, was normalized along the sample dimension. Then the means and variances used in the normalization became invariant constants. Additionally, images for training were randomly cropped into $512 \times 512$ patches.

Both ESN and BPN were trained with Adam optimizer[46], following the procedures shown in Algorithm 1. The ESN was trained first with the learning rate slowly descending from $2 \times 10^{-4}$ to $1 \times 10^{-5}$. After about 300 epochs, when the loss on training sets became stable, the parameters in ESN were made untrainable. Then BPN was appended to the training graph and trained with the same learning rate. Considering BPN may require global information to decide the proper exposure time, we used bigger patches but avoid training with the whole image to prevent overfitting. Although the loss must propagate through the deep ESN before reaching BPN, the training process went smoothly and completed after 100 epochs.
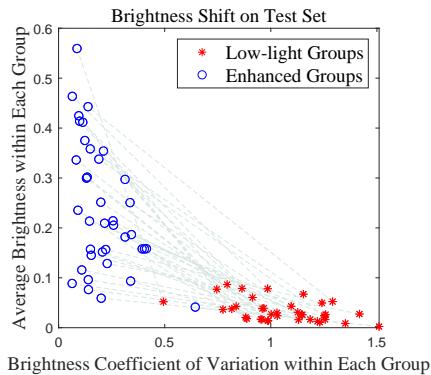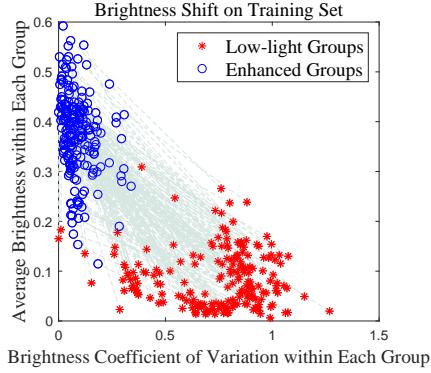
### 5.2 Content-Based Brightness Control

Firstly, we highlight the adaptive brightness control feature of our method. Here we explicitly define the brightness as the mean value of an image $\mathcal{Y}$, denoted as $Br(\mathcal{Y})$. We expect a dynamic adjustment of the image brightness based on the image content. More specifically, while extreme low-light images suffer from severe noise along with color distortion and are barely visible before post-processing, other mild low-light images only have moderate noise and relatively lower brightness compared with normal images. Thus, the enhancement algorithm needs to make adjustments adaptively for different input images. Currently the classic methods, such as algorithms based on Retinex and Dehazing, have poor performance on the extreme low-light images, On the other hand, learning-based methods targeting at end-to-end denoising can process all kinds of low-light images, but they lack the flexibility and adaptability because they either need brightness amplification hand-tuning, or simply serve as a denoising procedure in other low-light methods. The proposed model has both the advantages and is able to process low-light images of different levels.

In order to evaluate the adaptability, we apply the trained model to enhance all images in the test set. Recall that the image group is the basic unit of our CID dataset, accordingly, instead of evaluating single images, the collaborative characteristics across different images
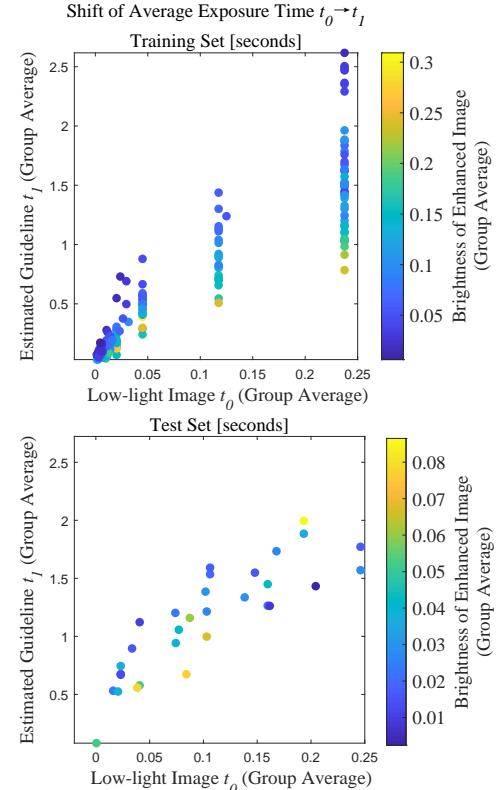
$$\alpha = 0 \qquad \alpha = 0.15 \qquad \alpha = 0.30$$
$$L_1 = L_{MAE} \qquad L_1 = 0.85L_{MAE} + 0.15L_{SSIM} \qquad L_1 = 0.7L_{MAE} + 0.3L_{SSIM}$$

**Fig. 7**: Controlled experiments on the impact of SSIM loss.



**Fig. 8**: The shift of average brightness and brightness coefficient of variation (CV) within each group. Each Group contain 3–6 images captured in a brust with different exposure time, e.g. Fig.5. In both test set and training set, an increase in brightness and reduction in CV is observed. Note that the CV decreased considerably, indicating the highly uniformed brightness for images within each group after enhancement.



**Fig. 9**: The comparison of image original exposure time $t_0$ and BPN assigned time $t_1$. BPN infers $t_1$ based on not only the original time $t_0$ but also the content of the scene. The color indicates the mean global brightness value of enhanced images within each group.

are investigated within each image group. Note that the group identity of an image is NOT involved in model training and only serves as a tag to gather resulting images for further analysis.
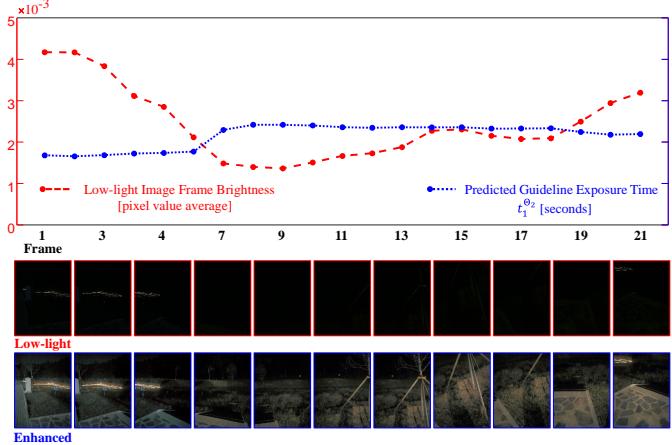
For a sequence of multi-exposed images $(\mathcal{X}_{R,1}, \mathcal{X}_{R,2}, ...)$ with increasing exposure time, it is expected that after enhancement the resulting images possess similar image brightness, namely $Br(\hat{\mathcal{y}}_1) \approx Br(\hat{\mathcal{y}}_2) \approx ...$, where $\hat{\mathcal{y}}$ denotes the enhanced image. Fig.5(a) illustrates a busy-road scene, the first row displays the RGB images produced by the camera, the second displays the enhanced results by utilizing our method on individual raw images. Note that the first image in this group is exposed for 1/20 seconds, causing motion blur to the truck, while the last image is exposed for only

1/100 seconds. It is shown that all four images have approximately uniform brightness after the image enhancement, despite the changing $t_0$ and the moving vehicle headlight. The performance of our method on extreme low-light conditions is shown in Fig.5(b). The last image in the sequence suffers from both low environment illumination and short exposure time and the color and shape of objects in the image are drowned out by noise. However, the resulting image has the noise reduced and still maintain acceptable global brightness.

As mentioned above, for quantitative analysis, we gathered images that belong to the same groups and calculated the brightness of all images in order to observe the shift of brightness. Let $(Br_1, Br_2, ...)_k$ be the image brightness of each unenhanced low-light image in $k$-th image group and let $(Br'_1, Br'_2, ...)_k$ be that of the enhanced ones. Then, the mean value and the CV (coefficient of variation) of $(Br_1, Br_2, ...)$ is computed and denoted as $(\mu, c_v)_k$.

**Fig. 10**: Applying our method to dynamic low-light video enhancement.

Correspondingly, $(\mu', c'_v)_k$ is calculated from $(Br'_1, Br'_2, ...)_k$. In Fig.8, we plotted $(\mu, c_v)$ and $(\mu', c'_v)$ for all groups. To illustrate the change after enhancement, for each group $(\mu, c_v)$ (marked with star) and $(\mu', c'_v)$ (marked with circle) are linked with a dotted line. It can be observed from Fig.8 that for different groups that captured from diverse scenes, the enhanced images all end up with increased image brightness and significantly reduced CV. And yet the brightness growth among the groups can be very different, ranging from around 60% to over 1,000%. Considering together with Fig.5, Fig.8 indicates that the BP sub-model is able to adaptively control the global image brightness not only based on the brightness of the original image, but also according to the content of the scene. On the other hand, the dramatic reduction of CV suggests that the enhanced images within each group share almost identical brightness in spite of the varying exposure time $t_0$. To be precise, our model can control group brightness CV under 0.29 on 76.3% test set groups and 98.5% training set groups considering the influence of the diverse scenes, illumination conditions, ISO etc.

When designing our model, we expect the BPN to extract features of the scene content and then estimates a reasonable $t_1$ accordingly. In order to see through clearly how this black box works, the relationships among the real exposure time $t_0$, the BPN-estimated guideline exposure time $t_1$ and the brightness of enhanced images $Br(\hat{\mathcal{Y}})$ are shown in Fig.9. Similar to Fig.8, we use group average to simplify the figure, in which each point corresponds to an image group. As shown in Fig.9, the network tends to estimate $t_1$ positively associated with $t_0$. Furthermore, for groups with identical average $t_0$, the contents of the scenes have an impact on $t_1$ prediction. To be specific, if the scene of a group has more relatively bright areas, a smaller $t_1$ is more likely to be adopted to prevent overexposure and vice versa. And yet a few exceptions exist, suggesting the BPN sub-model has learned some more sophisticated rules in the neural network black box.

## 5.3    Denoising Evaluation

We compare our method with four classic and state-of-art methods, including multiscale Retinex (MSR)[5], illumination map estimation based (LIME)[7], deep Retinex decomposition (Retinex-Net, R-Net)[9] and learning-to-see-in-the-dark (SID)[15]. For fair comparison, 3D transform-domain filtering (BM3D)[10] is applied on MSR, LIME, R-Net, since those methods do not model noise and require extra denoising process. Additionally, we train the SID model on the CID dataset and manually select ratio $\gamma$ for the SID model since it does not provide automatic brightness estimation.

Recall that the BPN utilizes a reduced reference evaluation index as the loss in its training. In this way, the exposure parameters of the ground-truth image, which is chosen by highly subjective human judgment, will have much less influence on BPN model training. As a result of this novel technique, however, PSNR and SSIM evaluation

cannot be adopted because there is no reference image. We provide perceptual comparisons in Fig.6 and evaluate our method with no-reference quality metrics.

In Table.1, we compare the proposed method with MSR, LIME and R-Net(Retinex-Net), using several no-reference quality metrics including Statistical Noise Level Estimation (SNLE) [47], Noise Level Estimation (NLE) [48], Noise Variance (NV) [49] and image entropy. SNLE, NLE and NV metrics can estimate noise level from a single image, and clean images have relatively smaller values than the noisy ones. The SNLE estimates noise variance by observing the eigenvalues of images and the NLE by applying the principal component analysis (PCA) on special image patches. The NV metric is motivated by the fact that clean and noisy images have different sensitivity to the Laplacian operation. The image entropy reflects the average amount of information in an image, it is the most simple image evaluation metric. Images with proper brightness have more visible information, thus have larger image entropy compared to under-exposed or over-exposed ones. We run evaluations on about 400 CID test images and the SID method is excluded because it requires an external manually-selected amplification ratio for every image. As it is shown in Table.1, our method has superior performance on Noise Level Estimation. Moreover, the no-reference metrics use very different standards to assess image quality. MSR+BM3D, R-Net+BM3D and LIME+BM3D methods all have undesired scores on one or multiple indices, while our method has good performance on all metrics.

It is illustrated in Fig.6 that our method and SID method provide more significant improvements in noise removal than the classic method BM3D and can adaptively denoise images with diverse noise levels. Furthermore, our results benefit from the white balance index introduced in IEV and avoid abnormal color in extreme low-light images.

However, the CNN-based method tends to blur the object with sharp edges (e.g. text). By introducing SSIM into loss function, our method shows more promising results in images with text and has advantages over SID in keeping more edge information, e.g. the images in the first row of Fig.6. But on the other hand, in controlled experiments shown in Fig.7, it is revealed that the ratio of SSIM component $\alpha$ in loss function $L_{ES}$ has to be contained to avoid the noise being interpreted into texture and edge by ESN.

## 5.4    Low-light Video Processing

The advantage of our work lies in getting rid of the brightness control ratio while achieving state-of-the-art low-light noise suppression. Without any handcrafting parameters as network input, we can directly apply our method to raw low-light video processing.

A raw-format video can occupy a very large amount of storage space. Limited by devices, we only developed a relatively simple way to verify our method on raw low-light videos. In Fig.10, we experiment on 21 continuous raw images captured in a burst session, shot with short exposure time and identical ISO, to discuss the potential applications in raw video processing. In this scene, the camera starts from a bridge with lighting decoration, then turns to the dark riverbank, and finally to the dimly illuminated sidewalk. As is shown in Fig.10, though the enhancement process of each frame is completely independent of other frames, neither the estimated $t_1$ nor the resulting images show any sudden fluctuation. However, for practical considerations, it is also convenient to implement guideline time filters between ESN and BPN, which allows the value of $t_1$ to be dynamically filtered or amplified according to different demands.

## 5.5    Computational Cost

The proposed algorithm is evaluated on the ModelArts cloud service, which is equipped with 8 vCPUs and an Nvidia-P100 GPU (16GB GPU memory). The execution time for a high-resolution image ($3968 \times 2976$) is 0.27 seconds on average, which only shows a minor increase compared with the SID algorithm (0.21 seconds). More specifically, the execution time of BPN and ESN sub-model is 0.05 and 0.22 respectively. For comparison, it takes 4.32 seconds

for BM3D (GPU implementation [50]) to complete the denoising process, therefore other low-light methods that depend on BM3D for denoising, including MSR, LIME, R-Net et al., are much less efficient on devices with GPUs.

## 6 Conclusion and Discussion

### 6.1 Conclusion

In this paper, an adaptive raw image enhancement model is proposed for extreme low-light image processing. The two-stage framework provides adaptability to images with different scenes and exposure parameters. We also presented the CID dataset for model training and evaluation. The experimental test shows that our model can provide the state-of-the-art low-light image denoising as well as adaptive global brightness control.

The proposed method has many advantages over existing approaches. In the first place, no external parameters are needed after model training, a single raw image and its Exif metadata (e.g. white balance, ISO and exposure time) are sufficient to complete the process. Therefore, our model is able to process a large batch of raw images without parameter-handcrafting. Secondly, our model significantly outperforms other methods in which the denoising stage works as a separable step. Because instead of simply implementing noise-suppression algorithms (e.g. BM3D) before or after the main low-light enhancement procedure, we make denoising procedure an integrated part in our model, allowing the ESN module to learn exposure shifting and denoising simultaneously in an end-to-end way. In quantitative experiments, we adopted no-reference image quality metrics including SNLE, NLE, NV and image entropy to test the denoising performance. Our model has the smallest noise variance on the NLE metric compared to BM3D-based low-light methods and obtains near-optimal scores on the other three indices.

We also reveal the potential application in low-light video processing. The BPN module makes it possible for our model to process images with diverse exposure levels, from extreme low-light images to mild ones. But to be applied in video processing, it should face more challenges, e.g. the brightness of frames should change continuously and avoid fluctuation. We experiment on continuous raw image series, despite that no information of adjacent frames is used, the output of BPN changes continuously with the illumination condition and shows no fluctuation.

### 6.2 Limitations and Future Works

The model proposed is limited by dataset, which is a common problem in data-driven methods. On the one hand, the model is trained by raw-format and multi-exposed images, which significantly improved image quality. But on the other hand, all reference ground-truth images are captured in the real scenes, as a result, they sometimes have minor defects such as noise, halos around the light source, local over-exposure and white balance deviation. The denoising performance may be further improved if those drawbacks can be avoided. Then the image groups in CID have not covered enough nature scenes. For example, it is found that green grass and trees tend to have low color saturation in the resulting images, that is because most images in the training set are collected in winter when there is hardly any green plants in scenes. For future study, the raw-format CID dataset can be extended by unprocessing [14] technique and generating raw-format images from other available online datasets.

Another limitation is the memory consumption problem. The batch size is constrained because it takes too much GPU memory to train the proposed model. Therefore training the model can be time-consuming on GPUs with small memory. We are planning to improve network architecture and investigate the possibility of application on mobile devices.

As a future expectation, the model can be further modified by introducing HDRI (High Dynamic Range Imaging). Since our CID dataset is composed of multi-exposed images, we can calculate an HDR image of each scene. Finally, the HDR images can be used as ground truth to participate in the ESN training.

## 7 References

1 Pizer, S.M., Amburn, E.P., Austin, J.D., Cromartie, R., Geselowitz, A., Greer, T., et al.: 'Adaptive histogram equalization and its variations', *Computer vision, graphics, and image processing*, 1987, **39**, (3), pp. 355–368

2 Huang, S.C., Cheng, F.C., Chiu, Y.S.: 'Efficient contrast enhancement using adaptive gamma correction with weighting distribution', *IEEE Trans. on image processing*, 2012, **22**, (3), pp. 1032–1041

3 Land, E.H.: 'The retinex', *American Scientist*, 1964, **52**, (2), pp. 247–264

4 Jobson, D.J., Rahman, Z., Woodell, G.A.: 'Properties and performance of a center/surround retinex', *IEEE Trans. on Image Processing*, 1997, **6**, (3), pp. 451–462

5 Jobson, D.J., Rahman, Z.u., Woodell, G.A.: 'A multiscale retinex for bridging the gap between color images and the human observation of scenes', *IEEE Trans. on Image processing*, 1997, **6**, (7), pp. 965–976

6 Fu, X., Zeng, D., Huang, Y., Liao, Y., Ding, X., Paisley, J.: 'A fusion-based enhancing method for weakly illuminated images', *Signal Processing*, 2016, **129**, pp. 82–96

7 Guo, X., Li, Y., Ling, H.: 'Lime: Low-light image enhancement via illumination map estimation.', *IEEE Trans. Image Processing*, 2017, **26**, (2), pp. 982–993

8 Dong, X., Wang, G., Pang, Y., Li, W., Wen, J., Meng, W., et al.: 'Fast efficient algorithm for enhancement of low lighting video'. IEEE Int. Conf. on Multimedia and Expo, Barcelona, Spain, 2011. pp. 1–6

9 Wei, C., Wang, W., Yang, W., Liu, J.: 'Deep retinex decomposition for low-light enhancement', *arXiv preprint arXiv:180804560*, 2018

10 Kostadin, D., Alessandro, F., Vladimir, K., Karen, E.: 'Image denoising by sparse 3-d transform-domain collaborative filtering', *IEEE Trans. on Image Processing*, 2007, **16**, (8), pp. 2080–2095

11 Burger, H.C., Schuler, C.J., Harmeling, S.: 'Image denoising: Can plain neural networks compete with bm3d?'. IEEE Conf. on computer vision and pattern recognition, Providence, RI, USA, 2012. pp. 2392–2399

12 Lore, K.G., Akintayo, A., Sarkar, S.: 'Llnet: A deep autoencoder approach to natural low-light image enhancement', *Pattern Recognition*, 2017, **61**, pp. 650–662

13 Zhang, K., Zuo, W., Chen, Y., Meng, D., Zhang, L.: 'Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising', *IEEE Trans. on Image Processing*, 2017, **26**, (7), pp. 3142–3155

14 Brooks, T., Mildenhall, B., Xue, T., Chen, J., Sharlet, D., Barron, J.T.: 'Unprocessing images for learned raw denoising'. Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, Long Beach, CA, USA, 2019. pp. 11036–11045

15 Chen, C., Chen, Q., Xu, J., Koltun, V.: 'Learning to see in the dark'. Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition, Salt Lake City, UT, USA, 2018. pp. 3291–3300

16 Finlayson, G., Hordley, S., Schaefer, G.: 'Illuminant and device invariant colour using histogram equalisation', *Pattern Recognition*, 2005, **38**, pp. 179–190

17 Pizer, S.M., Amburn, E.P., Austin, J.D., Cromartie, R., Geselowitz, A., Greer, T., et al.: 'Adaptive histogram equalization and its variations', *Graphical Models graphical Models and Image Processing computer Vision, Graphics, and Image Processing*, 1987, **39**, (3), pp. 355–368

18 Zuiderveld, K.J.: 'Contrast limited adaptive histogram equalization', *Graphics gems*, 1994

19 Naik, S.K., Murthy, C.A.: 'Hue-preserving color image enhancement without gamut problem', *IEEE Trans. on Image Processing*, 2003, **12**, (12), pp. 1591–1598

20 Wang, Y., Chen, Q., Zhang, B.: 'Image enhancement based on equal area dualistic sub-image histogram equalization method', *IEEE Trans. on Consumer Electronics*, 1999, **45**, (1), pp. 68–75

21 Kim, Y.: 'Contrast enhancement using brightness preserving bi-histogram equalization', *IEEE Trans. on Consumer Electronics*, 1997, **43**, (1), pp. 1–8

22 Li, L., Wang, R., Wang, W., Gao, W.: 'A low-light image enhancement method for both denoising and contrast enlarging', , 2015

23 Malm, H., Oskarsson, M., Warrant, E.J., Clarberg, P., Hasselgren, J., Lejdfors, C.: 'Adaptive enhancement and noise reduction in very low light-level video', , 2007

24 Jobson, D.J., Rahman, Z., Woodell, G.A.: 'Properties and performance of a center/surround retinex', *IEEE Trans. on Image Processing*, 1997, **6**, (3), pp. 451–462

25 Fu, X., Zeng, D., Huang, Y., Liao, Y., Ding, X., Paisley, J.: 'A fusion-based enhancing method for weakly illuminated images', *Signal Processing*, 2016, **129**, pp. 82–96

26 Wang, S., Zheng, J., Hu, H., Li, B.: 'Naturalness preserved enhancement algorithm for non-uniform illumination images', *IEEE Trans. on Image Processing*, 2013, **22**, (9), pp. 3538–3548

27 Kimmel, R., Elad, M., Shaked, D., Keshet, R., Sobel, I.: 'A variational framework for retinex', *Int. Journal of Computer Vision*, 2003, **52**, (1), pp. 7–23

28 Hao, S., Feng, Z., Guo, Y.: 'Low-light image enhancement with a refined illumination map', *Multimedia Tools and Applications*, 2018, **77**, (22), pp. 29639–29650

29 Li, M., Liu, J., Yang, W., Sun, X., Guo, Z.: 'Structure-revealing low-light image enhancement via robust retinex model', *IEEE Transactions on Image Processing*, 2018, **27**, (6), pp. 2828–2841

30 Gu, Z., Chen, C., Zhang, D.y.: 'A low-light image enhancement method based on image degradation model and pure pixel ratio prior', *Mathematical Problems in Engineering*, 2018, **2018**, pp. 1–19

31 Park, S., Moon, B., Ko, S., Yu, S., Paik, J.: 'Low-light image enhancement using variational optimization-based retinex model', , 2017

32 Fu, G., Duan, L., Xiao, C.: 'A hybrid l2-lp variational model for single low-light image enhancement with bright channel prior', , 2019. pp. 1925–1929

33 Zhang, K., Zuo, W., Chen, Y., Meng, D., Zhang, L.: 'Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising', *IEEE Trans. on Image Processing*, 2017, **26**, (7), pp. 3142–3155

34 Guo, S., Yan, Z., Zhang, K., Zuo, W., Zhang, L.: 'Toward convolutional blind denoising of real photographs', , 2019

35 Shen, L., Yue, Z., Feng, F., Chen, Q., Liu, S., Ma, J.: 'Msr-net: Low-light image enhancement using deep convolutional network', *arXiv preprint arXiv:171102488*,

2017

36  Cai, J., Gu, S., Zhang, L.: 'Learning a deep single image contrast enhancer from multi-exposure images', *IEEE Trans. on Image Processing*, 2018, **27**, (4), pp. 2049–2062

37  Hasinoff, S.W., Sharlet, D., Geiss, R., Adams, A., Barron, J.T., Kainz, F., et al.: 'Burst photography for high dynamic range and low-light imaging on mobile cameras', *ACM Trans. on Graphics (TOG)*, 2016, **35**, (6), pp. 192

38  Xu, J., Li, H., Liang, Z., Zhang, D., Zhang, L.: 'Real-world noisy image denoising: A new benchmark', *arXiv preprint arXiv:180402603*, 2018

39  Anaya, J., Barbu, A.: 'Renoir - a dataset for real low-light noise image reduction', *arXiv preprint arXiv:14098230*, 2014

40  Loh, Y.P., Chan, C.S.: 'Getting to know low-light images with the exclusively dark dataset', *Computer Vision and Image Understanding*, 2019, **178**, pp. 30–42

41  Reinhard, E., Ward, G., Pattanaik, S., Debevec, P.: 'High dynamic range imaging : acquisition, display, and image-based lighting'. (Morgan Kaufmann, 2005)

42  Ronneberger, O., Fischer, P., Brox, T.: 'U-net: Convolutional networks for biomedical image segmentation'. Int. Conf. on Medical image computing and computer-assisted intervention, Munich, Germany: Springer, 2015. pp. 234–241

43  Wang, Z., Simoncelli, E.P., Bovik, A.C.: 'Multiscale structural similarity for image quality assessment'. The Thrity-Seventh Asilomar Conf. on Signals, Systems & Computers, 2003, vol. 2. Pacific Grove, CA, USA: Ieee, 2003. pp. 1398–1402

44  Fu, L., Qi, Y.: 'Camera response function estimation and application with a single image'. Yang, D., editor. Informatics in Control, Automation and Robotics, Berlin, Heidelberg: Springer Berlin Heidelberg, 2012. pp. 149–156

45  Ioffe, S., Szegedy, C.: 'Batch normalization: Accelerating deep network training by reducing internal covariate shift', *arXiv preprint arXiv:150203167*, 2015

46  Kingma, D.P., Ba, J.: 'Adam: A method for stochastic optimization', *arXiv preprint arXiv:14126980*, 2014

47  Chen, G., Zhu, F., Heng, P.A.: 'An efficient statistical method for image noise level estimation'. 2015 IEEE Int. Conf. on Computer Vision (ICCV), , 2015. pp. 477–485

48  Liu, X., Tanaka, M., Okutomi, M.: 'Single-image noise level estimation for blind denoising', *IEEE Trans. on image processing*, 2013, **22**, (12), pp. 5226–5237

49  Immerkaer, J.: 'Fast noise variance estimation', *Computer vision and image understanding*, 1996, **64**, (2), pp. 300–302

50  HonzÃątko, D., KruliÅą, M.: 'Accelerating block-matching and 3d filtering method for image denoising on gpus', *Journal of Real-Time Image Processing*, 2017