# Decoupled Cross-Scale Cross-View Interaction for Stereo Image Enhancement in The Dark

Huan Zheng[1], Zhao Zhang[1*], Jicong Fan[2], Richang Hong[1], Yi Yang[3], and Shuicheng Yan[4]

[1]Hefei University of Technology, China
[2]The Chinese University of Hong Kong (Shenzhen), China
[3]Zhejiang University, China
[4]Sea AI Lab, Singapore

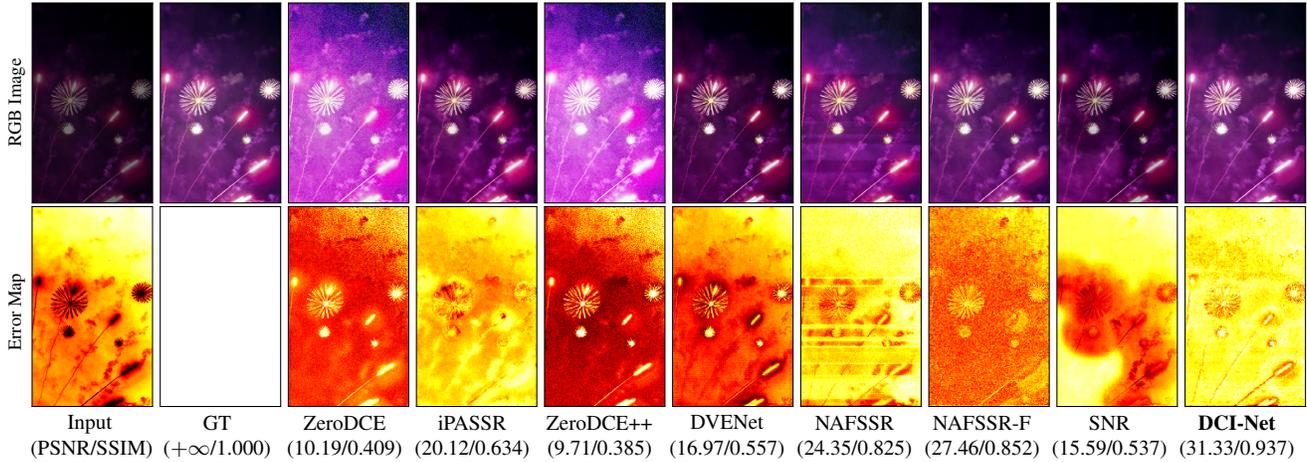| Input (PSNR/SSIM) | GT ($+\infty$/1.000) | ZeroDCE (10.19/0.409) | iPASSR (20.12/0.634) | ZeroDCE++ (9.71/0.385) | DVENet (16.97/0.557) | NAFSSR (24.35/0.825) | NAFSSR-F (27.46/0.852) | SNR (15.59/0.537) | **DCI-Net** (31.33/0.937) |

Figure 1. Visualization of the enhanced images and the corresponding error maps of each method based on Flickr1024 dataset, including NAFSSR [3], NAFSSR-F [3], iPASSRNet [43], SNR [48], DVENet [13], ZeroDCE [10], ZeroDCE++ [20] and our DCI-Net. Whiter and brighter pixels in the error maps indicate smaller errors. It is clear that other compared methods obtain darker pixels in the error maps than our DCI-Net, which means that our method is capable of preserving consistent color and recovering the textures more accurately.

## Abstract

*Low-light stereo image enhancement (LLSIE) is a relatively new task to enhance the quality of visually unpleasant stereo images captured in dark condition. However, current methods achieve inferior performance on detail recovery and illumination adjustment. We find it is because: 1) the insufficient single-scale inter-view interaction makes the cross-view cues unable to be fully exploited; 2) lacking long-range dependency leads to the inability to deal with the spatial long-range effects caused by illumination degradation. To alleviate such limitations, we propose a LLSIE model termed Decoupled Cross-scale Cross-view Interaction Network (DCI-Net). Specifically, we present a decoupled interaction module (DIM) that aims for sufficient dual-view information interaction. DIM decouples the dual-view information exchange into discovering multi-scale cross-view correlations and further exploring cross-scale information flow. Besides, we present a spatial-channel information mining block (SIMB) for intra-view feature extraction, and the benefits are twofold. One is the long-range dependency capture to build spatial long-range relationship, and the other is expanded channel information refinement that enhances information flow in channel dimension. Extensive experiments on Flickr1024, KITTI 2012, KITTI 2015 and Middlebury datasets show that our method obtains better illumination adjustment and detail recovery, and achieves SOTA performance compared to other related methods. Our codes, datasets and models will be publicly available.*

## 1. Introduction

Single image processing and understanding have made great achievements across a wide range of application areas, such as image classification [12], object detection [33] and semantic segmentation [29]. Recently, with the growing application of dual cameras, stereo vision has attracted much attention in various fields, e.g., mobile phones and autonomous driving cars [23]. However, the stereo images captured in dark environments usually suffer from low-
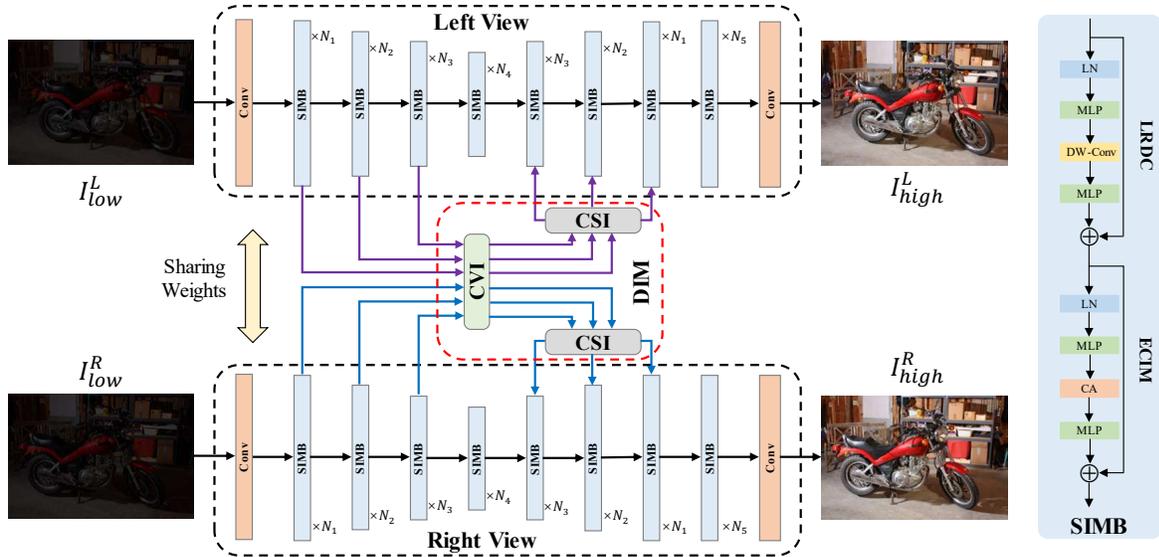
Figure 2. The overall framework of our proposed DCI-Net for LLSIE, which contains two weights-shared branches to process left and right views respectively. Besides, DCI-Net includes two main modules, i.e., DIM and SIMB. To be specific, DIM completes sufficient inter-view information interaction and flow across two views, which includes cross-view interaction (CVI) and cross-scale interaction (CSI); SIMB enhances intra-view feature representation, whose structure is shown on the right.

contrast, weak illumination and various noise [55]. As a consequence, there are obviously negative effects on subsequent high-vision tasks. Hence, low-light stereo image enhancement (LLSIE) is proposed to enhance dark stereo images [14]. To be specific, LLSIE is a task of enhancing the illumination and recovering the hidden details in the dark, via utilizing the stereo images from left and right views.

Compared with LLSIE, single low-light image enhancement (LLIE) methods aim to refine the illumination of single-view images in the dark, which can be divided into traditional and deep learning-based ones [19]. Traditional LLIE methods build prior-based optimization models to adjust the illumination and enhance the contrast [11]. While these methods are relatively simple or highly rely on the hand-crafted priors, which may cause low-quality enhanced results [26]. With the great development of deep learning, convolutional neural network (CNN)-based methods have achieved impressive performance in various low-level vision tasks [5, 7, 45, 56, 59, 60]. More importantly, CNN-based deep models also show superior capability for LLIE [2, 21, 27, 41, 42, 44, 47, 54]. These deep LLIE methods use CNN as backbones to establish a neural network to learn a map from low-light image to normal-light image.

Stereo image restoration is the task of recovering high-quality stereo images from diverse degradations. In comparison to single image restoration, parallax in stereo image pairs is a key point for stereo image restoration. Some recent methods have been proposed to restore the lost information using the correlations between two views [3, 13, 39, 43]. iPASSRNet is firstly proposed for stereo image super-

resolution via exploring symmetry cues between two views [43]. Inspired by iPASSRNet, DVENet is presented, which is the most representative method exploring the LLSIE task [13]. However, the generated illumination-improved images are still unsatisfactory for both low-light enhancement and stereo image restoration methods, as shown in Fig. 1. Therefore, we ask: *what makes the enhanced stereo images suffer from undesired and inaccurate contents*? We attempt to answer this question from two respects:

(1) **Insufficient dual-view information interaction**. The single LLIE methods do not consider the relationship between two views as stereo image at all. In contrast, stereo image restoration methods clearly need to exploit the cues between stereo image pairs. Nevertheless, current dual-view information interaction strategies are still weak. Because existing methods only explore the cross-view correlations at single scale, while ignoring the cues at different scales and missing cross-scale interaction, causing the inability of achieving sufficient cross-scale cross-view interaction.

(2) **Lack of long-range dependency in intra-view learning**. Current methods usually adopt CNNs with a kernel size of 3×3 to build a neural network for image restoration. However, small kernel size design may limit the learning ability of CNNs, because the convolutional operation can only extract information from small regions, and prevents capturing long-range dependencies. Poor illumination has a great influence on the entire image. To handle the spatial long-range effects caused by the degradation, it is important to build

the long-range relationship for LLSIE.

In this paper, we therefore explore effective strategies to facilitate interaction and enhance the stereo images in the dark, and propose a decoupling strategy to complete cross-scale cross-view information interaction. The main contributions of this paper are summarized as follows:

(1) **DCI-Net: LLSIE by Decoupled Cross-scale Cross-view Interaction Network**. We propose DCI-Net to address the issue of weak cross-view information interaction for LLSIE. Specifically, DCI-Net aims at improving the enhancement process by refining both intra-view feature extraction and cross-view information interaction. Experiments on Flickr1024, KITTI 2012, KITTI 2015 and Middlebury datasets demonstrate that our method can better adjust the illumination, recover the details and obtain SOTA performance.

(2) **Decoupled Interaction Module (DIM)**. To enable sufficient dual-view information interaction, namely, cross-scale cross-view information interaction, DIM decouples the above process into two levels, i.e., cross-view interactions at multiple scales and further cross-scale interaction. The first level aims at discovering multi-scale cross-view cues, and the second level focuses on exploring cross-scale information flow for further interaction. Hence, DIM can make full use of the correlations between stereo image pairs.

(3) **Spatial-channel Information Mining Block (SIMB)**. We design the novel module SIMB for intra-view feature extraction. To be specific, SIMB is based on the structure of vision transformer (ViT) so that it can possess a strong learning ability, but our core idea departs from ViT. Instead of using multi-head self-attention, large-kernel design is incorporated into the process of long-range dependency capture (LRDC) for discovering spatial long-range correlations. In addition, expanded channel information refinement (ECIR) is also developed for enhancing channel information flow.

## 2. Related Work

We briefly review the recent progress on single low-light image enhancement and stereo image restoration.

### 2.1. Low-light single image enhancement

For traditional LLIE methods, we mainly introduce the retinex-based and histogram equalization (HE)-based ones. Inspired by the retinex theory [17], retinex models are proposed to decompose the low-light images and reconstruct the normal-light images [1, 22, 35]. HE-based methods aim at adjusting the dynamic range of the low-light image to enhance the contrast [18]. Deep LLIE methods can be further fallen into end-to-end and retinex-based modes. To be specific, end-to-end method directly learns a map to reverse the

illumination degradation, which takes the low-light image as input and directly outputs the enhanced normal-light image [10, 20, 30, 34, 48, 49, 58]; Retinex-based deep methods decouple an image into the illumination map and reflectance map [27, 44, 46, 53]. Differently, deep retinex-based methods employ deep neural networks for image decomposition, in comparison to the traditional retinex-based methods.

### 2.2. Stereo image restoration

Recently, stereo vision has been attracting much attention. A few stereo image restoration methods are also studied. For stereo image super-resolution, Jeon et al. [15] proposed the first work that uses the shift operation to compensate for the parallax between two views. Wang et al. [43] developed a novel parallax attention to make full use of the correlations between stereo image pairs. The latest state-of-the-art method is NAFSSR, which is the champion of the NTIRE 2022 Stereo Image Super-resolution Challenge [3, 38]. For stereo image deraining, Zhang et al. [52] incorporated semantic priors into deraining process for better rain removal. For stereo deblurring, Zhou et al. [61] presented a depth-aware and view aggregated method. Li et al. [24] delivered a novel stereo image debluring model by exploring the dual-pixel alignment. For stereo image dehazing, Nie et al. [32] proposed SRDNet which aims at better exploiting the stereo information from two views. There are also few works for LLSIE [13, 16, 25]. Specifically, DVENet [13] is the most representative method, which incorporates retinex theory into the overall framework in a coarse-to-fine manner. Besides, parallax attention model is also used to explore the correlations between two views.

## 3. Proposed Method

We introduce the proposed DCI-Net in detail in this section. We first illustrate the overall architecture of DCI-Net. Then, we describe the detailed structures of the designed modules. Finally, the used loss functions are discussed.

### 3.1. Overall framework

An overview of the proposed DCI-Net is shown in Fig. 2. Clearly, our model takes a pair of low-light stereo images as input, enhances the illumination of both views, and outputs the enhanced normal-light stereo images. The pipeline of DCI-Net can be divided into three stages: shallow feature extraction, deep feature extraction and stereo image reconstruction. To be specific, we use two convolutional layers in the head and tail, where the first one extracts shallow features and the last one reconstructs the enhanced normal-light stereo images. Given a pair of low-light stereo images, the above processes can be formulated as follows:

$$I_{high}^L, I_{high}^R = \text{H}_{\text{SR}}(\text{H}_{\text{DF}}(\text{H}_{\text{SF}}(I_{low}^L, I_{low}^R))), \quad (1)$$

$\otimes$ Matrix multiplication  $\oplus$ Element-wise summation
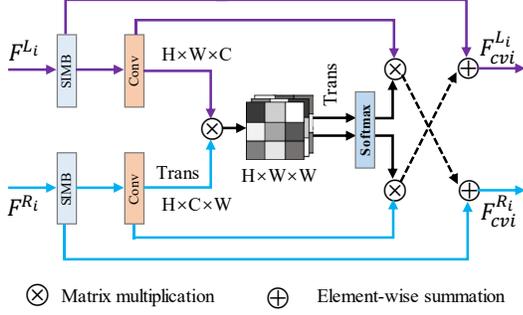
Figure 3. The detailed process of cross-view interaction in DIM. To be specific, CVI explores the cues between both left and right views at multiple scales. Note that only the process at single scale is shown as example.

where $I_{low}^L$, $I_{low}^R$, $I_{high}^L$ and $I_{high}^R$ denote the low-light left-view image, low-light right-view image, enhanced left-view image and right-view image, $H_{SF}(\cdot)$, $H_{DF}(\cdot)$ and $H_{SR}(\cdot)$ denote the transformations for shallow feature extraction, deep feature extraction and stereo image reconstruction respectively. Deep feature extraction can be further fallen into intra-view feature extraction and dual-view interaction. For dual-view interaction, we deliver a decoupled interaction module (DIM) to explore synchronous cross-view and cross-scale interaction. For intra-view feature extraction, we construct a spatial-channel information mining block (SIMB)-based U-Net to obtain stronger feature representation. It is worth noting that the weights of shallow feature extraction, intra-view feature extraction in deep feature extraction and stereo image reconstruction are always shared.

## 3.2. Decoupled interaction module (DIM)

Different from single image processing, one key point of LLSIE is exploring the correlations between two views to promote illumination enhancement. Hence, methods for low-light single image enhancement do not well in enhancing the stereo images, since they only consider one view. Some previous attempts have done to discover and exploit the cues between a pair of stereo images [3, 13, 15, 40, 43]. Nevertheless, these methods lack considering the cross-view interaction at different scales, which makes the cross-scale interaction be ignored. Note that a lot of existing studies on CNN and ViT show the importance of multi-scale information interaction [9, 37]. To alleviate this issue, we propose DIM to decouple cross-scale cross-view information interaction into studying inter-view correlations at multiple scales, and further cross-scale interaction.

**Cross-view interaction at multiple scales**. Previous methods only use the correlations at a single scale [3, 13, 15, 40, 43]. In contrast, we explore cross-view cues at multiple scales. We incorporate a cross-view interaction (CVI) module into different scales for discovering the multi-scale cross-view cues. The detailed structure of CVI is shown
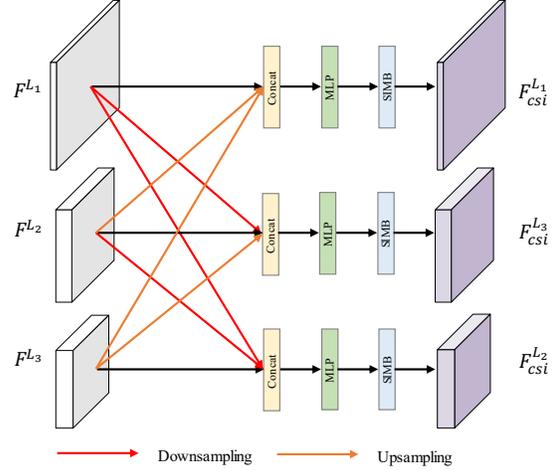


Figure 4. The detailed process of cross-scale interaction in DIM. Specifically, CSI completes the interaction among different scales. Note that only the process of left view is shown.

in Fig. 3. Given the input stereo feature maps from the $i$-th stage, CVI uses the matrix multiplication to compute the cross-view correlations for information interaction. This process can be formulated as follows:

$$M_{cor} = LR^T, \tag{2}$$

where $L \in \mathbb{R}^{H \times W \times C}$ and $R \in \mathbb{R}^{H \times C \times W}$ denote the input stereo feature maps, and $M_{cor} \in \mathbb{R}^{H \times W \times W}$ denotes the correlations matrix. It is noted that there are only horizontal shifts for stereo images. Hence, we mainly pay attention to the horizontal correlations between two views.

**Cross-scale interaction**. CVI has extracted the multi-scale correlations between two views, but the information interaction among different scales is still missing. Cross-scale interaction (CSI) is therefore presented to handle this issue by enhancing the cross-scale information flow. The overall structure of CSI is shown in Fig. 4. Given the left-view multi-scale feature maps $F^{L_1}$, $F^{L_2}$ and $F^{L_3}$ obtained by CVI as an example, CSI firstly uses scaling operations to densely concatenate them, and then utilizes a multilayer perceptron (MLP) to process the concatenated feature maps. The purpose is twofold, one is to reduce channels and the other is to exchange and fuse information in channel dimension. In the end, a spatial-channel information mining block is incorporated into the tail for further spatial information interaction. The above processes can be formulated as

$$F_{csi}^{L_1}, F_{csi}^{L_2}, F_{csi}^{L_3} = \text{SIMB}(\text{MLP}(\text{DC}(F^{L_1}, F^{L_2}, F^{L_3}))) \tag{3}$$

where $\text{DC}(\cdot)$ denotes the densely concatenate operation in CSI, and $F_{csi}^{L_1}$, $F_{csi}^{L_2}$ and $F_{csi}^{L_3}$ denote the cross-scale interacted feature maps.
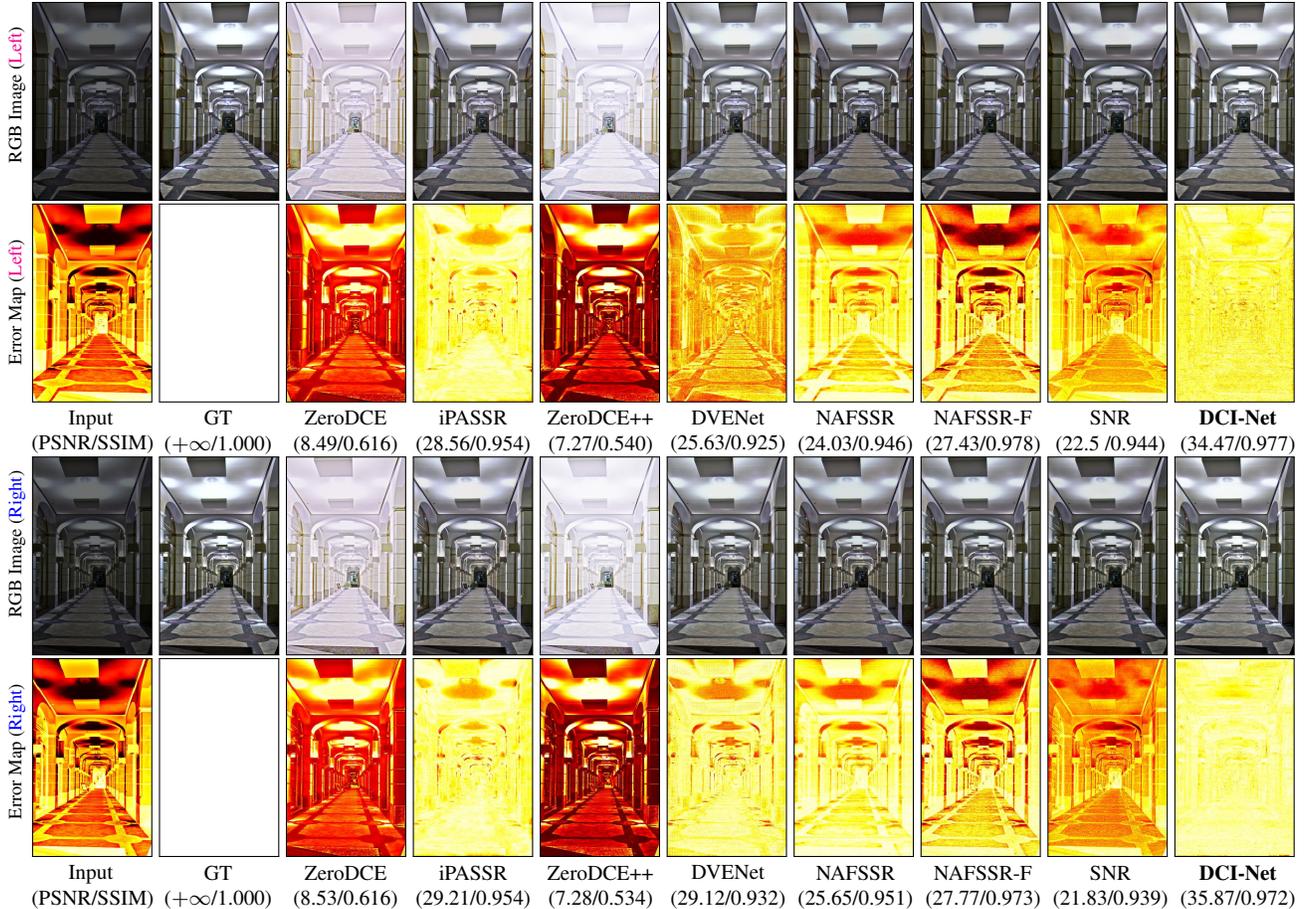
Figure 5. Visualization of the enhanced images and corresponding error maps of each method based on Flickr1024 dataset, including the results of NAFSSR [3], NAFSSR-F [3], iPASSRNet [43], SNR [48], DVENet [13], ZeroDCE [10], ZeroDCE++ [20] and our DCI-Net. Whiter and brighter pixels in the error maps indicate smaller errors. Clearly, our DCI-Net achieves better illumination adjustment and color correction, and obtains smaller errors and superior performance than other methods, which can also be seen from the shown metrics.

## 3.3. Spatial-channel information mining block (SIMB)

Recalling that previous single low-light image enhancement and stereo image restoration methods usually adopt CNNs with a kernel size of 3×3 to construct a deep neural network. As a result, the long-range dependency cannot be well captured. Recently, transformer-based models have achieved impressive performance in diverse computer vision tasks due to the capacity of building long-range dependency [6,28]. But vision transformer is computationally expensive. Some recent works have shown that large kernel convolutional layers can also obtain long-range correlations [4]. Besides, PoolFormer shows that the overall architecture plays an important role in vision transformer [50]. Inspired by these works, we therefore propose SIMB. Although SIMB inherits the structure of vision transformer, it replaces the multi-head self-attention with large kernel convolutional layers, and further explores the information flow in channel dimension. The structure of SIMB is shown on

the right of Fig. 2. As can be seen, three are two core steps in SIMB: long-range dependency capture (LRDC) and expanded channel information refinement (ECIR).

**Long-range dependency capture (LRDC).** To overcome the shortage that CNNs with small kernel size cannot build long-range relationship, we use large kernel convolutional layers for long-range dependency capture. The shifted kernel in CNN can be regarded as the shifted window in Swintransformer [28]. Nevertheless, it is computationally expensive for the vanilla CNN with increased kennel size. Hence, large kernel depth-wise convlolutional (DW-Conv) layer and MLP are used to approach the effect of vanilla CNN. There are two advantages for this design. Firstly, large kernel design can discover long-range correlations; secondly, DW-Conv layer can significantly reduce the computational cost. We set the kernel size of DW-Conv layer to 7 that is consistent with the window size of Swintransformer [28]. From Fig. 2, the transformation of LRDC
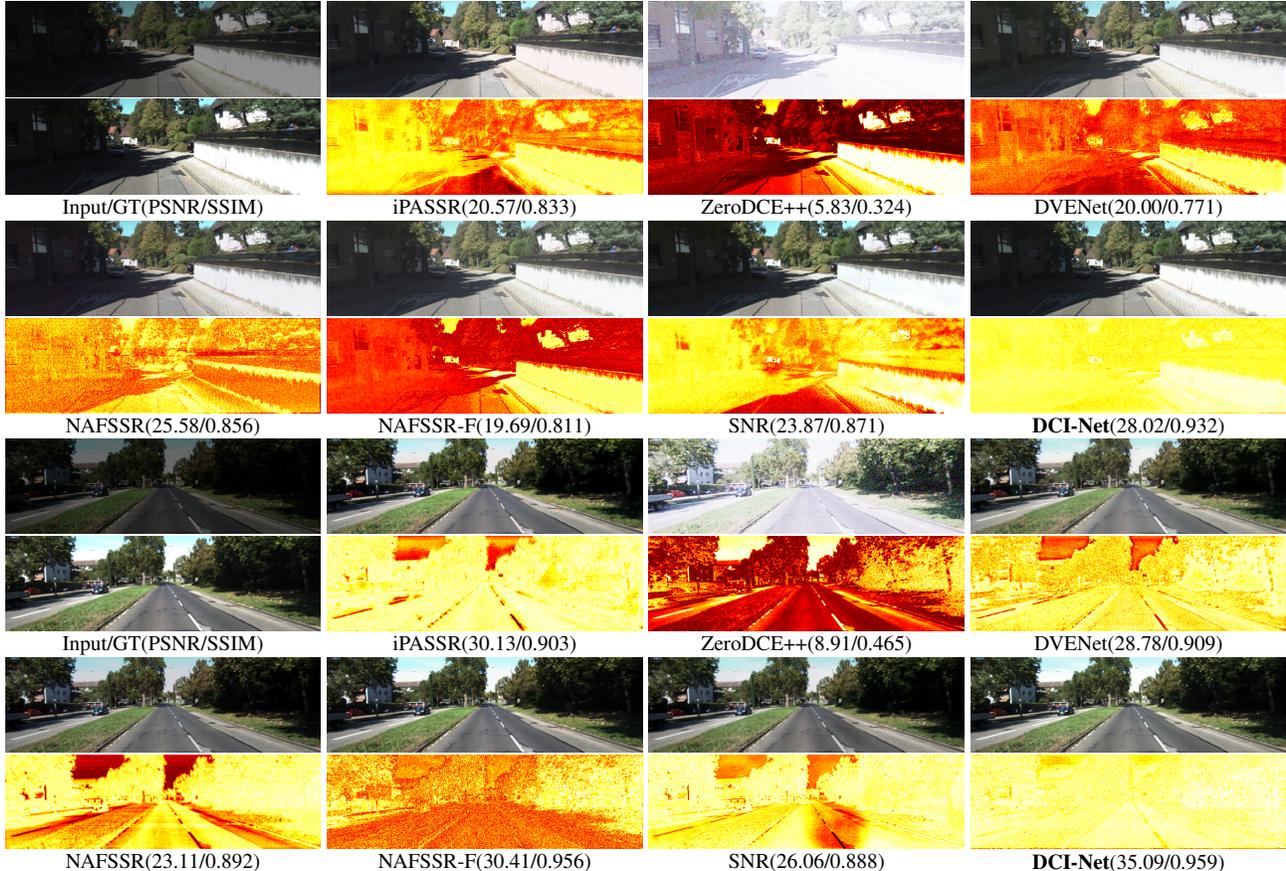
5

Figure 6. Visualization of the enhanced images and corresponding error maps of each method based on the KITTI 2012 and KITTI 2015 datasets, including NAFSSR [3], NAFSSR-F [3], iPASSRNet [43], SNR [48], DVENet [13], ZeroDCE++ [20] and our DCI-Net. Whiter and brighter pixels in the error maps indicate smaller errors. Compared with other methods, our DCI-Net obtains the smallest errors as can been seen from the error maps. This indicates that our proposed DCI-Net can better adjust the illumination and restore the details.

can be described as follows:

$$F_{lr} = \text{MLP}(\text{DWConv}(\text{MLP}(\text{LN}(F)))) + F, \quad (4)$$

where $\text{LN}(\cdot)$, $\text{MLP}(\cdot)$ and $\text{DWConv}(\cdot)$ denote the layer normalization (LN), MLP and DW-Conv, $F$ and $F_{lr}$ represent the input and output feature maps.

**Expanded channel information refinement (ECIR).** In the first step of SIMB, we mainly explore spatial information. But there less attention has been paid to channel information. Hence, we develop ECIR. The core idea of ECIR is simple yet effective, which completes channel information mixing in higher dimensional space. Note that this design can be easily implemented by incorporating channel attention (CA) into the second stage of vision transformer. As shown in Fig. 2, the process of ECIR can be formulated by

$$F_{sb} = \text{MLP}(\text{CA}(\text{MLP}(\text{LN}(F_{lr})))) + F_{lr}, \quad (5)$$

where $\text{CA}(\cdot)$ denotes the transformation of channel attention and $F_{sb}$ is the processed result of SIMB. It should be noted that the first MLP is used to expand the channel of feature maps to the higher dimension, while the second MLP

re-maps the higher dimensional feature maps to the original channel size. Since CA can compress the spatial information and fully focuses on the channel information, it can be used to refine the channel information.

### 3.4. Loss function

The total loss function $\mathcal{L}$ of our DCI-Net contains two losses, i.e., frequency-domain reconstruction loss $\mathcal{L}_{fre}$ and smooth loss $\mathcal{L}_{tv}$, which are illustrated as follows:

$$\mathcal{L} = \mathcal{L}_{fre} + \lambda \mathcal{L}_{tv}, \quad (6)$$

where $\lambda$ is a hyper-parameter that is set to 0.1 in this paper. To be specific, frequency-domain reconstruction loss is used to guide our method to reconstruct the normal-light stereo image, which can be described as

$$\mathcal{L}_{fre} = \|\text{FFT}(I_{high}^L) - \text{FFT}(I_{gt}^L)\|_1 + \\ \|\text{FFT}(I_{high}^R) - \text{FFT}(I_{gt}^R)\|_1, \quad (7)$$

where $I_{gt}^L$ and $I_{gt}^R$ denote the ground-truth stereo normal-light images, and $\text{FFT}(\cdot)$ denotes the fast fourier transform.
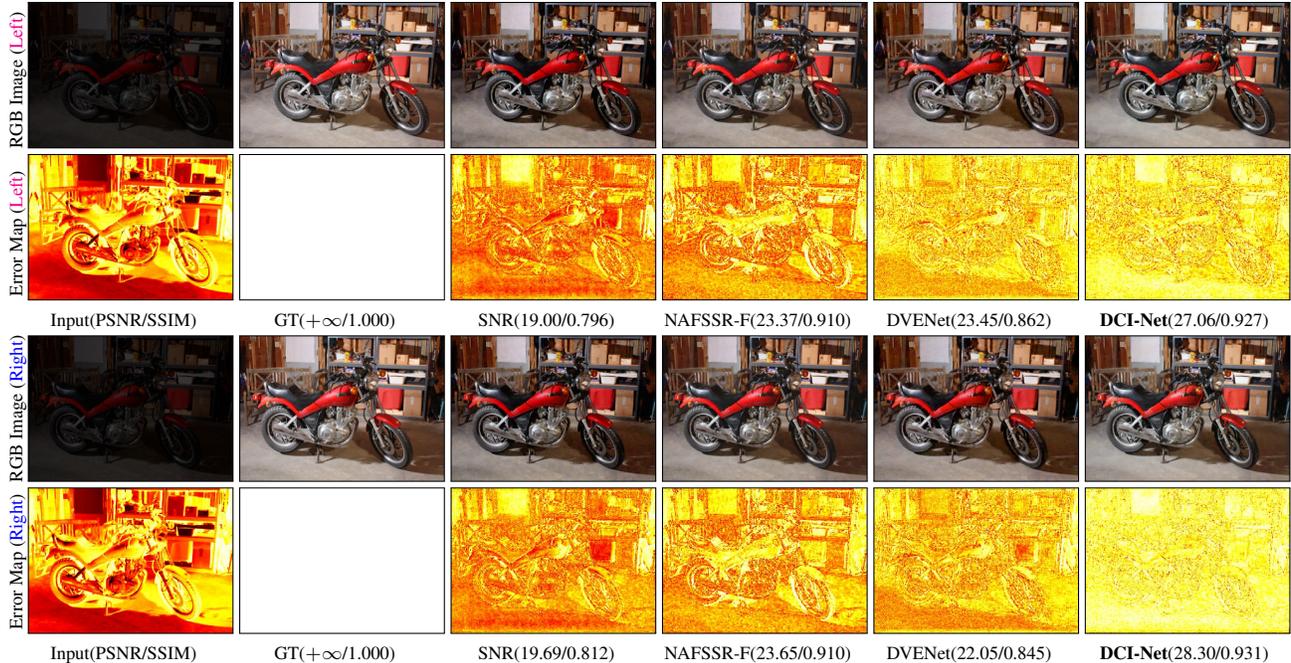
Figure 7. Visual results of each method based on Middlebury dataset, including NAFSSR-F [3], SNR [48], DVENet [13] and our DCI-Net. Whiter and brighter pixels in the error maps indicate smaller errors. We can see that there is a motorcycle in the error maps of all compared methods. But in contrast, it is hard to see the shape of the motorcycle for our DCI-Net, i.e., less information is lost in our method.

Table 1. Evaluation results of each method on Flickr1024, KITTI 2012, KITTI 2015 and Middlebury datasets. Note that PSNR/SSIM values achieved on both the left images (i.e., Left) and a pair of stereo images (i.e., (Left + Right) /2) are reported. The **bold** denotes the best. It is clear that our DCI-Net achieves SOTA performance among all compared methods.

| Method | Venue | Left | | | | (Left + Right) / 2 | | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Flickr1024 | KITTI 2012 | KITTI 2015 | Middlebury | Flickr1024 | KITTI 2012 | KITTI 2015 | Middlebury |
| ZeroDCE [10] | cvpr'20 | 10.49/0.420 | 9.21/0.403 | 10.17/0.441 | 11.34/0.591 | 10.49/0.420 | 9.18/0.401 | 10.17/0.439 | 11.27/0.592 |
| iPASSRNet [43] | cvprw'21 | 22.27/0.777 | 23.32/0.820 | 21.76/0.802 | 21.50/0.872 | 22.20/0.777 | 23.49/0.827 | 21.69/0.804 | 21.67/0.871 |
| ZeroDCE++ [20] | tpami'21 | 10.74/0.418 | 9.15/0.388 | 10.08/0.427 | 10.28/0.563 | 10.74/0.418 | 9.14/0.386 | 10.07/0.425 | 10.19/0.563 |
| DVENet [13] | tmm'22 | 21.82/0.751 | 22.14/0.787 | 21.05/0.767 | 20.75/0.855 | 21.66/0.751 | 21.96/0.787 | 21.05/0.769 | 21.09/0.856 |
| NAFSSR [3] | cvprw'22 | 21.68/0.734 | 21.43/0.775 | 20.68/0.767 | 22.25/0.857 | 21.74/0.735 | 21.70/0.778 | 20.64/0.771 | 22.39/0.858 |
| NAFSSR-F [3] | cvprw'22 | 24.21/0.815 | 24.22/0.856 | 22.29/0.837 | 23.49/0.920 | 24.26/0.816 | 23.93/0.854 | 22.45/0.842 | 23.43/0.921 |
| SNR [48] | cvpr'22 | 21.34/0.749 | 22.11/0.794 | 20.93/0.769 | 21.43/0.851 | 21.40/0.751 | 22.16/0.793 | 20.90/0.770 | 21.65/0.856 |
| DCI-Net | - | **25.47/0.832** | **26.69/0.883** | **26.26/0.862** | **24.27/0.931** | **25.36/0.832** | **26.90/0.884** | **26.03/0.864** | **24.26/0.929** |

The smooth loss is based on the total variation prior to obtain better and smoother results.

## 4. Experimental Results and Analysis

In this section, we first introduce the experimental setting. Then, the experimental results and detailed analysis of our DCI-Net will be illustrated.

### 4.1. Experimental setting

**Datasets and evaluation metrics**. Following the setting of iPASSRNet [43], we use 60 stereo image pairs from Middlebury [36] and 800 stereo image pairs from Flickr1024 [39] as the training set. For testing set, we pick 112 stereo image pairs from Flickr1024, 20 stereo image pairs from

KITTI 2015 [31], 20 stereo image pairs from KITTI 2012 [8] and 5 stereo image pairs from Middlebury. Note that we clearly follow [51] to synthesize the low-light stereo images. For evaluations, we employ two widely used image quality assessment metrics: PSNR and SSIM. The greater PSNR and SSIM, the better quality of the enhanced results. As for the structure of our DCI-Net, we set $N_1 = 4$, $N_2 = 4$, $N_3 = 2$, $N_4 = 2$, and $N_5 = 4$.

**Implementation details**. All experiments are conducted by using Pytorch with RTX 2080Ti GPUs. We employ Adam optimizer with batch size of 16. All the training images are randomly cropped into $128 \times 128$ pixels. The initial learning rate is set to 0.0002 and reduced by half per 500 epochs. DCI-Net is trained for 2000 epochs totally. Some related methods are compared, including NAFSSR

[3], NAFSSR-F [3], iPASSRNet [43], SNR [48], DVENet [13], ZeroDCE [10] and ZeroDCE++ [20]. For iPASSR-Net, NAFSSR and NAFSSR-F, we set the upscale factor to 1. Note that we use frequency-domain reconstruction loss to retrain NAFSSR, which is denoted as NAFSSR-F.

## 4.2. Quantitative results

We evaluate the performance of our DCI-Net on four datasets, i.e., Flickr1024, KITTI 2012, KITTI 2015 and Middlebury datasets. Table 1 shows the numerical results of both single view and a pair of stereo images. We see that our DCI-Net obtains better performance on all evaluated datasets than other compared methods. Overall, the two zero-shot methods ZeroDCE and ZeroDCE++ obtain the worst results, since no supervised data drives the training process. For supervised single low-light image enhancement method, SNR is inferior to the stereo image restoration methods, because single LLIE methods do not consider the cross-view cues. To be specific, greater SSIM values suggests that our DCI-Net can better recover the structural information, compared with other competitors. The best PSNR demonstrates that our DCI-Net is capable of restoring the details to reconstruct the normal-light stereo images.

## 4.3. Visual analysis results

For better comparison, we visualize the enhanced images and corresponding error maps of each method on the Flickr1024 dataset in Fig.5. Note that the process of computing the error maps can be referred to [57]. We see that our DCI-Net is capable of recovering more consistent color and accurate illumination than other competitors. Based on the error maps, DCI-Net is the lightest one, indicating that our model can recover more details. Fig.6 shows the visual comparison on KITTI 2012 and KITTI 2015 datasets. According to the displayed PSNR and SSIM, the proposed DCI-Net obtains maximum values, which suggests that the enhanced images of our model are of higher quality. Similar results can also be found from the error maps. The errors produced by our method are obviously smaller than others. The visual results on Middlebury dataset are shown in Fig. 7. We can find that our DCI-Net can better reconstruct the normal-light image with less information loss.

## 4.4. Ablation study

We perform ablation studies to demonstrate the effectiveness of the designed modules, losses and kernel sizes in our DCI-Net. The experiments are performed on Flickr1024 dataset. The numerical results are shown in Table 2.

**Effect of DIM**. To show the effectiveness of DIM, we design three models as shown in Table 2. Specifically, W/o CVI and W/o CSI denote removing CVI and CSI from our DCI-Net respectively. W/o DIM denotes removing DIM from DCI-Net, which is also regarded as removing CVI

Table 2. Ablation study on the effects of losses, kernel size, DIM and SIMB over the Flickr2014 dataset. **Bold** denotes the best.

| Model | Left | | (Left + Right) / 2 | |
|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM |
| W/o DIM | 24.67 | 0.809 | 24.63 | 0.811 |
| W/o CVI | 25.26 | 0.828 | 25.23 | 0.829 |
| W/o CSI | 25.07 | 0.824 | 25.08 | 0.826 |
| W/o LRDC | 22.07 | 0.795 | 22.13 | 0.796 |
| W/o ECIM | 25.01 | 0.817 | 25.06 | 0.818 |
| W/ $\mathcal{L}_1$ | 24.17 | 0.788 | 24.08 | 0.789 |
| W/ $\mathcal{L}_2$ | 21.56 | 0.655 | 21.54 | 0.655 |
| W/ $\mathcal{L}_{ssim}$ | 24.01 | 0.819 | 23.93 | 0.821 |
| KS of $3 \times 3$ | 24.91 | 0.814 | 24.85 | 0.814 |
| KS of $5 \times 5$ | 25.07 | 0.823 | 25.01 | 0.823 |
| DCI-Net | **25.47** | **0.832** | **25.36** | **0.832** |

and CSI simultaneously. From Table 2, we see that there is significant performance reduction when we remove any of them. Because the proposed DIM has the ability to complete the cross-scale cross-view information interaction.

**Effect of SIMB**. We further verify the role of SIMB by dropping LRDC and ECIM from SIMB in our DCI-Net, which are denoted as W/o LRDC and W/o ECIM in Table 2. When LRDC is deleted, the numerical results decrease significantly. The performance degradation suggests that LRDC is an important component that can capture long-range dependency. Removing ECIM also makes worse performance since ECIM is able to mine the channel information and refine the feature representation.

**Effect of losses**. We evaluate the performance with different losses. As shown in Table 2, models W/ $\mathcal{L}_1$ loss, W/ $\mathcal{L}_2$ loss and W/ $\mathcal{L}_{ssim}$ denotes using $l_1$ loss, $l_2$ loss and $ssim$ loss to replace the frequency-domain reconstruction loss $\mathcal{L}_{fre}$ to train our model. The performance degrades for the three models, which means the frequency-domain reconstruction loss is superior to others for our model.

**Effect of kernel size**. We finally study the impact of different kernel sizes of the DW-Conv layer in SIMB on the results. We test the performance with kernel size of $3 \times 3$ and $5 \times 5$, denoted as KS $3 \times 3$ and KS $5 \times 5$. Note that the kernel size in our DCI-Net is set to $7 \times 7$. From Table 2, we see that the performance gets better with larger kernels, and our proposed DCI-Net obtains the best result.

## 5. Conclusion

We explored effective strategies to address the issues of weak cross-view information interaction and lacking of long-range dependencies in intra-view learning to deal with the spatial long-range effects for stereo image enhancement in the dark. Technically, we proposed a novel decoupled cross-scale cross-view interaction network (DCI-Net). To be specific, we present a decoupled interaction module

to improve the information flow via exploring cross-scale cross-view information interaction in a decoupling manner, namely, cross-view interaction at different scales and cross-scale interaction. In addition, we further propose a spatial-channel information mining block to capture long-range dependency and improve the feature representation. Extensive experiments show that our DCI-Net achieves state-of-the-art quantitative performance and obtains better visual results in terms of more accurate texture-recovery and less detail loss, compared to other related methods. In the future, designing more effective and more efficient cross-view interaction method is worth studying. Besides, incorporating the low-light stereo image enhancement into some high-level vision tasks is also an interesting future work.

## 6. Acknowledgment

## References

[1] Bolun Cai, Xian shun Xu, Kailing Guo, Kui Jia, B. Hu, and Dacheng Tao. A joint intrinsic-extrinsic prior model for retinex. In *ICCV*, 2017. 3

[2] Chen Chen, Qifeng Chen, Jia Xu, and Vladlen Koltun. Learning to see in the dark. In *CVPR*, 2018. 2

[3] Xiaojie Chu, Liangyu Chen, and Wenqing Yu. Nafssr: Stereo image super-resolution using nafnet. In *CVPR Workshops*, 2022. 1, 2, 3, 4, 5, 6, 7, 8

[4] Xiaohan Ding, Xiangyu Zhang, Jungong Han, and Guiguang Ding. Scaling up your kernels to 31x31: Revisiting large kernel design in cnns. In *CVPR*, 2022. 5

[5] Chao Dong, Chen Change Loy, Kaiming He, and Xiaoou Tang. Image super-resolution using deep convolutional networks. *IEEE TPAMI*, 2016. 2

[6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2020. 5

[7] Ruicheng Feng and Chongyi Li et al. Mipi 2022 challenge on under-display camera image restoration: Methods and results. In *ECCV Workshops*, 2022. 2

[8] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In *CVPR*, 2012. 7

[9] Jiaqi Gu, Hyoukjun Kwon, Dilin Wang, Wei Ye, Meng Li, Yu-Hsin Chen, Liangzhen Lai, Vikas Chandra, and David Z Pan. Multi-scale high-resolution vision transformer for semantic segmentation. In *CVPR*, 2022. 4

[10] Chunle Guo, Chongyi Li, Jichang Guo, Chen Change Loy, Junhui Hou, Sam Kwong, and Runmin Cong. Zero-reference deep curve estimation for low-light image enhancement. In *CVPR*, 2020. 1, 3, 5, 7, 8

[11] Xiaojie Guo, Yu Li, and Haibin Ling. Lime: Low-light image enhancement via illumination map estimation. *IEEE TIP*, 2016. 2

[12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1

[13] Jie Huang, Xueyang Fu, Zeyu Xiao, Fengmei Zhao, and Zhiwei Xiong. Low-light stereo image enhancement. *IEEE TMM*, 2022. 1, 2, 3, 4, 5, 6, 7, 8

[14] Jie Huang, Xueyang Fu, Zeyu Xiao, Feng Zhao, and Zhiwei Xiong. Low-light stereo image enhancement. *IEEE TMM*, 2022. 2

[15] Daniel S Jeon, Seung-Hwan Baek, Inchang Choi, and Min H Kim. Enhancing the spatial resolution of stereo images using a parallax prior. In *CVPR*, 2018. 3, 4

[16] Eunah Jung, Nan Yang, and Daniel Cremers. Multi-frame gan: Image enhancement for stereo visual odometry in low light. In *CRL*, 2020. 3

[17] Edwin H Land. The retinex theory of color vision. *Scientific American*, 1977. 3

[18] Chulwoo Lee, Chulwoo Lee, and Chang-Su Kim. Contrast enhancement based on layered difference representation of 2d histograms. *IEEE TIP*, 2013. 3

[19] Chongyi Li, Chunle Guo, Ling-Hao Han, Jun Jiang, Ming-Ming Cheng, Jinwei Gu, and Chen Change Loy. Low-light image and video enhancement using deep learning: a survey. *IEEE TPAMI*, 2021. 2

[20] Chongyi Li, Chunle Guo Guo, and Chen Change Loy. Learning to enhance low-light image via zero-reference deep curve estimation. In *IEEE TPAMI*, 2021. 1, 3, 5, 6, 7, 8

[21] Chongyi Li, Jichang Guo, Fatih Porikli, and Yanwei Pang. Lightennet: A convolutional neural network for weakly illuminated image enhancement. *PRL*, 2018. 2

[22] Mading Li, Jiaying Liu, Wenhan Yang, Xiaoyan Sun, and Zongming Guo. Structure-revealing low-light image enhancement via robust retinex model. *IEEE TIP*, 2018. 3

[23] Peiliang Li, Xiaozhi Chen, and Shaojie Shen. Stereo r-cnn based 3d object detection for autonomous driving. In *CVPR*, 2019. 1

[24] Yu Li, Yaling Yi, Dongwei Ren, Qince Li, and Wangmeng Zuo. Learning dual-pixel alignment for defocus deblurring. *ArXiv*, 2022. 3

[25] Jiawen Liao, Yanwei Pang, Jing Nie, Hanqing Sun, and Jiale Cao. No-reference enhancement of low-light images by stereo vision. *SSRN 4240548*, 2022. 3

[26] Jiaying Liu, Dejia Xu, Wenhan Yang, Minhao Fan, and Haofeng Huang. Benchmarking low-light image enhancement and beyond. *IJCV*, 2021. 2

[27] Risheng Liu, Long Ma, Jiaao Zhang, Xin Fan, and Zhongxuan Luo. Retinex-inspired unrolling with cooperative prior architecture search for low-light image enhancement. In *CVPR*, 2021. 2, 3

[28] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *CVPR*, 2021. 5

[29] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015. 1

[30] Long Ma, Tengyu Ma, Risheng Liu, Xin Fan, and Zhongxuan Luo. Toward fast, flexible, and robust low-light image enhancement. In *CVPR*, 2022. 3

[31] Moritz Menze and Andreas Geiger. Object scene flow for autonomous vehicles. In *CVPR*, 2015. 7

[32] Jing Nie, Yanwei Pang, J. Xie, Jing Pan, and Jungong Han. Stereo refinement dehazing network. *IEEE TCSVT*, 2022. 3

[33] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 1

[34] Wenqi Ren, Sifei Liu, Lin Ma, Qianqian Xu, Xiangyu Xu, Xiaochun Cao, Junping Du, and Ming-Hsuan Yang. Low-light image enhancement via a deep hybrid network. *IEEE TIP*, 2019. 3

[35] Xutong Ren, Wenhan Yang, Wen-Huang Cheng, and Jiaying Liu. Lr3m: Robust low-light enhancement via low-rank regularized retinex model. *IEEE TIP*, 2020. 3

[36] Daniel Scharstein, Heiko Hirschmüller, York Kitajima, Greg Krathwohl, Nera Nešić, Xi Wang, and Porter Westling. High-resolution stereo datasets with subpixel-accurate ground truth. In *GCPR*, 2014. 7

[37] Jingdong Wang, Ke Sun, Tianheng Cheng, Borui Jiang, Chaorui Deng, Yang Zhao, Dong Liu, Yadong Mu, Mingkui Tan, Xinggang Wang, et al. Deep high-resolution representation learning for visual recognition. *IEEE TPAMI*, 2020. 4

[38] Longguang Wang, Yulan Guo, Yingqian Wang, Juncheng Li, Shuhang Gu, Radu Timofte, Liangyu Chen, Xiaojie Chu, Wenqing Yu, Kai Jin, et al. Ntire 2022 challenge on stereo image super-resolution: Methods and results. In *CVPR Workshops*, 2022. 3

[39] Longguang Wang, Yingqian Wang, Zhengfa Liang, Zaiping Lin, Jungang Yang, Wei An, and Yulan Guo. Learning parallax attention for stereo image super-resolution. In *CVPR*, 2019. 2, 7

[40] Longguang Wang, Yingqian Wang, Zhengfa Liang, Zaiping Lin, Jungang Yang, Wei An, and Yulan Guo. Learning parallax attention for stereo image super-resolution. In *CVPR*, 2019. 4

[41] Ruixing Wang, Qing Zhang, Chi-Wing Fu, Xiaoyong Shen, Wei-Shi Zheng, and Jiaya Jia. Underexposed photo enhancement using deep illumination estimation. In *CVPR*, 2019. 2

[42] Yufei Wang, Renjie Wan, Wenhan Yang, Haoliang Li, Lap-Pui Chau, and Alex C Kot. Low-light image enhancement with normalizing flow. In *AAAI*, 2022. 2

[43] Yingqian Wang, Xinyi Ying, Longguang Wang, Jungang Yang, Wei An, and Yulan Guo. Symmetric parallax attention for stereo image super-resolution. In *CVPR Workshops*, 2021. 1, 2, 3, 4, 5, 6, 7, 8

[44] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. In *BMVC*, 2018. 2, 3

[45] Yanyan Wei, Zhao Zhang, Huan Zheng, Richang Hong, Yi Yang, and Meng Wang. Sginet: Toward sufficient interaction between single image deraining and semantic segmentation. In *MM*, 2022. 2

[46] Wen-Bin Wu, Jian Weng, Pingping Zhang, Xu Wang, Wenhan Yang, and Jianmin Jiang. Uretinex-net: Retinex-based deep unfolding network for low-light image enhancement. In *CVPR*, 2022. 3

[47] Ke Xu, Xin Yang, Baocai Yin, and Rynson WH Lau. Learning to restore low-light images via decomposition-and-enhancement. In *CVPR*, 2020. 2

[48] Xiaogang Xu, Ruixing Wang, Chiao Fu, and Jiaya Jia. Snr-aware low-light image enhancement. In *CVPR*, 2022. 1, 3, 5, 6, 7, 8

[49] Wenhan Yang, Shiqi Wang, Yuming Fang, Yue Wang, and Jiaying Liu. Band representation-based semi-supervised low-light image enhancement: Bridging the gap between signal fidelity and perceptual quality. *IEEE TIP*, 2021. 3

[50] Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *CVPR*, 2022. 5

[51] Fan Zhang, Yu Li, Shaodi You, and Ying Fu. Learning temporal consistency for low light video enhancement from single images. In *CVPR*, 2021. 7

[52] Kaihao Zhang, Wenhan Luo, Yanjiang Yu, Wenqi Ren, Fang Zhao, Changsheng Li, Lin Ma, Wei Liu, and Hongdong Li. Beyond monocular deraining: Parallel stereo deraining network via semantic prior. *IJCV*, 2022. 3

[53] Yonghua Zhang, Jiawan Zhang, and Xiaojie Guo. Kindling the darkness: A practical low-light image enhancer. In *MM*, 2019. 3

[54] Zhao Zhang, Huan Zheng, Richang Hong, Jicong Fan, Yi Yang, and Shuicheng Yan. Frc-net: A simple yet effective architecture for low-light image enhancement. *TechRxiv*, 2022. 2

[55] Zhao Zhang, Huan Zheng, Richang Hong, Mingliang Xu, Shuicheng Yan, and Meng Wang. Deep color consistent network for low-light image enhancement. In *CVPR*, 2022. 2

[56] Suiyi Zhao, Zhao Zhang, Richang Hong, Mingliang Xu, Yi Yang, and Meng Wang. Fcl-gan: A lightweight and real-time baseline for unsupervised blind image deblurring. In *MM*, 2022. 2

[57] Chuanjun Zheng, Daming Shi, and Yukun Liu. Windowing decomposition convolutional neural network for image enhancement. In *MM*, 2021. 8

[58] Chuanjun Zheng, Daming Shi, and Wentian Shi. Adaptive unfolding total variation network for low-light image enhancement. In *ICCV*, 2021. 3

[59] Huan Zheng, Zhao Zhang, Yang Wang, Zheng Zhang, Mingliang Xu, Yi Yang, and Meng Wang. Gcm-net: Towards effective global context modeling for image inpainting. In *MM*, 2021. 2

[60] Huan Zheng, Zhao Zhang, Haijun Zhang, Yi Yang, Shuicheng Yan, and Meng Wang. Deep multi-resolution mutual learning for image inpainting. In *MM*, 2022. 2

[61] Shangchen Zhou, Jiawei Zhang, Wangmeng Zuo, Haozhe Xie, Jinshan Pan, and Jimmy S Ren. Davanet: Stereo deblurring with view aggregation. In *CVPR*, 2019. 3