

# Learning to Enhance Low-Light Image via Zero-Reference Deep Curve Estimation

Chongyi Li, Chunle Guo, and Chen Change Loy, *Senior Member, IEEE*

**Abstract**— This paper presents a novel method, Zero-Reference Deep Curve Estimation (Zero-DCE), which formulates light enhancement as a task of image-specific curve estimation with a deep network. Our method trains a lightweight deep network, DCE-Net, to estimate pixel-wise and high-order curves for dynamic range adjustment of a given image. The curve estimation is specially designed, considering pixel value range, monotonicity, and differentiability. Zero-DCE is appealing in its relaxed assumption on reference images, i.e., it does not require any paired or even unpaired data during training. This is achieved through a set of carefully formulated non-reference loss functions, which implicitly measure the enhancement quality and drive the learning of the network. Despite its simplicity, we show that it generalizes well to diverse lighting conditions. Our method is efficient as image enhancement can be achieved by an intuitive and simple nonlinear curve mapping. We further present an accelerated and light version of Zero-DCE, called Zero-DCE++, that takes advantage of a tiny network with just 10K parameters. Zero-DCE++ has a fast inference speed (1000/11 FPS on a single GPU/CPU for an image of size 1200×900×3) while keeping the enhancement performance of Zero-DCE. Extensive experiments on various benchmarks demonstrate the advantages of our method over state-of-the-art methods qualitatively and quantitatively. Furthermore, the potential benefits of our method to face detection in the dark are discussed. The source code will be made publicly available at [https://li-chongyi.github.io/Proj\\_Zero-DCE++.html](https://li-chongyi.github.io/Proj_Zero-DCE++.html).

**Index Terms**—Computational photography, low-light image enhancement, curve estimation, zero-reference learning.



## 1 INTRODUCTION

MANY photos are often captured under suboptimal lighting conditions due to inevitable environmental and/or technical constraints. These include inadequate and unbalanced lighting conditions in the environment, incorrect placement of objects against extreme back light, and under-exposure during image capturing. Such low-light photos suffer from compromised aesthetic quality and unsatisfactory transmission of information. The former affects viewers' experience while the latter leads to wrong message being communicated, such as inaccurate object/face detection and recognition. In addition, although deep neural networks have shown impressive performance on image enhancement and restoration [1], [2], [3], [4], [5], they inevitably lead to high memory footprint and long inference time due to massive parameter space. The low computational cost and fast inference speed of deep models are desired in practical applications, especially for resource-limited and real-time devices, such as mobile platforms.

In this study, we present a novel deep learning-based method, Zero-Reference Deep Curve Estimation (Zero-DCE), for low-light image enhancement. It can cope with diverse lighting conditions including nonuniform and poor lighting cases. Instead of performing image-to-image mapping, we reformulate the task as an image-specific curve estimation problem. In particular, the proposed method takes a low-light image as input and produces high-order curves as its output. These curves are then used for pixel-wise adjustment on the dynamic range of the input to obtain an enhanced image. The curve estimation is carefully formulated so that it maintains the range of the enhanced image and preserves the

contrast of neighboring pixels. Importantly, it is differentiable, and thus we can learn the adjustable parameters of the curves through a deep convolutional neural network. The proposed network is lightweight and the designed curve can be iteratively applied to approximate higher-order curves for more robust and accurate dynamic range adjustment.

A unique advantage of our deep learning-based method is **zero-reference**, i.e., it does not require any paired or even unpaired data in the training process as in existing CNN-based [6], [7], [8] and GAN-based methods [9], [10]. This is made possible through a set of specially designed non-reference loss functions including spatial consistency loss, exposure control loss, color constancy loss, and illumination smoothness loss, all of which take into consideration multi-factor of light enhancement. We show that even with zero-reference training, Zero-DCE can still perform competitively against other methods that require paired or unpaired data for training. The proposed method is flexible. We provide options to balance the enhancement performance and the computational cost for the practical application of Zero-DCE and propose an accelerated and light version Zero-DCE++. This is achieved by re-designing the network structure, reformulating the curve estimation, and controlling the sizes of input image.

An example of enhancing a low-light image comprising nonuniform illumination is shown in Figure 1. Comparing to state-of-the-art methods, both Zero-DCE and Zero-DCE++ brighten up the image while preserving the inherent color and details. In contrast, both CNN-based method [6] and GAN-based method [9] yield under-(the face) and over-(the cabinet) enhancement. We show in this paper that our method obtains state-of-the-art performance both in qualitative and quantitative metrics. In addition, it is capable of improving high-level visual tasks, e.g., face detection, without inflicting high computational burden.

Our **contributions** are summarized as follows.

C. Li and C. C. Loy are with S-Lab, Nanyang Technological University (NTU), Singapore (e-mail: chongyi.li@ntu.edu.sg and ccloy@ntu.edu.sg).

C. Guo is with the College of Computer Science, Nankai University, Tianjin, China (e-mail: guochunle@nankai.edu.cn).

C. Li and C. Guo contribute equally.

C. C. Loy is the corresponding author.



Fig. 1: Visual comparisons on a typical low-light image comprising nonuniform illumination. The proposed Zero-DCE and Zero-DCE++ achieve visually pleasing results in terms of brightness, color, contrast, and naturalness, while existing methods either fail to cope with the extreme back light or generate color artifacts. In contrast to other deep learning-based methods, our approach is trained without any reference image.

- We propose the first low-light enhancement network that is independent of paired and unpaired training data, thus avoiding the risk of overfitting. As a result, our method generalizes well to various lighting conditions.
- We design an image-specific curve that is able to approximate pixel-wise and higher-order curves by iteratively applying itself. Such an image-specific curve can effectively perform mapping within a wide dynamic range.
- We show the potential of training a deep image enhancement model in the absence of reference images through task-specific non-reference loss functions that indirectly evaluate enhancement quality.
- The proposed Zero-DCE can be accelerated considerably while still keeping impressive enhancement performance. We provide multiple options to balance the enhancement performance and the cost of computational resources.

This work is an extension of our earlier conference version that has appeared in CVRP2020 [11]. In comparison to the conference version, we have introduced a significant amount of new materials. 1) We investigate the relations between the enhancement performance and the network structure, curve estimation, and input sizes. According to the investigation, we re-design the network structure, reformulate the curve formation, and control the sizes of input image, and thus present an accelerated and light version, called Zero-DCE++, which is more suitable for real-time enhancement on resource-limited devices. 2) Comparing to our earlier work, without compromising the enhancement performance, the trainable parameters (79K) and floating point operations (FLOPs) (84.99G) for an input image of size  $1200 \times 900 \times 3$  of Zero-DCE are reduced to 10K and 0.115G on Zero-DCE++. This translates to two times in runtime speed up, from 500 FPS in Zero-DCE to 1000 FPS in Zero-DCE++, for processing an image of size  $1200 \times 900 \times 3$  on a single NVIDIA 2080Ti GPU. In addition, even only with Intel Core i9-10920X CPU@3.5GHz, the processing time of Zero-DCE also can be reduced from 10s to 0.09s on Zero-DCE++, a 111 times speed up on a single CPU setting. The training time is also reduced from 30 minutes to 20 minutes. 3) We perform more experiments, design analysis, and ablation studies to demonstrate the advantages of zero-reference learning for low-light image enhancement and show the effectiveness of our method over existing state-of-the-art methods. 4) We conduct a more comprehensive literature survey on low-light image enhancement and discuss the advantages and limitations of current methods.

## 2 RELATED WORK

Our work is a new attempt for low-light image enhancement by combining zero-reference learning with deep curve estimation, which is rarely touched in the previous works. In what follows, we review the low-light image enhancement related works, including conventional methods and data-driven methods.

**Conventional Methods.** Histogram Equalization (HE)-based methods perform light enhancement through expanding the dynamic range of an image. Histogram distribution of images is adjusted at both global [12], [13] and local levels [14], [15]. There are also various methods adopting the Retinex theory [16] that typically decomposes an image into reflectance and illumination. The reflectance component is commonly assumed to be consistent under any lighting conditions; thus, light enhancement is formulated as an illumination estimation problem. Building on the Retinex theory, several methods have been proposed. Wang *et al.* [17] designed a naturalness- and information-preserving method when handling images of nonuniform illumination; Fu *et al.* [18] proposed a weighted variation model to simultaneously estimate the reflectance and the illumination of an image. The estimated reflectance is treated as the enhanced result; Guo *et al.* [19] first estimated a coarse illumination map by searching the maximum intensity of each pixel position, then refining the illumination map by a structure prior; Li *et al.* [20] proposed a new Retinex model that takes noise into consideration. The illumination map was estimated by solving an optimization problem.

Contrary to the conventional methods that fortuitously change the distribution of image histogram or that rely on potentially inaccurate physical models, the proposed method produces an enhanced result through image-specific curve mapping. Such a strategy enables light enhancement on images without creating unrealistic artifacts. Yuan and Sun [21] proposed an automatic exposure correction method, where the S-shaped curve for a given image is estimated by a global optimization algorithm and each segmented region is pushed to its optimal zone by curve mapping. Different from [21], our method is purely data-driven and takes multiple light enhancement factors into consideration in the design of the non-reference loss functions, and thus enjoys better robustness, wider image dynamic range adjustment, and lower computational burden.

**Data-Driven Methods.** Data-driven methods are largely categorized into two branches, namely Convolutional Neural Network (CNN)-based and Generative Adversarial Network (GAN)-based methods. Most CNN-based solutions rely on paired data for supervised training, therefore they are resource-intensive. Often time, the paired data are exhaustively collected through automatic

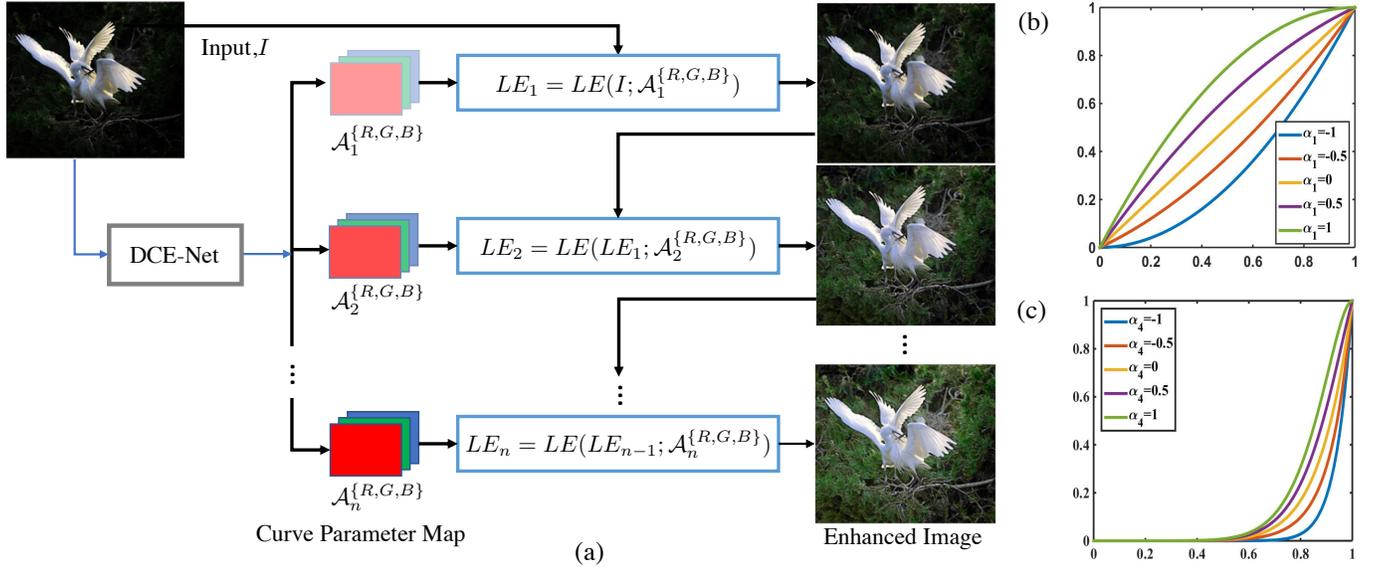


Fig. 2: (a) The framework of Zero-DCE. A DCE-Net is devised to estimate a set of best-fitting Light-Enhancement curves (LE-curves) that iteratively enhance a given input image (*i.e.*, takes the enhanced image as the input of next iteration and the input is enhanced in a progressive manner). (b, c) LE-curves with different adjustment parameters  $\alpha$  and numbers of iteration  $n$ . In (c),  $\alpha_1$ ,  $\alpha_2$ , and  $\alpha_3$  are equal to -1 while  $n$  is equal to 4. In each subfigure, the horizontal axis represents the input pixel values while the vertical axis represents the output pixel values.

light degradation, changing the settings of cameras during data capturing, or synthesizing data via image retouching. For example, LL-Net [22] and MBLLEN [23] were trained on data simulated on random Gamma correction; the LOL dataset [7] of paired low/normal light images was collected through altering the exposure time and ISO during image acquisition; the MIT-Adobe FiveK dataset [24] comprises 5,000 raw images, each of which has five retouched images produced by trained experts. MIT-Adobe FiveK dataset was originally collected for image global retouching; the SID [25] provides paired low/normal light raw data; a dataset of raw low-light videos with the corresponding normal light videos captured at video rate was collected in [26].

Inspired by the Retinex model, recent deep models design the networks to estimate the reflectance and illumination of an input image by supervised learning with paired data. Ren *et al.* [27] proposed a deep hybrid network for low-light image enhancement, which consists of two streams to learn the global content and the salient structures in a unified network. Wang *et al.* [6] proposed an underexposed photo enhancement network by estimating the illumination map. This network was trained on paired data that were retouched by three experts. Zhang *et al.* [28] built a network for kindling the darkness of an image, called KinD, which decomposes images into two components. The illumination component is responsible for the light adjustment while the reflectance component is for degradation removal. Retinex model-based deep models still suffer from the same limitations as the conventional Retinex-based methods, such as ideal assumption.

More recently, Xu *et al.* [8] proposed a frequency-based decomposition-and-enhancement model for low-light image enhancement. This model first recovers the image content in the low-frequency layer, then enhances high-frequency details based on the recover image content. This model is trained on a low-light dataset of real noisy low-light and ground truth sRGB image pairs.

Understandably, light enhancement solutions based on paired

data are impractical in many ways, considering the high cost involved in collecting sufficient paired data as well as the inclusion of factitious and unrealistic data in training the deep models. Such constraints are reflected in the poor generalization capability of CNN-based methods. Artifacts and color casts are commonly generated, when these methods are presented with real-world images of various light intensities.

Unsupervised GAN-based methods have the advantage of eliminating paired data for training. An unsupervised GAN-based method, EnlightenGAN [9], learns to enhance low-light images using unpaired low/normal light data. The network was trained by taking elaborately designed discriminators and loss functions into account. However, unsupervised GAN-based solutions usually require careful selection of unpaired training data.

To integrate the superiorities of CNNs and GANs, Yang *et al.* [29] proposed a semi-supervised model for low-light image enhancement, which performs enhancement in two stages. In the first stage, a coarse-to-fine band representation is learned and different band signals are inferred jointly in a recursive process with paired data. In the second stage, the band representation is recomposed via adversarial learning. Although the semi-supervised learning framework can effectively improve the generalization capability of the deep model, it still takes the risk of overfitting on paired training data and induces high memory footprint.

The proposed method is superior to existing data-driven methods in three aspects. First, it explores a new learning strategy, *i.e.*, one requires *zero reference*, hence eliminating the need for paired and unpaired data. Second, the network is trained by taking carefully defined non-reference loss functions into account. This strategy allows output image quality to be implicitly evaluated, the results of which would be reiterated for network learning. Third, our method is highly efficient and cost-effective. The accelerated and light version Zero-DCE++ only contains 10K trainable parameters and 0.115G FLOPs, achieves 1000/11 FPS inference time

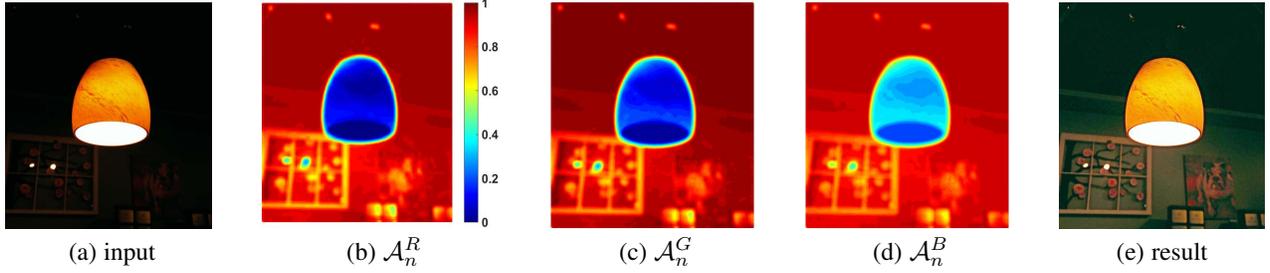


Fig. 3: An example of the pixel-wise curve parameter maps. For visualization, we average the curve parameter maps of all iterations ( $n = 8$ ) and normalize the values to the range of  $[0, 1]$ .  $\mathcal{A}_n^R$ ,  $\mathcal{A}_n^G$ , and  $\mathcal{A}_n^B$  represent the averaged best-fitting curve parameter maps of R, G, and B channels, respectively. The maps in (b), (c), and (d) are represented by heatmaps.

on a single GPU/CPU, and needs 20 minutes for training. The efficiency of our method precedes current deep models [6], [7], [9], [23] by a large margin. These advantages benefit from our zero-reference learning framework, lightweight network structure, and effective non-reference loss functions.

### 3 METHODOLOGY

We show the framework of Zero-DCE in Figure 2. A Deep Curve Estimation Network (DCE-Net) is devised to estimate a set of best-fitting Light-Enhancement curves (LE-curves) given an input image. The framework then maps all pixels of the input’s RGB channels by applying the curves iteratively for obtaining the final enhanced image. In what follows, we detail the key components, namely LE-curve, DCE-Net, and non-reference loss functions.

#### 3.1 Light-Enhancement Curve

Inspired by curve adjustment used in photo editing software, we design a kind of curve that can map a low-light image to its enhanced version automatically, where the self-adaptive curve parameters are solely dependent on the input image. There are three objectives in the design of such a curve:

- 1) each pixel value of the enhanced image should fall in the normalized range of  $[0, 1]$  to avoid information loss induced by overflow truncation;
- 2) this curve should be monotonous to preserve the differences (contrast) of neighboring pixels; and
- 3) the form of this curve should be as simple as possible and differentiable in the process of gradient backpropagation.

To achieve these three objectives, we design a quadratic curve, which can be expressed as:

$$LE(I(\mathbf{x}); \alpha) = I(\mathbf{x}) + \alpha I(\mathbf{x})(1 - I(\mathbf{x})), \quad (1)$$

where  $\mathbf{x}$  denotes pixel coordinates,  $LE(I(\mathbf{x}); \alpha)$  is the enhanced version of the given input  $I(\mathbf{x})$ , the trainable curve parameter  $\alpha \in [-1, 1]$  adjusts the magnitude of LE-curve and also controls the exposure level. Each pixel of input is normalized to the range of  $[0, 1]$  and all operations are pixel-wise. We separately apply the LE-curve to three RGB channels instead of solely on the luminance channel. The three-channel adjustment can better preserve the inherent color and reduce the risk of over-saturation. We report more details in the ablation study.

An illustration of LE-curves with different adjustment parameters  $\alpha$  is shown in Figure 2(b). It is clear that the LE-curve complies with the three aforementioned objectives. Thus,

each pixel value of enhanced images is in the range of  $[0, 1]$ . In addition, the LE-curve enables us to increase or decrease the dynamic range of an input image. This capability is conducive to not only enhancing low-light regions but also removing over-exposure artifacts. We choose a specific single-parameter form for the quadratic because 1) the single-parameter form can reduce the computational cost and speed up our method and 2) the specially designed quadratic meets the three objectives of our designs and already achieves satisfactory enhancement performance.

**Higher-Order Curve.** The LE-curve defined in Equation (1) can be applied iteratively to enable more versatile adjustment to cope with challenging low-light conditions. Specifically,

$$LE_n(\mathbf{x}) = LE_{n-1}(\mathbf{x}) + \alpha_n LE_{n-1}(\mathbf{x})(1 - LE_{n-1}(\mathbf{x})), \quad (2)$$

where  $n$  is the number of iteration, which controls the curvature. In this paper, we set the value of  $n$  to 8, which can deal with most cases satisfactorily. Equation (2) can be degraded to Equation (1) when  $n$  is equal to 1. Figure 2(c) provides an example showing high-order curves with different  $\alpha$  and  $n$ . Such high-order curves offer more powerful adjustment capability (*i.e.*, greater curvature) than the curves in Figure 2(b).

**Pixel-Wise Curve.** In comparison to a single-order curve, a higher-order curve adjusts an image within a wider dynamic range. Nonetheless, it is still a global adjustment since  $\alpha$  is used for all pixels. A global mapping tends to over-/under- enhance local regions. To address this problem, we formulate  $\alpha$  as a pixel-wise parameter, *i.e.*, each pixel of the given input image has a corresponding curve with the best-fitting  $\alpha$  to adjust its dynamic range. Hence, Equation (2) can be reformulated as:

$$LE_n(\mathbf{x}) = LE_{n-1}(\mathbf{x}) + \mathcal{A}_n(\mathbf{x}) LE_{n-1}(\mathbf{x})(1 - LE_{n-1}(\mathbf{x})), \quad (3)$$

where  $\mathcal{A}$  is a parameter map with the same size as the given image. Here, pixels in a local region are assumed to having the same intensity (also the same adjustment curves), and thus the neighboring pixels in the output result still preserve the monotonous relations. In this way, the pixel-wise higher-order curves also comply with the three aforementioned objectives. As a result, each pixel value of enhanced images is still in the range of  $[0, 1]$ .

We present an example of the estimated curve parameter maps in Figure 3. As shown, the best-fitting parameter maps of different color channels have similar adjustment tendency but different values, indicating the relevance and difference among the three channels of a low-light image. The curve parameter map accurately indicates the brightness of different regions (*e.g.*, the two glitters on the wall). With the fitting maps, the enhanced version

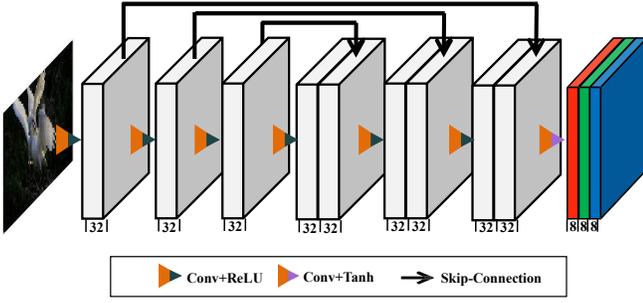


Fig. 4: The architecture of DCE-Net.

image can be directly obtained by pixel-wise curve mapping. As shown in Figure 3(e), the enhanced version reveals the content in dark regions and preserves the bright regions.

### 3.2 DCE-Net

To learn the mapping between an input image and its best-fitting curve parameter maps, we propose a Deep Curve Estimation Network (DCE-Net). In Figure 4, we present the detailed network architecture and parameter settings of DCE-Net.

The input to the DCE-Net is a low-light image while the outputs are a set of pixel-wise curve parameter maps for corresponding higher-order curves. Instead of employing fully connected layers that require fixed input sizes, we employ a plain CNN of seven convolutional layers with symmetrical skip concatenation. In the first six convolutional layers, each convolutional layer consists of 32 convolutional kernels of size  $3 \times 3$  and stride 1 followed by the ReLU activation function. The last convolutional layer consists of 24 convolutional kernels of size  $3 \times 3$  and stride 1 followed by the Tanh activation function, which produces 24 curve parameter maps for eight iterations, where each iteration generates three curve parameter maps for the three channels (*i.e.*, RGB channels). We discard the down-sampling and batch normalization layers that break the relations of neighboring pixels. It is noteworthy that DCE-Net only has 79K trainable parameters and 85G FLOPs for an input image of size  $1200 \times 900 \times 3$ , which is already smaller than existing low-light image enhancement deep models, such as RetinexNet [7]: 555K/587G, EnlightenGAN [9]: 8M/273G), and MBLLN [23]: 450K/301G.

### 3.3 Non-Reference Loss Functions

To enable zero-reference learning in DCE-Net, we propose a set of differentiable non-reference losses that allow us to evaluate the quality of enhanced images. The following four types of losses are adopted to train our DCE-Net.

**Spatial Consistency Loss.** The spatial consistency loss  $L_{spa}$  encourages spatial coherence of the enhanced image through preserving the difference of neighboring regions between the input image and its enhanced version:

$$L_{spa} = \frac{1}{K} \sum_{i=1}^K \sum_{j \in \Omega(i)} (|(Y_i - Y_j)| - |(I_i - I_j)|)^2, \quad (4)$$

where  $K$  is the number of local region, and  $\Omega(i)$  is the four neighboring regions (top, down, left, right) centered at the region  $i$ . We denote  $Y$  and  $I$  as the average intensity value of the local region in the enhanced version and input image, respectively. We

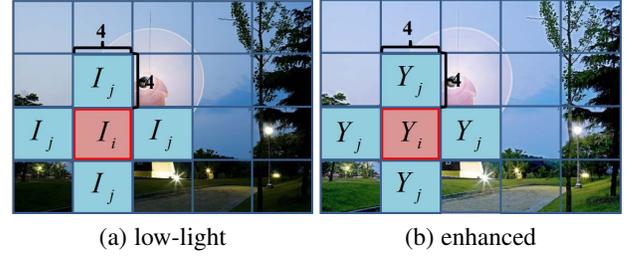


Fig. 5: An illustration of the spatial consistency loss.

empirically set the size of the local region to  $4 \times 4$ . This loss is stable given other region sizes. We illustrate the process of computing the spatial consistency loss in Figure 5.

**Exposure Control Loss.** To restrain under-/over-exposed regions, we design an exposure control loss  $L_{exp}$  to control the exposure level. The exposure control loss measures the distance between the average intensity value of a local region to the well-exposedness level  $E$ . We follow existing practices [30], [31] to set  $E$  as the gray level in the RGB color space. We empirically set  $E$  to 0.6 in our experiments. The loss  $L_{exp}$  can be expressed as:

$$L_{exp} = \frac{1}{M} \sum_{k=1}^M |Y_k - E|, \quad (5)$$

where  $M$  represents the number of non-overlapping local regions of size  $16 \times 16$ , the average intensity value of a local region in the enhanced image is represented as  $Y$ .

**Color Constancy Loss.** Following the Gray-World color constancy hypothesis [32] that color in each sensor channel averages to gray over the entire image, we design a color constancy loss to correct the potential color deviations in the enhanced image and also build the relations among the three adjusted channels. The color constancy loss  $L_{col}$  can be expressed as:

$$L_{col} = \sum_{\forall (p,q) \in \varepsilon} (J^p - J^q)^2, \varepsilon = \{(R, G), (R, B), (G, B)\}, \quad (6)$$

where  $J^p$  denotes the average intensity value of  $p$  channel in the enhanced image, a pair of channels is represented as  $(p, q)$ .

**Illumination Smoothness Loss.** To preserve the monotonicity relations between neighboring pixels, we add an illumination smoothness loss to each curve parameter map  $\mathcal{A}$ . The illumination smoothness loss  $L_{tv_{\mathcal{A}}}$  is defined as:

$$L_{tv_{\mathcal{A}}} = \frac{1}{N} \sum_{n=1}^N \sum_{c \in \xi} (|\nabla_x \mathcal{A}_n^c| + |\nabla_y \mathcal{A}_n^c|)^2, \xi = \{R, G, B\}, \quad (7)$$

where  $N$  is the number of iteration, the horizontal and vertical gradient operations are represented as  $\nabla_x$  and  $\nabla_y$ , respectively.

**Total Loss.** The total loss can be expressed as:

$$L_{total} = L_{spa} + L_{exp} + W_{col} L_{col} + W_{tv_{\mathcal{A}}} L_{tv_{\mathcal{A}}}, \quad (8)$$

where the weights  $W_{col}$  and  $W_{tv_{\mathcal{A}}}$  are used for balancing the scales of different losses.

## 4 ZERO-DCE++

Though Zero-DCE is already faster and smaller than existing deep learning-based models [7], [9], [23], reduced computational cost and faster inference speed are still desired for practical

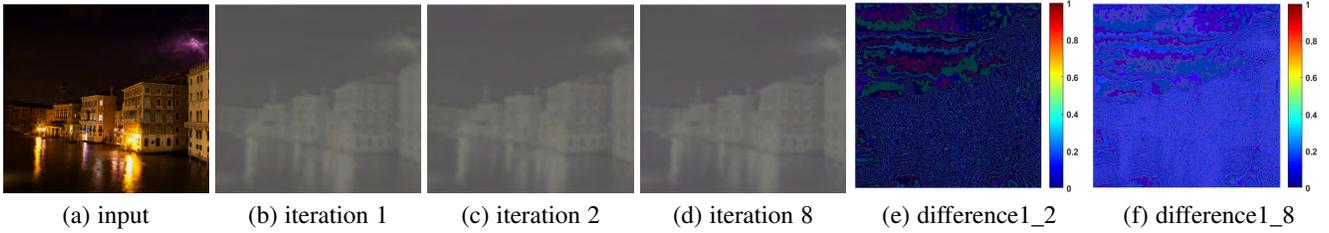


Fig. 6: The estimated curve parameter maps in different iteration stages. Subfigure (b), (c), and (d) depict the estimated curve parameter maps in iterations 1, 2, and 8, respectively. Subfigure (e) and (f) show the difference maps between iteration 1 and iteration 2 as well as between iteration 1 and iteration 8, respectively. For visualization, we normalize the curve parameter maps and amplify the intensity of the difference maps by 30 times.

applications, especially when dealing with large images captured by modern mobile devices. We propose an accelerated and light version of Zero-DCE, called Zero-DCE++, to achieve the above-mentioned characteristics.

To this end, we carefully investigate the relations between enhancement performance and network structure, curve estimation, and input sizes. We observed that 1) the convolutional layers used in DCE-Net can be replaced with the more efficient depthwise separable convolutions [33] that are commonly used in computer vision tasks [34], [35], [36] for reducing network parameters without compromising the performance much; 2) the estimated curve parameters in different iteration stages (a total of eight iterations in Zero-DCE) are similar in most cases. In Figure 6, we show an example of the estimated curve parameter maps in different iteration stages and their difference maps. As observed, the curve parameter maps are similar and the values in the difference maps are small. Such results manifest that the curve parameter map can be reused in different iteration stages to handle most cases, thus we can reduce the estimated curve parameter maps from 24 to 3; and 3) our method is not sensitive to the sizes of input image. Consequently, we can use the downsampled input as the input of curve parameter estimation network and then upsample the estimated curve parameter maps back to the original resolutions for image enhancement. The low-resolution input can significantly reduce the computational cost. Based on these observations, we modify the Zero-DCE from three aspects.

First, we re-design the DEC-Net by replacing the convolutional layers with depthwise separable convolutions for reducing the network parameters. Each depthwise separable convolutional layer consists of a depthwise convolution with kernels of size  $3 \times 3$  and stride 1 and a pointwise convolution with kernels of size  $1 \times 1$  and stride 1.

Second, we reformulate the curve estimation and only estimate 3 curve parameter maps, then reuse them in different iteration stages instead of estimating 24 parameter maps across eight iterations. Thus, Equation (3) can be reformulated as

$$LE_n(\mathbf{x}) = LE_{n-1}(\mathbf{x}) + \mathcal{A}(\mathbf{x})LE_{n-1}(\mathbf{x})(1 - LE_{n-1}(\mathbf{x})), \quad (9)$$

where the same curve parameter map  $\mathcal{A}$  is used to adjust the curves in different iteration stages. Although we reuse the curve parameter maps, it still retains the high-order property thanks to the iteration process.

Third, we can use the downsampled image as the input of our network to estimate the curve parameter maps. By default, we downsample an input by a factor of 12 in Zero-DCE++ to balance the enhancement performance and computational cost. Even with

the extreme downsampling factors, our method maintains a good performance. The reasons are briefly explained as follows. Firstly, although we adopt the downsampled input to estimate the curve parameters, we resize the small curve parameter maps back to the same size as the original input image based on the assumption that the pixels in a local region have the same intensity (also the same adjustment curves). The mapping from the input image to the enhanced image is conducted on the original resolution. Secondly, the proposed spatial consistency loss encourages the results to preserve the content of the input image. Thirdly, the adopted losses in our framework are region-wise but not pixel-wise. We present more discussions and results in the ablation study.

These modifications offer Zero-DCE++ the advantages of having a tiny network (10K trainable parameters, 0.115G FLOPs for an image of size  $1200 \times 900 \times 3$ ), real-time inference speed (1000/11 FPS on a single GPU/CPU for an image of size  $1200 \times 900 \times 3$ ), and fast training (20 minutes).

## 5 EXPERIMENTS

### 5.1 Implementation Details

CNN-based models usually use self-captured paired data for network training [7], [25] while GAN-based models elaborately select unpaired data [9], [37]. To bring the capability of wide dynamic range adjustment into full play, we incorporate both low-light and over-exposed images into our training set. To this end, we employ 360 multi-exposure sequences from the Part1 of SICE dataset [38] to train our model. The dataset is also used as a part of the training data in EnlightenGAN [9]. We randomly split 3,022 images of different exposure levels in the Part1 subset [38] into two parts (2,422 images for training and the rest for validation). We resize the training and testing images to the size of  $512 \times 512 \times 3$ .

We implement our framework with MindSpore on an NVIDIA 2080Ti GPU. A batch size of 8 is applied. The filter weights of each layer are initialized with standard zero mean and 0.02 standard deviation Gaussian function. Bias is initialized as a constant. We use ADAM optimizer with default parameters and fixed learning rate  $1e^{-4}$  for our network optimization. The weights  $W_{col}$  and  $W_{tv_A}$  are set to 0.5, and 20, respectively, to balance the scale of losses. Zero-DCE and Zero-DCE++ adopt the same training dataset and configurations during training.

### 5.2 Experimental Settings

We compare our method with several state-of-the-art methods: three conventional methods (SRIE [18], LIME [19], Li *et al.* [20]), four CNN-based methods (Wang *et al.* [6], RetinexNet [7],

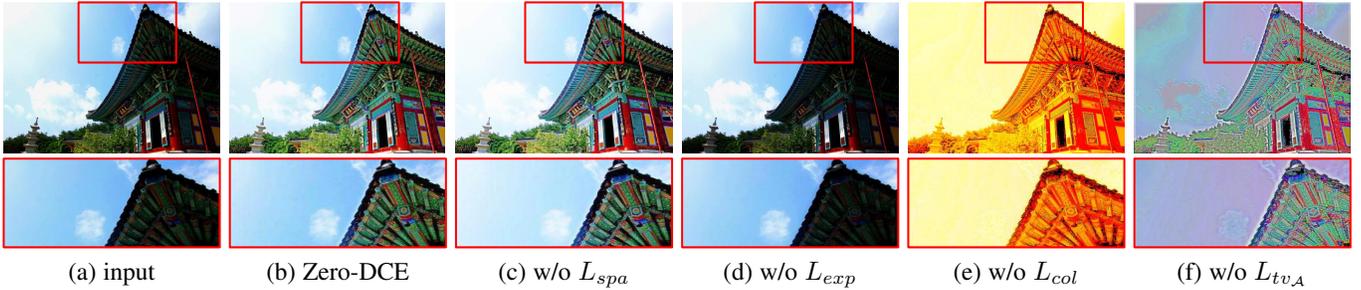


Fig. 7: Ablation study of the contribution of each loss (spatial consistency loss  $L_{spa}$ , exposure control loss  $L_{exp}$ , color constancy loss  $L_{col}$ , illumination smoothness loss  $L_{tv_A}$ ). Red boxes indicate the obvious differences and amplified details.

LightenNet [39], MBLLEN [23]), and one GAN-based method (EnlightenGAN [9]). The results are reproduced using publicly available source codes with recommended parameters.

We perform qualitative and quantitative experiments on standard image sets used by previous works including NPE [17] (84 images), LIME [19] (10 images), MEF [40] (17 images), DICM [41] (64 images), and VV<sup>3</sup> (24 images). Besides, we quantitatively validate our method on the Part2 subset of SICE dataset [38], which consists of 229 multi-exposure sequences and the corresponding reference image for each multi-exposure sequence. For a fair comparison, we only use the low-light images of Part2 subset [38] for testing, since baselines cannot handle over-exposed images well. Specifically, we choose the first three (resp. four) low-light images if there are seven (resp. nine) images in a multi-exposure sequence and resize all images to a size of  $1200 \times 900 \times 3$ . Finally, we obtain 767 paired low/normal light images, denoted as Part2 testing set.

The low/normal light image dataset mentioned in [42] was discarded because the training datasets of RetinexNet [7] and EnlightenGAN [9] consist of some images from this dataset. We did not use the MIT-Adobe FiveK dataset [24] as it is not primarily designed for underexposed photos enhancement and still contains some low-light images in the ground truth set. Note that this paper only focuses on low-light image enhancement on RGB images, thus we did not include the methods that require raw data as inputs and were designed for general photo enhancement.

### 5.3 Ablation Study

We perform ablation studies to demonstrate the effectiveness of each component of Zero-DCE. Additionally, the comparisons between Zero-DCE and Zero-DCE++ are carried out to analyze the advantages and disadvantages of the accelerated and light version at the end of this section.

**Contribution of Each Loss.** We present the results of Zero-DCE trained by various combinations of losses in Figure 7. The result without spatial consistency loss  $L_{spa}$  has relatively lower contrast (e.g., the cloud regions) than the full result. This shows the importance of  $L_{spa}$  in preserving the difference of neighboring regions between the input and the enhanced image. Removing the exposure control loss  $L_{exp}$  fails to recover the low-light region. Severe color casts emerge when the color constancy loss  $L_{col}$  is discarded. This variant ignores the relations among three channels when curve mapping is applied. Finally, removing the illumination smoothness loss  $L_{tv_A}$  hampers the correlations

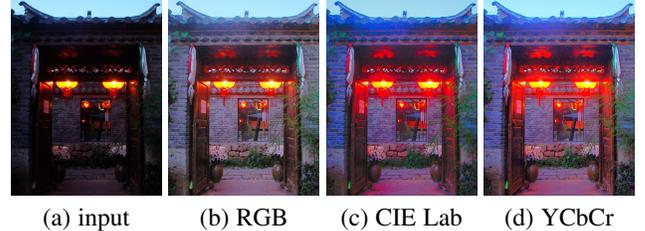


Fig. 8: Ablation study of the advantage of three channels adjustment (RGB, CIE Lab, and YCbCr color spaces).

between neighboring regions leading to obvious artifacts. Such results demonstrate that each loss used in our zero-reference learning framework plays a significant role in achieving the final visually pleasing results.

**Advantage of Three Channels Adjustment.** To demonstrate the advantage of three channels adjustment, we try to adjust the illumination related channel only in CIE Lab and YCbCr color spaces using the same configurations as the adjustment in RGB color space, except removing the color constancy loss which is only available for three channels adjustment.

Specifically, we first transfer the input from RGB color space to CIE Lab (YCbCr) color space, then feed the L (Y) component to the DCE-Net for estimating a set of curve parameter maps, where we compute each loss in L (Y) channel in the phase of training. At last, we adjust the L (Y) component using Equation 3 with the estimated curve parameters. After the adjustment of L (Y) component, the corresponding ab (CbCr) components are adjusted accordingly (equal proportion adjustment). In Figure 8, we show an example to demonstrate the advantage of three channels adjustment. As observed, all the results show improved brightness and contrast, suggesting the effectiveness of both the single channel adjustment (CIE Lab and YCbCr color spaces) and the three-channel adjustment (RGB color space) in improving the brightness of the given low-light image. However, the results adjusted in CIE Lab and YCbCr color spaces as shown in Figure 8(c) and (d) have obvious color deviations (e.g., the color of wall) and over-saturation (e.g., the region of lantern). The visual comparison suggests that three channels adjustment can better preserve the inherent color and reduce the risk of over-saturation.

**Effect of Parameter Settings.** We evaluate the effect of parameters in Zero-DCE, consisting of the depth and width of the DCE-Net and the number of iterations. A visual example is presented in Figure 9. As observed in Figure 9(b), with just three convolutional

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



Fig. 9: Ablation study of the effect of parameter settings.  $l$ - $f$ - $n$  represents the proposed Zero-DCE with  $l$  convolutional layers,  $f$  feature maps of each layer (except the last layer), and  $n$  iterations.



Fig. 10: Ablation study of the impact of training data. Zero-DCE<sub>Low</sub> represents that the Zero-DCE was trained on only 900 low-light images out of 2,422 images in the original training set. Zero-DCE<sub>LargeL</sub> represents that the Zero-DCE was trained on 9,000 unlabeled low-light images provided in the DARK FACE dataset [42]. Zero-DCE<sub>LargeLH</sub> represents that the Zero-DCE was trained on 4800 multi-exposure images from the data augmented combination of Part1 and Part2 subsets in the SICE dataset [38]. (b) suggests that Zero-DCE has a good balance between over-enhancement and under-enhancement.

layers, Zero-DCE<sub>l3-f32-n8</sub> can already produce satisfactory results, suggesting the effectiveness of zero-reference learning. The Zero-DCE<sub>l7-f32-n8</sub> and Zero-DCE<sub>l7-f32-n16</sub> produce the most visually pleasing results with natural exposure and proper contrast. By reducing the number of iterations to 1, an obvious decrease in performance is observed on Zero-DCE<sub>l7-f32-n1</sub> as shown in Figure 9(d). This is because the curve with only single iteration has limited adjustment capability. This suggests the need for higher-order curves in our method.

The same tendency also can be found in the quantitative comparisons in Table 1. The comparison between the input and the enhanced results by Zero-DCE<sub>l3-f32-n8</sub> suggests the effectiveness of the proposed method despite the network only contains three convolutional layers. The Zero-DCE<sub>l7-f32-n1</sub> achieves the worst quantitative performance due to the limited adjustment capability of only one-time iteration (*i.e.*,  $n=1$ ), suggesting the importance of using more iterations. When we increase the number of feature maps of each layer from 16 to 32, the quantitative performance is improved (*i.e.*, Zero-DCE<sub>l7-f16-n8</sub> and Zero-DCE<sub>l7-f32-n8</sub>). Increasing the number of iterations from 8 to 16 only boosts the average PSNR value marginally (*i.e.*, Zero-DCE<sub>l7-f32-n8</sub> and Zero-DCE<sub>l7-f32-n16</sub>). Consequently, we choose Zero-DCE<sub>f7-l32-n8</sub> as the final model based on its good trade-off between efficiency and restoration performance.

**Impact of Training Data.** To test the impact of training data, we retrain the Zero-DCE on datasets different from that described in Sec. 5.1. As shown in Figure 10(c) and (d), after removing the over-exposed training data, Zero-DCE tends to over-enhance the well-lit regions (*e.g.*, the cup in the results of Zero-DCE<sub>Low</sub> and Zero-DCE<sub>LargeL</sub>), in spite of using more low-light images, (*i.e.*, Zero-DCE<sub>LargeL</sub>). Such results indicate the rationality and

TABLE 1: Quantitative comparisons in terms of Peak Signal-to-Noise Ratio (PSNR, dB), Structural Similarity (SSIM) [43], and Mean Absolute Error (MAE). These comparisons are carried out on Part2 testing set.  $l$ - $f$ - $n$  represents the proposed Zero-DCE with  $l$  convolutional layers,  $f$  feature maps of each layer (except the last layer), and  $n$  iterations.

Method	PSNR $\uparrow$	SSIM $\uparrow$	MAE $\downarrow$
input	10.71	0.33	209.65
l3-f32-n8	14.50	0.56	119.40
l7-f16-n8	15.67	0.58	111.01
l7-f32-n1	12.21	0.42	172.89
l7-f32-n8	16.57	0.59	98.78
l7-f32-n16	16.79	0.57	98.70

TABLE 2: Quantitative comparisons in terms of PSNR, SSIM, and MAE. These comparisons are carried out on Part2 testing set.

Method	PSNR $\uparrow$	SSIM $\uparrow$	MAE $\downarrow$
input	10.71	0.33	209.65
Zero-DCE <sub>E0.4</sub>	11.82	0.45	177.87
Zero-DCE <sub>E0.5</sub>	15.19	0.56	114.47
Zero-DCE <sub>E0.6</sub>	16.57	0.59	98.78
Zero-DCE <sub>E0.7</sub>	13.40	0.55	156.83

necessity of the usage of multi-exposure training data in the training process of our network. In addition, the Zero-DCE can better recover the dark regions (*e.g.*, the roses) when more multi-exposure training data are used (*i.e.*, Zero-DCE<sub>LargeLH</sub>), as shown in Figure 10(e). The proposed Zero-DCE as shown in Figure 10(b) has a good balance between over-enhancement and under-enhancement. For a fair comparison with other deep learning-based methods, we use a comparable amount of training data with them although more training data can bring better visual



Fig. 11: A visual comparison among the results generated by the Zero-DCE trained using different well-exposedness level,  $E$ , in exposure control loss (see Equation (5)).

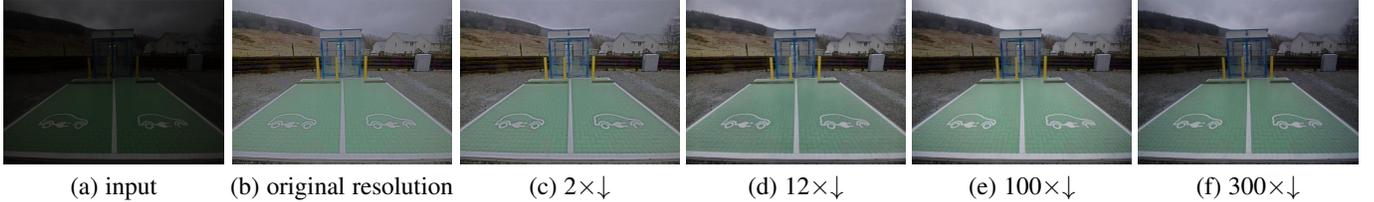


Fig. 12: A set of results by feeding different sizes of input to the modified framework. Even  $300\times$  downsampling does not hamper the performance much when compared with the original resolution as input. Here,  $\downarrow$  represents the downsampling operation.

TABLE 3: The statistic relations between enhancement performance and input sizes measured in PSNR and FLOPs. The FLOPs (in G) is computed for an image of size  $1200\times 900\times 3$ . “number $\times\downarrow$ ” indicates the times of downsampling the input image. These comparisons are carried out on Part2 testing set.

Metrics	original resolution	$2\times\downarrow$	$4\times\downarrow$	$6\times\downarrow$	$12\times\downarrow$	$20\times\downarrow$	$50\times\downarrow$	$75\times\downarrow$	$100\times\downarrow$	$300\times\downarrow$
PSNR	16.09	16.17	16.29	16.37	16.42	16.33	15.85	15.56	15.35	14.53
FLOPs	11.442	2.887	0.749	0.352	0.115	0.064	0.040	0.037	0.036	0.035

TABLE 4: Ablation study between Zero-DCE and Zero-DCE++ in terms of PSNR, trainable parameters (#P), and FLOPs (in G). The FLOPs is computed for an image of size  $1200\times 900\times 3$ . The average PSNR values are computed on Part2 testing set.

Zero-DCE	DSconv	Pshared	PSNR $\uparrow$	#P $\downarrow$	FLOPs $\downarrow$
✓			16.57	79,416	84.99
	✓		16.51	11,926	0.375
		✓	16.24	67,299	0.540
	✓	✓	16.42	10,561	0.115

performance to our approach.

**Effect of Well-Exposedness Level.** We study the effect of well-exposedness level  $E$  used in the exposure control loss on the enhancement performance of our method. We set four different well-exposedness levels  $E$  (*i.e.*, 0.4, 0.5, 0.6, 0.7) to train our network, denoted as Zero-DCE $_{E0.4}$ , Zero-DCE $_{E0.5}$ , Zero-DCE $_{E0.6}$  (*i.e.*, our final Zero-DCE model), and Zero-DCE $_{E0.7}$ , respectively. A set of visual results are shown in Figure 11. The quantitative comparisons are presented in Table 2.

As shown in Table 2, Zero-DCE $_{E0.6}$  achieves the best quantitative scores. Zero-DCE $_{E0.5}$  obtains comparable performance to Zero-DCE $_{E0.6}$ . The quantitative performance of Zero-DCE $_{E0.4}$  and Zero-DCE $_{E0.7}$  is slightly inferior to that of Zero-DCE $_{E0.6}$  and Zero-DCE $_{E0.5}$ . As observed in Figure 11, Zero-DCE $_{E0.5}$  and Zero-DCE $_{E0.6}$  obtain visually pleasing brightness. In contrast, Zero-DCE $_{E0.4}$  produces under-exposure while Zero-DCE $_{E0.7}$  over-enhances the input image. Finally, we choose Zero-DCE $_{E0.6}$  as the final model based on its good qualitative and quantitative performance.

**Zero-DCE VS. Zero-DCE++.** We first analyze the effect of input

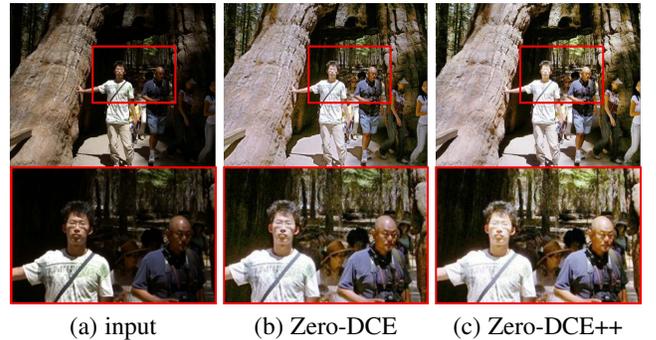


Fig. 13: A visual comparison between the results generated by Zero-DCE and Zero-DCE++. Zero-DCE shows better capability in handling extreme lighting conditions.

sizes on the enhancement performance in our method. As specified in Sec. 4, we first replace convolutional layers of DCE-Net with depthwise separation convolutions and reuse the curve parameter maps across eight iterations. Then, we feed the different sizes of input to the modified framework. The statistic relations between enhancement performance and input sizes are summarized in Table 3. We also show several results by feeding different sizes of input to the modified framework in Figure 12. As shown in Table 3 and Figure 12, downsampling the sizes of input has unnoticeable effect on the enhancement performance but significantly saves computational cost (measured in FLOPs). As shown, results of  $12\times\downarrow$  achieve the highest average PSNR value, thus we adopt it as the default operation in Zero-DCE++.

TABLE 5: User study (US) $\uparrow$ /Perceptual index (PI) $\downarrow$  scores on the image sets (NPE, LIME, MEF, DICM, VV). Higher US score indicates better human subjective visual quality while lower PI value indicates better perceptual quality. The best result is in red whereas the second best one is in blue under each case.

Method	NPE	LIME	MEF	DICM	VV	Average
SRIE [18]	3.65/2.79	3.50/2.76	3.22/2.61	3.42/3.17	2.80/3.37	3.32/2.94
LIME [19]	3.78/3.05	3.95/3.00	3.71/2.78	3.31/3.35	3.21/3.03	3.59/3.04
Li <i>et al.</i> [20]	3.80/3.09	3.78/3.02	2.93/3.61	3.47/3.43	2.87/3.37	3.37/3.72
LightenNet [39]	3.76/2.88	3.02/2.84	3.07/2.51	3.11/3.13	2.55/3.29	2.70/2.93
MBLLEN [23]	3.81/2.77	3.77/3.18	3.21/3.04	3.07/3.19	2.72/3.63	3.33/3.16
RetinexNet [7]	3.30/3.18	2.32/3.08	2.80/2.86	2.88/3.24	1.96/2.95	2.58/3.06
Wang <i>et al.</i> [6]	3.83/2.83	3.82/2.90	3.13/2.72	3.44/3.20	2.95/3.42	3.43/3.01
EnlightenGAN [9]	3.90/2.96	3.84/2.83	3.75/2.45	3.50/3.13	3.17/4.71	3.63/3.22
Zero-DCE	3.81/2.84	3.80/2.76	4.13/2.43	3.52/3.04	3.24/3.33	3.70/2.88
Zero-DCE++	3.79/2.93	3.81/2.97	4.10/2.50	3.48/3.21	3.26/3.31	3.69/2.98

TABLE 6: Quantitative comparisons in terms of PSNR, SSIM, and MAE on the Part2 testing set. The best result is in red whereas the second best one is in blue under each case.

Method	PSNR $\uparrow$	SSIM $\uparrow$	MAE $\downarrow$
SRIE [18]	14.41	0.54	127.08
LIME [19]	16.17	0.57	108.12
Li <i>et al.</i> [20]	15.19	0.54	114.21
RetinexNet [7]	15.99	0.53	104.81
LightenNet [39]	13.17	0.55	140.92
MBLLEN [23]	15.02	0.52	119.14
Wang <i>et al.</i> [6]	13.52	0.49	142.01
EnlightenGAN [9]	16.21	0.59	102.78
Zero-DCE	16.57	0.59	98.78
Zero-DCE++	16.42	0.58	102.87

Then, we conduct an ablation study to compare the network structures between Zero-DCE and Zero-DCE++ by replacing the modified component. The ablated models include the Zero-DCE with the depthwise separable convolutions (denoted as DSconv) and the Zero-DCE that shares the curve parameter maps in different iteration stages (denoted as Pshared). The input of Zero-DC is the original resolution image while  $12\times$  downsampling operation as default is used in “DSconv” and “Pshared”. The quantitative comparison results of the ablated models are presented in Table 4.

As shown in Table 4, both “DSconv” and “Pshared” have fewer trainable parameters and FLOPs than Zero-DCE. Introducing “DSconv” and “Pshared” slightly decreases the PSNR values. The trainable parameters and FLOPs of the combination of “DSconv” and “Pshared” (*i.e.*, Zero-DCE++) are significantly decreased with negligible decrease in PSNR value. The results suggest the effectiveness of such modifications. In Figure 13, Zero-DCE still outperforms Zero-DCE++ in some challenging cases. For example, Zero-DCE can more effectively handle challenging lighting without introducing over-/under-exposure when compared with Zero-DCE++. One could choose between Zero-DCE and Zero-DCE++ according to the specific requirements on quality and efficiency.

## 5.4 Benchmark Evaluations

In this section, we conduct qualitative and quantitative experiments to compare different methods. We also investigate the performance of different methods on face detection in the dark.

### 5.4.1 Visual and Perceptual Comparisons

We present the visual comparisons on typical low-light images in Figure 14. For challenging back-lit regions (*e.g.*, the face in Figure 14(a)), Zero-DCE yields natural exposure and clear details

while SRIE [18], LIME [19], LightenNet [39], MBLLEN [23], Wang *et al.* [6], and EnlightenGAN [9] cannot recover the face clearly. In comparison, RetinexNet [7] produces over-exposed artifacts. In the second example featuring an indoor scene, our method enhances dark regions and preserves color of the input image simultaneously. The result is visually pleasing without obvious noise and color casts. In contrast, Li *et al.* [20] and MBLLEN [23] over-smooth the details while other baseline methods amplify noise and even produce color deviation (*e.g.*, the color of wall). Overall, Zero-DCE++ obtains the comparable performance to Zero-DCE in both indoor and outdoor scenes.

We also show the results of different methods on the image sampled from the Part2 subset testing set. The comparison results are presented in Figure 15. Compared with the results of other methods that remain the low-light regions or introduce obvious artifacts, the proposed Zero-DCE and Zero-DCE++ not only produce more clear details but also do not introduce blocking artifacts. Our method tends to generate the results with proper contrast, clear details, vivid color, and less noise.

We perform a user study to quantify the subjective visual quality of various methods. We process low-light images from the image sets (NPE, LIME, MEF, DICM, VV) by different methods. For each enhanced result, we display it on a screen and provide the input image as a reference. A total of 15 human subjects are invited to independently score the visual quality of the enhanced image. These subjects are trained by observing the results from

- 1) whether the results contain over-/under-exposed artifacts or over-/under-enhanced regions;
- 2) whether the results introduce color deviation; and
- 3) whether the results have unnatural texture and obvious noise.

The scores of visual quality range from 1 to 5 (worst to best quality). The step size is set to 1. The average subjective scores for each image set are reported in Table 5. As summarized in Table 5, Zero-DCE achieves the highest average User Study (US) score for a total of 199 testing images from the above-mentioned image sets while Zero-DCE++ achieves the second-highest US score. The Zero-DCE and Zero-DCE++ obtain similar subject scores, which further indicates that the effectiveness and robustness of Zero-DCE++. For the MEF, DICM, and VV sets, our results are most favored by the subjects. All in all, the user study demonstrates that our method can produce a better performance on diverse low-light images from the human subjective visual perspective.

In addition to the US score, we employ a non-reference perceptual index (PI) [44], [45], [46] to evaluate the perceptual quality. The PI metric is originally used to measure perceptual quality in image super-resolution. It has also been used to assess the performance of other image restoration tasks, such as image

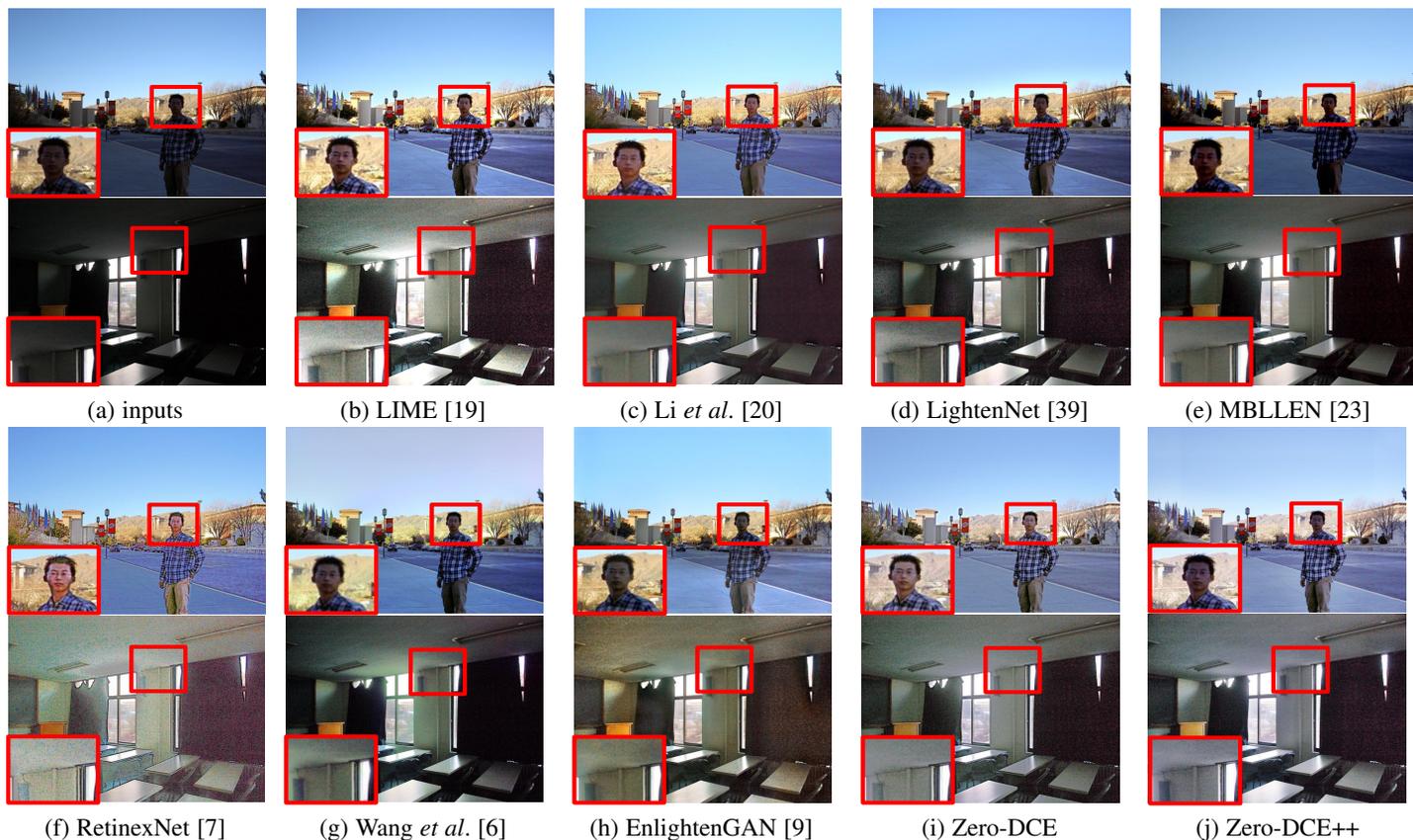


Fig. 14: Visual comparisons on typical low-light images.

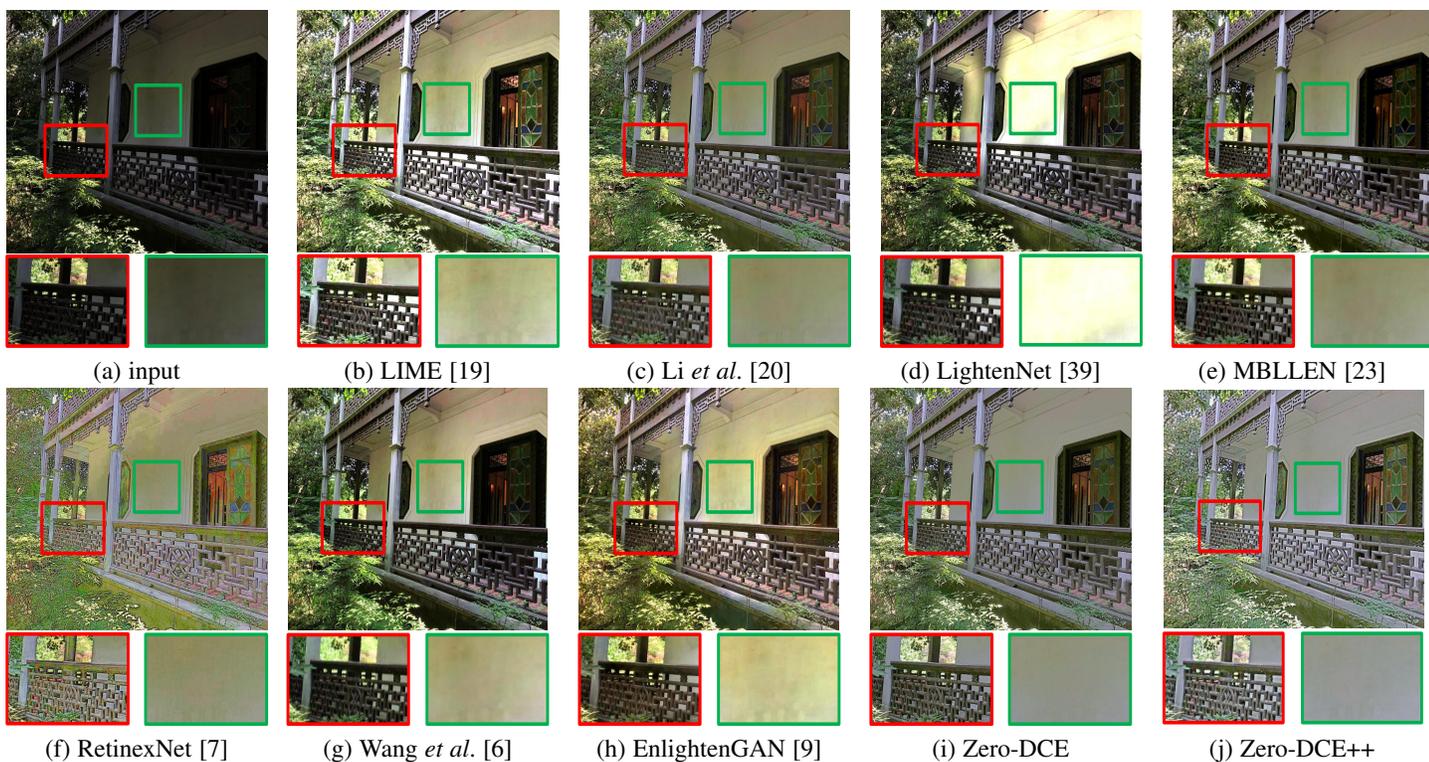


Fig. 15: Visual comparisons on a low-light image sampled from the Part2 subset testing set.

TABLE 7: Runtime (RT, in second), trainable parameters (#P), and FLOPs (in G) comparisons. “-” indicates that the result is not available. The best result is in red whereas the second best one is in blue under each case.

Method	RT	#P	FLOPs	Platform
SRIE [18]	12.1865	-	-	MATLAB (CPU)
LIME [19]	0.4914	-	-	MATLAB (CPU)
Li <i>et al.</i> [20]	90.7859	-	-	MATLAB (CPU)
LightenNet [39]	25.7716	29,532	30.54	MATLAB (CPU)
MBLLEN [23]	13.9949	450,171	301.12	TensorFlow (GPU)
RetinexNet [7]	0.1200	555,205	587.47	TensorFlow (GPU)
Wang <i>et al.</i> [6]	0.0210	998,816	0.19	TensorFlow (GPU)
EnlightenGAN [9]	0.0078	8,636,675	273.24	PyTorch (GPU)
Zero-DCE	0.0025	79,416	84.99	PyTorch (GPU)
Zero-DCE++	0.0012	10,561	0.12	PyTorch (GPU)

dehazing [47]. A lower PI value indicates better perceptual quality. The PI values are reported in Table 5 too. Similar to the user study, the proposed Zero-DCE is superior to other competing methods in terms of the average PI values. It obtains the best perceptual quality on LIME, MEF, and DICM sets. Zero-DCE++ also produces competing average PI values.

#### 5.4.2 Quantitative Comparisons

We employ the full-reference image quality assessment metrics PSNR, SSIM [43], and MAE metrics to quantitatively compare the performance of different methods on the Part2 testing set. A higher SSIM value indicates a result is closer to the ground truth in terms of structural properties. A higher PSNR (lower MAE) value indicates a result is closer to the ground truth in terms of pixel-level image content. In Table 6, the proposed Zero-DCE achieves the best values under all cases, despite that it does not use any paired or unpaired training data. In contrast, Zero-DCE++ obtains comparable performance to Zero-DCE, such as the second-best quantitative scores of PSNR and SSIM values on Part2 testing set.

Our Zero-DCE is computationally efficient, benefited from the simple curve mapping form and lightweight network structure. Further, Zero-DCE++ extremely speeds up Zero-DCE and only costs few computational resources. Table 7 shows the runtime<sup>4</sup>, trainable parameters, and FLOPs of different methods averaged on 32 images of size 1200×900×3. For conventional methods and LightenNet [39], only the codes of CPU version are available.

Compared with current methods, our method achieves the fastest runtime with a large margin (*i.e.*, Zero-DCE: 0.0025s and Zero-DCE++: 0.0012s). Moreover, the runtime of Zero-DCE++ is only 0.0012s, which is really faster than current methods. The runtime of Zero-DCE++ is 17.5 times and 6.5 times faster than recent deep learning-based methods Wang *et al.* [6] and EnlightenGAN [9], respectively. Zero-DCE++ only contains a tiny network structure that has 10,561 trainable parameters and costs 0.12G FLOPs, which are extremely suitable for practical applications.

#### 5.4.3 Face Detection in the Dark

We investigate the performance of low-light image enhancement methods on the face detection task under low-light conditions. Specifically, we use the DARK FACE dataset [42] that composes 10,000 images taken in the dark. Since the bounding boxes of test set are not publicly available, we perform an evaluation on the training and validation sets, which totally consists of 6,000 images.

4. Runtime is measured on a PC with an Nvidia GTX 2080Ti GPU and Intel I7 6700 CPU, except for Wang *et al.* [6], which has to run on GTX 1080Ti GPU.

TABLE 8: The average precision (AP) for face detection in the dark under different IoU thresholds (0.5, 0.7, 0.9). The best result is in red whereas the second best one is in blue under each case.

Method	IoU thresholds		
	0.5	0.7	0.9
input	0.231278	0.007296	0.000002
SRIE [18]	0.288193	0.012621	0.000007
LIME [19]	0.293970	0.013417	0.000007
Li <i>et al.</i> [20]	0.243714	0.008616	0.000003
LightenNet [39]	0.290128	0.012581	0.000005
MBLLEN [23]	0.289232	0.013696	0.000007
RetinexNet [7]	0.304933	0.017545	0.000005
Wang <i>et al.</i> [6]	0.280068	0.011107	0.000003
EnlightenGAN [9]	0.276574	0.013204	0.000009
Zero-DCE	0.303135	0.014772	0.000005
Zero-DCE++	0.297977	0.014587	0.000005

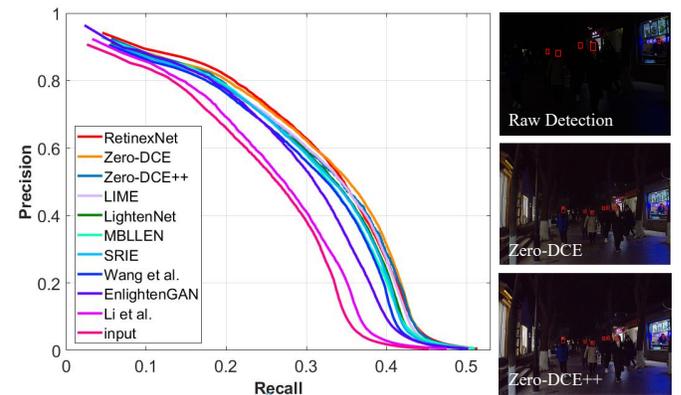


Fig. 16: The P-R curves of face detection in the dark. Best viewed on a color screen in high resolution with zoom in.

A state-of-the-art deep face detector, Dual Shot Face Detector (DSFD) [48], trained on WIDER FACE dataset [49], is used as the baseline model. We feed the results of different low-light image enhancement methods to DSFD [48]. We depict the precision-recall (P-R) curves under IoU threshold 0.5 in Figure 16 and compare the average precision (AP) under different IoU thresholds (*i.e.*, 0.5, 0.7, 0.9) using the evaluation tool<sup>5</sup> provided in DARK FACE dataset [42]. The AP results are presented in Table 8.

As shown in Figure 16, after image enhancement, the precision of DSFD [48] increases considerably compared to that using the input images without enhancement. Among different methods, RetinexNet [7], Zero-DCE, and Zero-DCE++ perform the best. These three methods are comparable but Zero-DCE and Zero-

5. [https://github.com/Ir1d/DARKFACE\\_eval\\_tools](https://github.com/Ir1d/DARKFACE_eval_tools)

DCE++ perform better in the high recall area. As presented in Table 8, the AP scores of all methods drop when we set higher IoU thresholds. When the IoU threshold is set to 0.9, the performance of all methods is extremely poor. Under the IoU thresholds of 0.5 and 0.7, Zero-DCE and Zero-DCE++ obtain similar AP scores that are just a little lower than the best result produced by RetinexNet [7]. However, the subjective and quantitative results of RetinexNet [7] are unsatisfactory as shown before. In contrast, our method does not require paired training data and balances the subject enhancement performance, application performance, and computational cost well. Observing the examples, our Zero-DCE and Zero-DCE++ lighten up the faces in the extremely dark regions and preserves the well-lit regions, thus improves the performance of face detector in the dark.

## 6 CONCLUSION

We proposed a deep network for low-light image enhancement. It can be trained end-to-end with zero reference images. This is achieved by formulating the low-light image enhancement task as an image-specific curve estimation problem, and devising a set of differentiable non-reference losses. By re-designing the network structure, reformulating the curve estimation, and controlling the sizes of input image, the proposed Zero-DCE can be further improved, which is significant light-weight and fast for practical applications. Our method excels in both enhancement performance and efficiency. Experiments demonstrate the superiority of our method against existing light enhancement methods.

## ACKNOWLEDGMENTS

This research was conducted in collaboration with SenseTime. This work is supported by A\*STAR through the Industry Alignment Fund - Industry Collaboration Projects Grant. It is also partially supported by Singapore MOE AcRF Tier 1 (2018-T1-002-056) and NTU SUG. Chunle Guo is sponsored by CAAI-Huawei MindSpore Open Fund.

## REFERENCES

- [1] J. Pan, D. Sun, H. Pfister, and M. H. Yang, "Deblurring images via dark channel prior," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 40, no. 10, pp. 2315–2328, 2018.
- [2] W. S. Lai, J. B. Huang, N. Ahuja, and M. H. Yang, "Fast and accurate image super-resolution with deep laplacian pyramid networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 41, no. 11, pp. 2599–2613, 2019.
- [3] S. Gu, S. Guo, W. Zuo, Y. Chen, R. Timofte, L. V. Gool, and L. Zhang, "Learned dynamic guidance for depth image reconstruction," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 10, pp. 2437–2452, 2020.
- [4] C. Li, C. Guo, W. Ren, R. Cong, J. Hou, S. Kwong, and D. Tao, "An underwater image enhancement benchmark dataset and beyond," *IEEE Transactions on Image Processing*, vol. 29, pp. 4376–4389, 2019.
- [5] C. Li, C. Guo, J. Guo, P. Han, H. Fu, and R. Cong, "PDR-Net: Perception-inspired single image dehazing network with refinement," *IEEE Transactions on Multimedia*, vol. 22, no. 3, pp. 704–716, 2019.
- [6] 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.
- [7] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," in *BMVC*, 2018.
- [8] K. Xu, X. Yang, B. Yin, and R. W. H. Lau, "Learning to restore low-light images via decomposition-and-enhancement," in *CVPR*, 2020, pp. 2281–2290.
- [9] 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," 2019, arXiv arXiv:1906.06972.
- [10] 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.
- [11] 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 *CVPR*, 2020, pp. 1780–1789.
- [12] D. Coltuc, P. Bolon, and J.-M. Chassery, "Exact histogram specification," *IEEE Transactions on Image Processing*, vol. 15, no. 5, pp. 1143–1152, 2006.
- [13] H. Ibrahim and N. S. P. Kong, "Brightness preserving dynamic histogram equalization for image contrast enhancement," *IEEE Transactions on Consumer Electronics*, vol. 53, no. 4, pp. 1752–1758, 2007.
- [14] J. A. Stark, "Adaptive image contrast enhancement using generalizations of histogram equalization," *IEEE Transactions on Image Processing*, vol. 9, no. 5, pp. 889–896, 2000.
- [15] C. Lee, C. Lee, and C.-S. Kim, "Contrast enhancement based on layered difference representation of 2d histograms," *IEEE Transactions on Image Processing*, vol. 22, no. 12, pp. 5372–5384, 2013.
- [16] E. H. Land, "The retinex theory of color vision," *Scientific American*, vol. 237, no. 6, pp. 108–128, 1977.
- [17] S. Wang, J. Zheng, H.-M. Hu, and B. Li, "Naturalness preserved enhancement algorithm for non-uniform illumination images," *IEEE Transactions on Image Processing*, vol. 22, no. 9, pp. 3538–3548, 2013.
- [18] 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.
- [19] 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, 2017.
- [20] 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.
- [21] L. Yuan and J. Sun, "Automatic exposure correction of consumer photographs," in *ECCV*, 2012, pp. 771–785.
- [22] 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.
- [23] F. Lv, F. Lu, J. Wu, and C. Lim, "MBLLEN: Low-light image/video enhancement using cnns," in *BMVC*, 2018.
- [24] V. Bychkovsky, S. Paris, E. Chan, and F. Durand, "Learning photographic global tonal adjustment with a database of input/output image pairs," in *CVPR*, 2011, pp. 97–104.
- [25] C. Chen, Q. Chen, J. Xu, and K. Vladlen, "Learning to see in the dark," in *CVPR*, 2018, pp. 3291–3300.
- [26] C. Chen, Q. Chen, M. N. Do, and V. Koltun, "Seeing motion in the dark," in *ICCV*, 2019, pp. 3185–3194.
- [27] W. Ren, S. Liu, L. Ma, Q. Xu, X. Xu, X. Cao, J. Du, and M.-H. Yang, "Low-light image enhancement via a deep hybrid network," *IEEE Transactions on Image Processing*, vol. 28, no. 9, pp. 4364–4375, 2019.
- [28] Y. Zhang, J. Zhang, and X. Guo, "Kindling the darkness: A practical low-light image enhancer," in *ACMMM*, 2019, pp. 1632–1640.
- [29] 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 *CVPR*, 2020, pp. 3063–3072.
- [30] T. Mertens, J. Kautz, and F. V. Reeth, "Exposure fusion," in *PCCGA*, 2007.
- [31] Mertens, J. Kautz, and F. V. Reeth, "Exposure fusion: A simple and practical alternative to high dynamic range photography," *Computer Graphics Forum*, vol. 28, no. 1, pp. 161–171, 2009.
- [32] G. Buchsbaum, "A spatial processor model for object colour perception," *J. Franklin Institute*, vol. 310, no. 1, pp. 1–26, 1980.
- [33] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *CVPR*, 2017, pp. 1251–1258.
- [34] L. C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *ECCV*, 2018, pp. 801–818.
- [35] H. Hu, J. Gu, Z. Zhang, J. Dai, and Y. Wei, "Relation networks for object detection," in *CVPR*, 2018, pp. 3588–3597.
- [36] C. Liu, L. C. Chen, F. Schroff, H. Adam, W. Hua, A. Yuille, and L. F. Fei, "Auto-deeplab: hierarchical neural architecture search for semantic image segmentation," in *CVPR*, 2019, pp. 82–92.
- [37] Y. Chen, Y. Wang, M. Kao, and Y. Chuang, "Deep photo enhancer: Unpaired learning for image enhancement from photographs with gans," in *CVPR*, 2018, pp. 6306–6314.
- [38] J. Cai, S. Gu, and L. Zhang, "Learning a deep single image contrast enhancer from multi-exposure image," *IEEE Transactions on Image Processing*, vol. 27, no. 4, pp. 2049–2026, 2018.

- [39] C. Li, J. Guo, F. Porikli, and Y. Pang, "LightenNet: a convolutional neural network for weakly illuminated image enhancement," *Pattern Recognition Letters*, vol. 104, pp. 15–22, 2018.
- [40] K. Ma, K. Zeng, and Z. Wang, "Perceptual quality assessment for multi-exposure image fusion," *IEEE Transactions on Image Processing*, vol. 24, no. 11, pp. 3345–3356, 2015.
- [41] C. Lee, C. Lee, and C.-S. Kim, "Contrast enhancement based on layered difference representation," in *ICIP*, 2012, pp. 965–968.
- [42] Y. Yuan, W. Yang, W. Ren, J. Liu, W. J. Scheirer, and W. Zhangyang, "UG+ Track 2: A collective benchmark effort for evaluating and advancing image understanding in poor visibility environments," 2019, arXiv:1904.04474.
- [43] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: From error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [44] Y. Blau and T. Michaeli, "The perception-distortion tradeoff," in *CVPR*, 2018, pp. 6228–6237.
- [45] C. Ma, C.-Y. Yang, X. Yang, and M.-H. Yang, "Learning a no-reference quality metric for single-image super-resolution," *Computer Vision and Image Understanding*, vol. 158, pp. 1–16, 2017.
- [46] 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, 2013.
- [47] Y. Qu, Y. Chen, J. Huang, and Y. Xie, "Enhanced pix2pix dehazing network," in *CVPR*, 2019, pp. 8160–8168.
- [48] J. Li, Y. Wang, C. Wang, Y. Tai, J. Qian, J. Yang, C. Wang, J. Li, and F. Huang, "Dsfed: Dual shot face detector," in *CVPR*, 2019, pp. 5060–5069.
- [49] S. Yang, P. Luo, C.-C. Loy, and X. Tang, "WIDER FACE: A face detection benchmark," in *CVPR*, 2016, pp. 5525–5533.



**Chen Change Loy** (Senior Member, IEEE) received the PhD degree in computer science from the Queen Mary University of London, in 2010. He is an associate professor with the School of Computer Science and Engineering, Nanyang Technological University. Prior to joining NTU, he served as a research assistant professor with the Department of Information Engineering, The Chinese University of Hong Kong, from 2013 to 2018. His research interests include computer vision and deep learning. He serves as an associate editor of the *IEEE Transactions on Pattern Analysis and Machine Intelligence* and the *International Journal of Computer Vision*. He also serves/served as an Area Chair of CVPR 2021, CVPR 2019, ECCV 2018, AAAI 2021 and BMVC 2018-2020.



**Chongyi Li** received the Ph.D. degree from the School of Electrical and Information Engineering, Tianjin University, Tianjin, China, in June 2018. From 2016 to 2017, he was a joint-training Ph.D. Student with Australian National University, Australia, under the supervision of Prof. Fatih Porikli. He was a postdoctoral fellow with the Department of Computer Science, City University of Hong Kong, working with Chair Prof. Sam Kwong. He is currently a research fellow with the School of Computer Science and En-

gineering, Nanyang Technological University (NTU), Singapore. His current research focuses on image processing, computer vision, and deep learning, particularly in the domains of image restoration and enhancement.



**Chunle Guo** received his PhD degree from Tianjin University in China under the supervision of Prof. Jichang Guo. He conducted the Ph.D. research as a Visiting Student with the School of Electronic Engineering and Computer Science, Queen Mary University of London (QMUL), UK. He continued his research as a Research Associate with Department of Computer Science, City University of Hong Kong (CityU), from 2018 to 2019. Now he is a postdoc research fellow working with Prof. Ming-Ming Cheng in Nankai

University. His research interests lies in image processing, computer vision, and deep learning.