

# Reti-Diff: Illumination Degradation Image Restoration with Retinex-based Latent Diffusion Model

Chunming He<sup>1,\*</sup>, Chengyu Fang<sup>1,\*</sup>, Yulun Zhang<sup>2,†</sup>, Tian Ye<sup>3</sup>, Kai Li<sup>4</sup>, Longxiang Tang<sup>1</sup>, Zhenhua Guo<sup>5</sup>, Xiu Li<sup>1,†</sup>, and Sina Farsiu<sup>6</sup>

<sup>1</sup> Shenzhen International Graduate School, Tsinghua University

<sup>2</sup> Shanghai Jiao Tong University

<sup>3</sup> The Hong Kong University of Science and Technology (Guangzhou)

<sup>4</sup> NEC Laboratories America

<sup>5</sup> Tianyi Traffic Technology

<sup>6</sup> Duke University

**Abstract.** Illumination degradation image restoration (IDIR) techniques aim to improve the visibility of degraded images and mitigate the adverse effects of deteriorated illumination. Among these algorithms, diffusion model (DM)-based methods have shown promising performance but are often burdened by heavy computational demands and pixel misalignment issues when predicting the image-level distribution. To tackle these problems, we propose to leverage DM within a compact latent space to generate concise guidance priors and introduce a novel solution called Reti-Diff for the IDIR task. Reti-Diff comprises two key components: the Retinex-based latent DM (RLDM) and the Retinex-guided transformer (RGformer). To ensure detailed reconstruction and illumination correction, RLDM is empowered to acquire Retinex knowledge and extract reflectance and illumination priors. These priors are subsequently utilized by RGformer to guide the decomposition of image features into their respective reflectance and illumination components. Following this, RGformer further enhances and consolidates the decomposed features, resulting in the production of refined images with consistent content and robustness to handle complex degradation scenarios. Extensive experiments show that Reti-Diff outperforms existing methods on three IDIR tasks, as well as downstream applications. The code will be released.

**Keywords:** Illumination degradation image restoration · Latent diffusion model · Retinex theory

## 1 Introduction

Illumination degradation image restoration (IDIR) seeks to enhance the visibility and contrast of degraded images while mitigating the adverse effects of deteriorated illumination, *e.g.*, indefinite noise and variable color deviation. IDIR has

---

\* Equal Contribution, † Corresponding Author

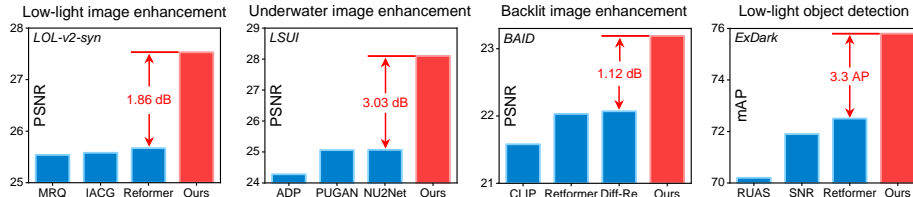
been investigated in various domains, including low-light image enhancement [3], underwater image enhancement [18], and backlit image enhancement [41]. By addressing illumination degradation, the enhanced images are expected to exhibit improved visual quality, making them more suitable for decision-making or subsequent tasks like nighttime object detection and segmentation.

Traditional IDIR approaches [14, 59] primarily rely on manually crafted enhancement techniques with limited generalization capabilities. Leveraging the robust feature extraction capabilities of convolutional neural networks and transformers, a series of deep learning-based methods [3, 29] have been proposed and have achieved remarkable success in the IDIR domain. However, as depicted in Fig. 1, they still face challenges in complex illumination degradation scenarios due to their constrained generative capacity.

To overcome these challenges, deep generative models, like generative adversarial networks [20] and variational autoencoder [24], have gained popularity in the IDIR task for their generative abilities. Recently, the diffusion model (DM) [73] has been introduced to the IDIR field for high-quality image restoration. However, existing DM-based methods, *e.g.*, Diff-Retinex [73] and GSAD [30], apply DM directly to image-level generation, leading to two main challenges: **(1)** These methods incur high computational costs, as predicting the image-level distribution requires a large number of inference steps. **(2)** The enhanced results may exhibit pixel misalignment with the original clean image in terms of restored details and local consistency.

To tackle the above problems, we propose introducing the latent diffusion model (LDM) to solve the IDIR problem. By applying DM in the low-dimensional compact latent space, we effectively alleviate the computational burden. Additionally, we incorporate LDM into transformers to prevent pixel misalignment in the generated image, which is often observed in existing deep generative models. Unlike existing LDM-based methods that solely use the priors extracted from the RGB domain, our method, tailored to the specific characteristics of IDIR tasks, empowers LDMs to extract Retinex information from both the reflectance and illumination domains. This adaptation allows our method to generate high-fidelity Retinex priors directly from low-quality input images. By doing so, this approach enables us to simultaneously enhance image details using the reflectance prior and correct color distortions with the illumination prior, resulting in visually appealing results with favorable downstream tasks.

With this inspiration, we present Reti-Diff, the first LDM-based solution to tackle the IDIR problem. Reti-Diff, depicted in Fig. 2, consists of two primary components: the Retinex-based LDM (RLDM) and the Retinex-guided transformer (RGformer). Initially, RLDM is employed to generate Retinex priors, which are then integrated into RGformer to produce visually appealing results. To ensure the generation of high-quality priors, we propose a two-phase training approach, wherein Reti-Diff undergoes initial pretraining followed by subsequent RLDM optimization. **In phase I**, we introduce a Retinex prior extraction (RPE) module to compress the ground-truth image into the highly compact Retinex priors, namely the reflectance prior and the illumination prior. These



**Fig. 1:** Our Reti-Diff outperforms cutting-edge techniques on three IDIR tasks and the low-light object detection task, where CLIP, Diff-Re, and SNR are short for CLIP-LIT [41], Diff-Retinex [73], and SNR-Net [69].

priors are then sent to RGformer to guide feature decomposition and the generation of reflectance and illumination features. Afterward, RGformer employs the Retinex-guided multi-head cross attention (RG-MCA) and dynamic feature aggregation (DFA) module to refine and aggregate the decomposed features, ultimately producing enhanced images with coherent content and ensuring robustness and generalization in extreme degradation scenarios. **In phase II**, we train RLDM in reflectance and illumination domains to estimate Retinex priors from the low-quality image, with the constraint of consistency with those extracted by RPE from the ground-truth image. Therefore, the extracted Retinex priors can guide the RGformer in detail enhancement and illumination correction, resulting in visually appealing results with favorable downstream performance.

Our contributions are summarized as follows:

- We propose a novel DM-based framework, Reti-Diff, for the IDIR task. To the best of our knowledge, this is the first application of the latent diffusion model to tackle the IDIR problem.
- We propose to let RLDM learn Retinex knowledge and generate high-quality reflectance and illumination priors from the low-quality input, which serve as critical guidance in detail enhancement and illumination correction.
- We propose RGformer to integrate extracted Retinex priors to decompose features into reflectance and illumination components and then utilize RG-MCA and DFA to refine and aggregate the decomposed features, ensuring robustness and generalization in complex illumination degradation scenarios.
- Extensive experiments on three IDIR tasks verify our superiority to existing methods in terms of image quality and favorability in downstream applications, including low-light object detection and segmentation.

## 2 Related Works

**Illumination Degradation Image Restoration.** Early IDIR methods mainly include three approaches: histogram equalization (HE) [1], gamma correction (GC) [27], and Retinex theory [35]. HE-based and GC-based methods focus on directly amplifying the low contrast regions but overlook illumination factors. Retinex-based variants [15, 40] propose the development of priors to constrain the solution space for reflectance and illumination maps. However, these methods still rely on hand-crafted priors, limiting their ability to generalize effectively. With the rapid development of deep learning, approaches based on CNNs and

transformers [3, 20, 29] have achieved remarkable success in IDIR. For instance, LLNet [46] proposed a sparse denoising structure to enhance illumination and suppress noise. DIE [60] integrated Retinex cues into a learning-based structure, presenting a one-stage Retinex-based solution for color correction. To enhance generative capacity, Diff-Retinex [73] and GSAD [30] introduced DM to the IDIR field by directly applying it to image-level generation. However, they entail significant computational costs and may lead to pixel misalignment with the original input, particularly concerning restored image details and local consistency.

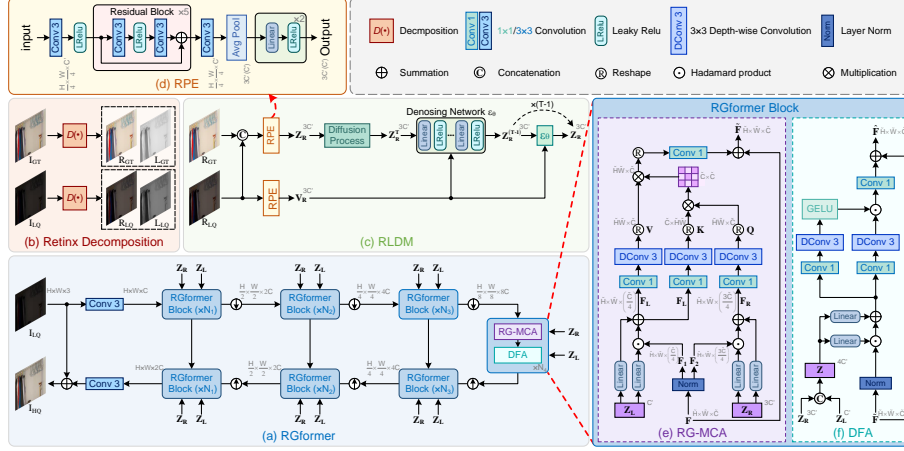
**Diffusion Models.** Diffusion models (DMs) have demonstrated considerable success in various domains, including density estimation [33] and data generation [24]. Such a probabilistic generative model adopts a parameterized Markov chain to optimize the lower variational bound on the likelihood function, enabling them to generate target distributions with greater accuracy than other generative models, *i.e.*, GAN and VAE. Recently, DMs have been introduced to solve the IDIR problem [30, 73]. However, when directly applied to image-level generation, these approaches introduce computational burdens and pixel misalignment. To overcome this, we propose employing LDM to estimate priors within a low-dimensional latent space. We then integrate these priors into the transformer-based framework, thus addressing the above problems. Besides, unlike existing LDM-based methods [6, 66] that solely rely on priors extracted from the RGB domain, our method, tailored to the specific characteristics of IDIR tasks, empowers LDMs to extract Retinex information from both the reflectance and illumination domains. This adaptation allows our method to generate high-fidelity Retinex priors directly from low-quality input images. By doing so, this novel approach enables us to simultaneously enhance image details using the reflectance prior and correct color distortions with the illumination prior, resulting in visually appealing results with favorable downstream tasks.

### 3 Methodology

In this paper, we propose Reti-Diff, the pioneering method based on Latent Diffusion Models (LDM) for IDIR tasks. Reti-Diff is specifically tailored to address the challenges inherent in IDIR tasks by leveraging Retinex priors extracted from both the illumination and reflectance domains to guide the restoration process. This innovative approach utilizes the extracted Retinex prior representation as dynamic modulation parameters, facilitating simultaneous enhancement of restoration details through the reflectance prior and correction of color distortion via the illumination prior. This ensures the generation of visually compelling results while positively impacting downstream tasks.

As shown in Fig. 2, our Reti-Diff comprises two parts: the Retinex-guided transformer (RGformer) and the Retinex-based latent diffusion model (RLDM). To ensure the generation of high-quality priors, Reti-Diff undergoes a two-phase training strategy, involving the initial pretraining of Reti-Diff and the subsequent optimization of RLDM. In this section, we provide an in-depth explanation of the two-phase training approach and elucidate the entire restoration process.





**Fig. 2:** Framework of Reti-Diff. We first pretrain Reti-Diff to ensure the robust learning of RLDM and then optimize RLDM to generate high-quality Retinex priors, which guide RGformer in detail enhancement and illumination correction. In (a), we omit the auxiliary decoder  $D_a(\cdot)$  for simplicity. In (c), we only give the example by using RLDM to extract the reflectance prior and the illumination prior can be extracted similarly.

### 3.1 Pretrain Reti-Diff

We first pretrain Reti-Diff to encode the clear image, termed ground truth, into compact priors with Retinex prior extraction (RPE) module and use the extracted Retinex priors to guide RGformer for restoration.

**Retinex prior extraction module.** Given the low-quality (LQ) image  $\mathbf{I}_{LQ} \in \mathbb{R}^{H \times W \times 3}$  and its corresponding ground truth  $\mathbf{I}_{GT} \in \mathbb{R}^{H \times W \times 3}$ , we initially decompose them into the reflectance image  $\mathbf{R} \in \mathbb{R}^{H \times W \times 3}$  and the illumination map  $\mathbf{L} \in \mathbb{R}^{H \times W}$  according to Retinex theory:

$$\mathbf{I}_{LQ} = \mathbf{R}_{LQ} \odot \mathbf{L}_{LQ}, \mathbf{I}_{GT} = \mathbf{R}_{GT} \odot \mathbf{L}_{GT}, \quad (1)$$

where  $\odot$  denotes the Hadamard product. Following Retformer [3], We use a pre-trained decomposing network  $D(\cdot)$  to decompose  $\mathbf{I}_{LQ}$  and  $\mathbf{I}_{GT}$ . Then we concatenate the corresponding components of ground truth and LQ image and use the RPE module  $\text{RPE}(\cdot)$  to encode them into Retinex priors  $\mathbf{Z}_R \in \mathbb{R}^{3C'}$ ,  $\mathbf{Z}_L \in \mathbb{R}^{C'}$ :

$$\begin{aligned} \mathbf{Z}_R &= \text{RPE}(\text{down}(\text{conca}(\mathbf{R}_{GT}, \mathbf{R}_{LQ}))), \\ \mathbf{Z}_L &= \text{RPE}(\text{down}(\text{conca}(\mathbf{L}_{GT}, \mathbf{L}_{LQ}))), \end{aligned} \quad (2)$$

where  $\text{conca}(\cdot)$  is concatenation and  $\text{down}(\cdot)$  is downsampling that is operated by PixelUnshuffle. We then send Retinex priors,  $\mathbf{Z}_R$  and  $\mathbf{Z}_L$ , to RGformer to serve as dynamic modulation parameters for detail restoration and color correction.

**Retinex-guided transformer.** RGformer mainly consists of two parts in each block, *i.e.*, Retinex-guided multi-head cross attention (RG-MCA) and dynamic feature aggregation (DFA) module. In RG-MCA, we first split the input feature  $\mathbf{F} \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times \tilde{C}}$  into two parts  $\mathbf{F}_1 \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times (3\tilde{C}/4)}$  and  $\mathbf{F}_2 \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times (\tilde{C}/4)}$  along the channel dimension. Afterwards, we integrated  $\mathbf{Z}_R$  and  $\mathbf{Z}_L$  as the corresponding dynamic modulation parameters to generate reflectance-guided feature

$\mathbf{F}_R \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times (3\tilde{C}/4)}$  and illumination-guided feature  $\mathbf{F}_L \in \mathbb{R}^{\tilde{H} \times \tilde{W} \times (\tilde{C}/4)}$ :

$$\begin{aligned}\mathbf{F}_R &= \text{Li}_1(\mathbf{Z}_R) \odot \text{Norm}(\mathbf{F}_1) + \text{Li}_2(\mathbf{Z}_R), \\ \mathbf{F}_L &= \text{Li}_1(\mathbf{Z}_L) \odot \text{Norm}(\mathbf{F}_2) + \text{Li}_2(\mathbf{Z}_L),\end{aligned}\quad (3)$$

where  $\text{Norm}(\cdot)$  is layer normalization.  $\text{Li}(\cdot)$  means linear layer. Afterward, we aggregate global spatial information by projecting  $\mathbf{F}_R$  into query  $\mathbf{Q} = \mathbf{W}_Q \mathbf{F}_R$  and key  $\mathbf{K} = \mathbf{W}_K \mathbf{F}_L$  and transforming  $\mathbf{F}_L$  into value  $\mathbf{V} = \mathbf{W}_V \mathbf{F}_L$ , where  $\mathbf{W}$  is the combination of a  $1 \times 1$  point-wise convolution and a  $3 \times 3$  depth-wise convolution. We then perform cross-attention and get the output feature  $\tilde{\mathbf{F}}$ :

$$\tilde{\mathbf{F}} = \mathbf{F} + \text{SoftMax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{\tilde{C}}}\right) \cdot \mathbf{V}. \quad (4)$$

By doing so, RG-MCA introduces explicit guidance to fully exploit Retinex knowledge at the feature level and use cross attention mechanism to implicitly model the Retinex theory and refine the decomposed features, which helps to restore missing details and correct color distortion.

Then we employ DFA for local feature aggregation. Apart from the  $1 \times 1$  convolution and  $3 \times 3$  depth-wise convolution used for information fusion, DFA also adopts GELU, termed  $\text{GELU}(\cdot)$ , to ensure the flexibility of aggregation [22]. Thus, given  $\tilde{\mathbf{F}}$  and  $\mathbf{Z}$ , where  $\mathbf{Z} = \text{conca}(\mathbf{Z}_R, \mathbf{Z}_L)$ , the output feature  $\hat{\mathbf{F}}$  is

$$\begin{aligned}\hat{\mathbf{F}} &= \tilde{\mathbf{F}} + \text{GELU}(\mathbf{W}_1 \mathbf{F}') \odot \mathbf{W}_2 \mathbf{F}', \\ \mathbf{F}' &= \text{Li}_1(\mathbf{Z}) \odot \text{Norm}(\tilde{\mathbf{F}}) + \text{Li}_2(\mathbf{Z}),\end{aligned}\quad (5)$$

**Optimization.** To facilitate the extraction of Retinex priors, the RPE module and RGformer are jointly trained by a reconstruction loss with  $L_1$  norm  $\|\cdot\|_1$ :

$$L_{Rec} = \|\mathbf{I}_{GT} - \mathbf{I}_{HQ}\|_1, \quad (6)$$

where  $\mathbf{I}_{HQ}$  is the enhanced result. In addition, to ensure that the separated features within RG-MCA effectively capture reflectance and illumination knowledge, we provide an auxiliary decoder  $D_a(\cdot)$  with the same structure as that in [44].  $D_a(\cdot)$  takes  $\tilde{\mathbf{F}}$  as input and outputs the reconstructed reflectance image  $\mathbf{R}_{Re}$  and illumination map  $\mathbf{L}_{Re}$ . For efficiency, we only apply  $D_a(\cdot)$  for the first transformer block in encoder to get  $\mathbf{R}_{Re}^I$  and  $\mathbf{L}_{Re}^I$  and for the last transformer block in decoder to get  $\mathbf{R}_{Re}^L$  and  $\mathbf{L}_{Re}^L$ .  $D_a(\cdot)$  is supervised by a Retinex loss  $L_R$ :

$$L_R = \|\mathbf{R}_{LQ} - \mathbf{R}_{Re}^I\|_1 + \|\mathbf{L}_{LQ} - \mathbf{L}_{Re}^I\|_1 + \|\mathbf{R}_{GT} - \mathbf{R}_{Re}^L\|_1 + \|\mathbf{L}_{GT} - \mathbf{L}_{Re}^L\|_1, \quad (7)$$

Eq. (7) serves to maintain crucial Retinex information throughout the network. Hence, the integration of Eq. (7) not only promotes the assimilation of Retinex theory by the split features but also amplifies the overall restoration capacity.

In Phase I, the final loss  $L_{P1}$  is defined as follows:

$$L_{P1} = L_{Rec} + \lambda_1 L_R, \quad (8)$$

where  $\lambda_1$  is a hyperparameter and  $\lambda_1 = 1$ .

### 3.2 Retinex-based Latent Diffusion Model

In Phase II, we train the RLDM to predict Retinex priors from the low-quality input, which are expected to be consistent with that extracted by RPE from the ground-truth image. Unlike conventional LDMs trained on the RGB domain, we introduce two RLDMs with a Siamese structure and train them on distinct

domains: the reflectance domain and the illumination domain. This approach, grounded in Retinex theory, equips our RLDM to generate a more generative reflectance prior  $\hat{\mathbf{Z}}_{\mathbf{R}}$  to enhance image details, and a more harmonized illumination prior  $\hat{\mathbf{Z}}_{\mathbf{L}}$  for color correction. Note that RLDM is constructed upon the conditional denoising diffusion probabilistic models, with both a forward diffusion process and a reverse denoising process. To simplify, we provide a detailed derivation for  $\hat{\mathbf{Z}}_{\mathbf{R}}$  herein, while that of  $\hat{\mathbf{Z}}_{\mathbf{L}}$  can be found in the appendix.

**Diffusion process.** In the diffusion process, we first use the pretrained RPE to extract the reflectance prior  $\mathbf{Z}_{\mathbf{R}}$ , which is treated as the starting point of the forward Markov process, *i.e.*,  $\mathbf{Z}_{\mathbf{R}} = \mathbf{Z}_{\mathbf{R}}^0$ . We then gradually add Gaussian noise to  $\mathbf{Z}_{\mathbf{R}}$  by  $T$  iterations and each iteration can be defined as:

$$q(\mathbf{Z}_{\mathbf{R}}^t | \mathbf{Z}_{\mathbf{R}}^{t-1}) = \mathcal{N}\left(\mathbf{Z}_{\mathbf{R}}^t; \sqrt{1 - \beta^t} \mathbf{Z}_{\mathbf{R}}^{t-1}, \beta^t \mathbf{I}\right), \quad (9)$$

where  $t = 1, \dots, T$ .  $\mathbf{Z}_{\mathbf{R}}^t$  denotes the noisy prior at time step  $t$ ,  $\beta^t$  is the predefined factor that controls the noise variance, and  $\mathcal{N}$  is the Gaussian distribution. Following [34], Eq. (9) can be simplified as follows:

$$q(\mathbf{Z}_{\mathbf{R}}^t | \mathbf{Z}_{\mathbf{R}}^0) = \mathcal{N}\left(\mathbf{Z}_{\mathbf{R}}^t; \sqrt{\bar{\alpha}^t} \mathbf{Z}_{\mathbf{R}}^0, (1 - \bar{\alpha}^t) \mathbf{I}\right), \quad (10)$$

where  $\alpha^t = 1 - \beta^t$  and  $\bar{\alpha}^t = \prod_{i=1}^t \alpha^i$ .

**Reverse process.** In the reverse process, RLDM aims to extract the reflectance prior from pure Gaussian noise. Thus, RLDM samples a Gaussian random noise map  $\mathbf{Z}_{\mathbf{R}}^T$  and then gradually denoise it to run backward from  $\mathbf{Z}_{\mathbf{R}}^T$  to  $\mathbf{Z}_{\mathbf{R}}^0$ :

$$p(\mathbf{Z}_{\mathbf{R}}^{t-1} | \mathbf{Z}_{\mathbf{R}}^t, \mathbf{Z}_{\mathbf{R}}^0) = \mathcal{N}\left(\mathbf{Z}_{\mathbf{R}}^{t-1}; \boldsymbol{\mu}^t(\mathbf{Z}_{\mathbf{R}}^t, \mathbf{Z}_{\mathbf{R}}^0), (\boldsymbol{\sigma}^t)^2 \mathbf{I}\right), \quad (11)$$

where mean  $\boldsymbol{\mu}^t(\mathbf{Z}_{\mathbf{R}}^t, \mathbf{Z}_{\mathbf{R}}^0) = \frac{1}{\sqrt{\alpha^t}} \left(\mathbf{Z}_{\mathbf{R}}^t - \frac{1 - \alpha^t}{\sqrt{1 - \bar{\alpha}^t}} \boldsymbol{\epsilon}\right)$  and variance  $(\boldsymbol{\sigma}^t)^2 = \frac{1 - \bar{\alpha}^{t-1}}{1 - \bar{\alpha}^t} \beta^t$ .  $\boldsymbol{\epsilon}$  denotes the noise in  $\mathbf{Z}_{\mathbf{R}}^t$  and is the only uncertain variable. Following previous practice [66], we employ a denoising network  $\epsilon_{\theta}(\cdot)$  to estimate  $\boldsymbol{\theta}$ . To operate in the latent space, we further introduce another RPE module  $\widetilde{\text{RPE}}(\cdot)$  to extract the conditional reflectance vector  $\mathbf{V}_{\mathbf{R}} \in \mathbb{R}^{3C'}$  from the reflectance image  $\mathbf{R}_{LQ}$  of the LQ image, *i.e.*,  $\mathbf{V}_{\mathbf{R}} = \widetilde{\text{RPE}}(\text{down}(\mathbf{R}_{LQ}))$ . Therefore, the denoising network can be represented by  $\epsilon_{\theta}(\mathbf{Z}_{\mathbf{R}}^t, \mathbf{V}_{\mathbf{R}}, t)$ . By setting the variance to  $1 - \alpha^t$ , we get

$$\mathbf{Z}_{\mathbf{R}}^{t-1} = \frac{1}{\sqrt{\alpha^t}} \left(\mathbf{Z}_{\mathbf{R}}^t - \frac{1 - \alpha^t}{\sqrt{1 - \bar{\alpha}^t}} \epsilon_{\theta}(\mathbf{Z}_{\mathbf{R}}^t, \mathbf{V}_{\mathbf{R}}, t)\right) + \sqrt{1 - \alpha^t} \boldsymbol{\epsilon}^t, \quad (12)$$

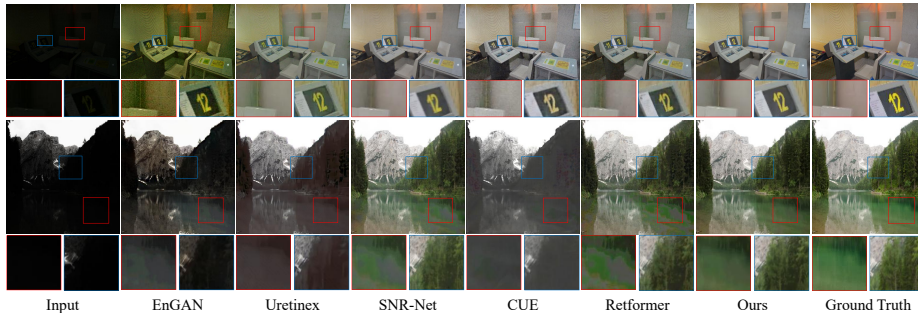
where  $\boldsymbol{\epsilon}^t \sim \mathcal{N}(0, \mathbf{I})$ . By using Eq. (12) for  $T$  iterations, we can get the predicted prior  $\hat{\mathbf{Z}}_{\mathbf{R}}$  and use it to guide RGformer for image restoration. Because the size of the predicted prior  $\hat{\mathbf{Z}}_{\mathbf{R}} \in \mathbb{R}^{3C'}$  is much smaller than the original reflectance image  $\mathbf{R}_{LQ} \in \mathbb{R}^{H \times W \times C}$ , RLDM needs much less iterations than those image-level diffusion models [73]. Thus, we run the complete  $T$  iterations for the prior generation rather than randomly selecting one time step.

**Optimization.** Given the predicted priors  $\hat{\mathbf{Z}}_{\mathbf{R}}$  and  $\hat{\mathbf{Z}}_{\mathbf{L}}$ , generated by two Siamese RLDMs with specific weights, we propose the diffusion loss to supervise them:

$$L_{Dif} = \|\mathbf{Z}_{\mathbf{R}} - \hat{\mathbf{Z}}_{\mathbf{R}}\|_1 + \|\mathbf{Z}_{\mathbf{L}} - \hat{\mathbf{Z}}_{\mathbf{L}}\|_1. \quad (13)$$

Methods	Sources	LOL-v1				LOL-v2-real				LOL-v2-synthetic				SID			
		PSNR ↑	SSIM ↑	FID ↓	BIQE ↓	PSNR ↑	SSIM ↑	FID ↓	BIQE ↓	PSNR ↑	SSIM ↑	FID ↓	BIQE ↓	PSNR ↑	SSIM ↑	FID ↓	BIQE ↓
MIRNet [75]	ECCV20	24.14	0.835	71.16	47.75	20.02	0.820	82.25	41.18	21.94	0.876	40.18	36.29	20.84	0.605	81.37	40.63
EnGAN [29]	TIP21	17.48	0.656	153.98	35.82	18.23	0.617	173.28	51.06	16.57	0.734	93.66	45.59	17.23	0.543	77.52	33.47
RUAS [42]	CVPR21	18.23	0.723	127.60	45.17	18.27	0.723	151.62	34.73	16.55	0.652	91.60	46.38	18.44	0.581	72.18	45.02
IPT [5]	CVPR21	16.27	0.504	158.83	29.35	19.80	0.813	97.24	31.17	18.30	0.811	76.79	42.15	20.53	0.618	70.58	36.71
URetinex [65]	CVPR22	21.33	0.835	85.59	30.37	20.44	0.806	76.74	28.85	24.73	0.897	33.25	33.46	22.09	0.633	71.58	38.44
UFormer [63]	CVPR22	16.36	0.771	166.69	41.06	18.82	0.771	164.41	40.36	19.66	0.871	58.69	39.75	18.54	0.577	100.14	42.13
Restormer [74]	CVPR22	22.43	0.823	78.75	33.18	19.94	0.827	114.35	37.27	21.41	0.830	46.89	35.06	22.27	0.649	75.47	32.49
SNR-Net [69]	CVPR22	24.61	0.842	66.47	28.73	21.48	0.849	68.56	28.83	24.14	0.928	30.52	33.47	22.87	0.625	74.78	30.08
SMG [70]	CVPR23	24.82	0.838	69.47	30.15	22.62	0.857	71.76	30.32	25.62	0.905	23.36	29.35	23.18	0.644	77.58	31.50
PyDiff [82]	IJCAI23	21.15	<b>0.857</b>	<b>49.47</b>	21.13	—	—	—	—	—	—	—	—	—	—	—	—
Retormer [3]	ICCV23	25.16	0.845	72.38	26.68	<b>22.80</b>	0.840	79.58	34.39	<b>25.67</b>	0.930	22.78	30.26	24.44	0.680	82.64	35.04
Diff-Retinex [73]	ICCV23	21.98	0.852	51.33	<b>19.62</b>	20.17	0.826	<b>46.67</b>	<b>24.18</b>	24.30	0.921	28.74	26.35	23.62	0.665	<b>58.93</b>	31.17
MRQ [43]	ICCV23	<b>25.24</b>	0.855	53.32	22.73	22.37	0.854	68.89	33.61	25.54	<b>0.940</b>	20.86	<b>25.09</b>	24.62	0.683	61.09	<b>27.81</b>
IAGC [62]	ICCV23	24.53	0.842	59.73	25.50	22.20	<b>0.863</b>	70.34	31.70	25.58	<b>0.941</b>	21.38	30.32	<b>24.80</b>	<b>0.688</b>	63.72	29.53
DIHR [66]	ICCV23	23.15	0.828	70.13	26.38	21.15	0.816	72.33	29.15	24.76	0.921	28.87	27.74	23.17	0.640	78.80	30.56
CUE [81]	ICCV23	21.86	0.841	69.83	27.15	21.19	0.829	67.05	28.83	24.41	0.917	31.33	33.83	23.25	0.652	77.38	28.85
GSAD [30]	NIPS23	23.23	0.852	51.64	19.96	20.19	0.847	46.77	28.85	24.22	0.927	<b>19.24</b>	25.76	—	—	—	—
Reti-Diff (Ours)	—	<b>25.35</b>	<b>0.866</b>	<b>49.14</b>	<b>17.75</b>	<b>22.97</b>	<b>0.858</b>	<b>43.18</b>	<b>23.66</b>	<b>27.53</b>	<b>0.951</b>	<b>13.26</b>	<b>15.77</b>	<b>25.53</b>	<b>0.692</b>	<b>51.66</b>	<b>25.58</b>

**Table 1:** Results on the low-light image enhancement task. The best two results are in red and blue fonts, respectively.



**Fig. 3:** Visual results on the low-light image enhancement task.

For restoration quality, we propose joint training RPE, RGformer, and RLDM. Thus, the loss in Phase II is formulated as follows:

$$L_{P2} = L_{Dif} + \lambda_2 L_{Rec} + \lambda_3 L_R, \quad (14)$$

where  $\lambda_2$  and  $\lambda_3$  are two hyper-parameters and are set as 1 in this paper.

### 3.3 Inference

In the inference phase, given the LQ input  $\mathbf{I}_{LQ}$ , Reti-Diff first uses  $\widehat{\text{RPE}}$  to extract the conditional vectors  $\mathbf{V}_R$  and  $\mathbf{V}_L$ , and then generates predicted Retinex priors  $\hat{\mathbf{Z}}_R$  and  $\hat{\mathbf{Z}}_L$  with two RLDMs. Under the guidance of the Retinex priors, RGformer generates the restored HQ image  $\mathbf{I}_{HQ}$ . Benefiting from our Retinex-based diffusion framework,  $\mathbf{I}_{HQ}$  enjoys richer texture details and more harmonized illumination, thereby further facilitating downstream tasks.

## 4 Experiment

### 4.1 Experimental Setup

Our Reti-Diff is implemented in PyTorch on four RTX3090TI GPUs and is optimized by Adam with momentum terms (0.9, 0.999). In phases I and II, we train the network for 300K iterations and the learning rate is initially set as  $2 \times 10^{-4}$  and gradually reduced to  $1 \times 10^{-6}$  with the cosine annealing [47]. Following [73], random rotation and flips are used for augmentation. Reti-Diff

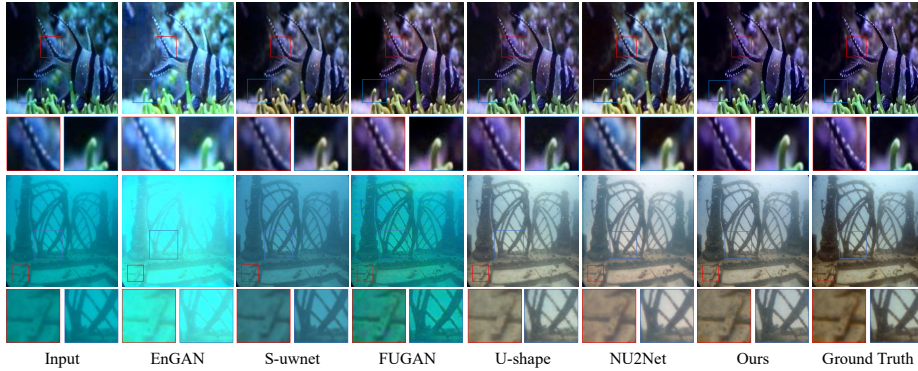


Fig. 4: Visual results on the underwater image enhancement task.

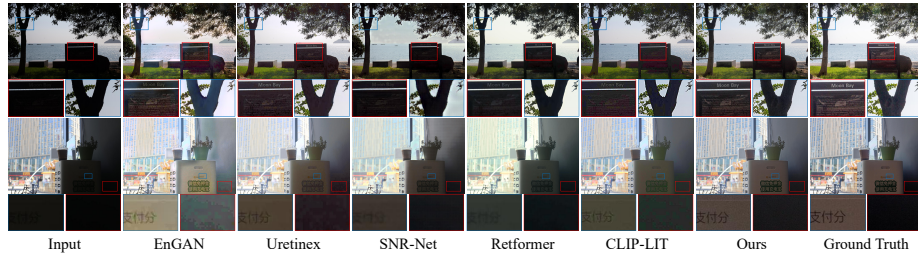


Fig. 5: Visual results on the backlit image enhancement task.

mainly comprises RLDM and RGformer. For RLDM, the channel number  $C'$  is set as 64. The total time step  $T$  is set to 4 and the hyperparameters  $\beta^{1:T}$  linearly increase from  $\beta^1 = 0.1$  to  $\beta^T = 0.99$ . RGformer adopts a 4-level cascade encoder-decoder structure. We set the number of transformer blocks, the attention heads, the channel number as  $[3, 3, 3, 3]$ ,  $[1, 2, 4, 8]$ ,  $[64, 128, 256, 512]$  from level 1 to 4.

## 4.2 Comparative Evaluation

**Low-light Image Enhancement.** We conduct a comprehensive evaluation on four datasets: *LOL-v1* [64], *LOL-v2-real* [72], *LOL-v2-syn* [72], and *SID* [4]. We adhere to the training manner outlined in [3]. Our assessment involves four metrics: PSNR, SSIM, FID [26], and BIQE [53]. Note that larger PSNR and SSIM scores, as well as smaller FID and BIQE scores, denote superior performance. We compare our approach against 17 cutting-edge enhancement techniques and report the results in Tab. 1. As depicted in Tab. 1, our method emerges as the top performer across all datasets and significantly outperforms the second-best method (Diff-Retinex) by 13.2%. These results underscore the superiority of our Reti-Diff. Fig. 3 presents qualitative results, showcasing our capacity to generate enhanced images with corrected illumination and enhanced texture, even in extremely challenging conditions. In contrast, existing methods struggle to achieve the same level of performance such as the boundaries of power lines, color harmonization of lakes, and detailed

	Diff-Retinex [73]	PyDiff [82]	GSAD [30]	Ours
Parameter (M)	56.88	97.89	17.17	26.11
MACs (G)	396.32	459.69	1340.63	156.55
FPS	4.25	3.63	2.33	12.27

Table 5: Efficiency analysis in diffusion model-based methods.

Methods	Sources	UIEB				LSUI				Methods	Sources	BAID			
		PSNR $\uparrow$	SSIM $\uparrow$	UCIQE $\uparrow$	UIQM $\uparrow$	PSNR $\uparrow$	SSIM $\uparrow$	UCIQE $\uparrow$	UIQM $\uparrow$			PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	FID $\downarrow$
FUGAN [28]	IRAL20	17.41	0.842	0.527	2.614	22.16	0.837	0.576	2.667	EnGAN [29]	TIP21	17.96	0.819	0.182	43.55
EnGAN [29]	TIP21	17.73	0.833	0.529	2.465	19.30	0.851	0.587	2.817	RUAS [42]	CVPR21	18.92	0.813	0.262	40.07
Ucolor [37]	TIP21	20.78	0.868	0.537	3.049	22.91	0.886	0.594	2.735	URetinetex [65]	CVPR22	19.08	0.845	0.206	42.26
S-uwnet [54]	AAA121	18.28	0.855	0.544	2.942	20.89	0.875	0.582	2.746	SNR-Net [69]	CVPR22	20.86	0.860	0.213	39.73
PUIE [16]	ECCV22	21.38	0.882	0.566	<b>3.021</b>	23.70	0.902	0.605	2.974	Restormer [74]	CVPR22	21.07	0.832	0.192	41.17
U-shape [56]	TIP23	22.91	<b>0.905</b>	0.592	2.896	24.16	<b>0.917</b>	0.603	3.022	RetFormer [3]	ICCV23	22.03	<b>0.862</b>	0.173	45.27
PUGAN [8]	TIP23	<b>23.05</b>	0.897	0.608	2.902	25.06	0.916	<b>0.629</b>	3.106	CLIP-LIT [41]	ICCV23	21.13	0.853	<b>0.159</b>	<b>37.30</b>
ADP [88]	LICV23	22.90	0.892	<b>0.621</b>	3.005	24.28	0.913	0.626	3.075	Diff-Retinex [73]	ICCV23	<b>22.07</b>	0.861	0.160	38.07
NU2Net [18]	AAA123	22.38	0.903	0.587	2.936	<b>25.97</b>	0.908	0.615	<b>3.112</b>	DiffR [66]	ICCV23	21.10	0.835	0.175	40.35
Reti-Diff (Ours)	—	<b>24.12</b>	<b>0.910</b>	<b>0.631</b>	<b>3.088</b>	<b>28.10</b>	<b>0.929</b>	<b>0.646</b>	<b>3.208</b>	Reti-Diff (Ours)	—	<b>23.19</b>	<b>0.876</b>	<b>0.147</b>	<b>27.47</b>

**Table 2:** Results on the underwater image enhancement task.**Table 3:** Results on the backlit image enhancement task.

Methods	Sources	DICM		LIME		MEF		NPE		VV	
		PI $\downarrow$	NIQE $\downarrow$	PI $\downarrow$	NIQE $\downarrow$	PI $\downarrow$	NIQE $\downarrow$	PI $\downarrow$	NIQE $\downarrow$	PI $\downarrow$	NIQE $\downarrow$
EnGAN [29]	TIP21	4.173	4.064	3.669	4.593	4.015	4.705	3.226	3.993	3.386	4.047
KimD++ [78]	IJCV21	3.835	3.898	3.785	4.908	4.016	4.557	3.179	3.915	3.773	3.822
SNR-Net [69]	CVPR22	3.585	4.715	3.753	5.937	3.677	6.449	3.278	6.446	3.503	9.506
DCC-Net [80]	CVPR22	3.630	3.709	<b>3.312</b>	<b>4.425</b>	<b>3.424</b>	4.598	<b>2.878</b>	<b>3.706</b>	3.615	<b>3.286</b>
UHDFor [38]	ICLR23	3.684	4.575	4.124	4.430	3.813	4.231	3.135	3.867	3.319	4.330
PairLIE [17]	CVPR23	3.685	4.034	3.387	4.587	4.133	4.065	3.726	4.187	<b>3.334</b>	3.574
GDP [12]	CVPR23	<b>3.552</b>	4.358	4.115	4.891	3.694	4.609	3.097	4.032	3.431	4.683
GSAD [30]	NIPS23	—	<b>3.465</b>	—	4.517	—	<b>3.815</b>	—	3.806	—	3.355
Reti-Diff (Ours)	—	<b>2.351</b>	<b>3.255</b>	<b>2.837</b>	<b>3.693</b>	<b>3.308</b>	<b>3.792</b>	<b>2.599</b>	<b>3.384</b>	<b>3.341</b>	<b>3.000</b>

**Table 4:** Results on the real-world illumination degradation image restoration task.

textures of wooded areas. Besides, we also compare the efficiency of the diffusion model-based methods. As presented in Tab. 5, despite having the second smallest parameter count, our Reti-Diff has the lowest MACs, highest FPS, and superior performance (see Tab. 1). This efficiency can be attributed to our utilization of the diffusion model within a low-dimensional compact latent space. For fairness, results from the compared methods are generated by their provided models under the same settings with no GT-mean strategy.

**Underwater Image Enhancement.** We extend our evaluation to encompass two widely-used underwater image enhancement datasets: *UIEB* [39] and *LSUI* [56]. In addition to PSNR and SSIM, we employ two metrics specifically tailored for underwater images, namely UCIQE [71] and UIQM [55], to assess the performance of the ten enhancement approaches. In all cases, higher values indicate superior performance. The results are presented in Tab. 2. As showcased in Tab. 2, our method achieved the highest overall performance and outperformed the second-best method (PUGAN) by 4.48%. A qualitative analysis is presented in Fig. 4, illustrating our method’s ability to correct underwater color aberrations and highlight fine texture details.

**Backlit Image Enhancement.** Following CLIP-LIT [41], we select the *BAID* [48] dataset for training the network with an image size of  $256 \times 256$ . In addition to PSNR and SSIM, our evaluation incorporates LPIPS [77] and FID [26] as metrics for evaluation, where lower LPIPS and FID denote superior performance. The evaluation results are reported in Tab. 3. As demonstrated in Tab. 3, our method excels in all metrics and generally outperformed the second-best method (CLIP-LIT) by 6.03%. Furthermore, a visual comparison in Fig. 5 provides additional evidence of our superiority in detail reconstruction and color correction.

**Real-world Illumination Degradation Image Restoration.** We also explore the applicability of our method in real-world IDIR tasks. Following the practice of CIDNet [13], we selected five commonly-used real-world datasets,



Datasets	Metrics	RGformer				Train w/o joint	Ours	Datasets	Metrics	w/o Reti-nex prior	w/ Reflect-ance prior	w/ Illumina-tion prior	w/ Retinex prior (Ours)
		RLDLM w/o RLDLM	w/o DFA	w/o RG-MCA	w/o $D_a(\cdot)$								
$L-v2-r$	PSNR	21.25	22.26	21.73	22.58	22.83	22.97	$L-v2-r$	PSNR	21.63	22.13	22.35	22.97
	SSIM	0.822	0.840	0.840	0.847	0.853	0.858		SSIM	0.830	0.842	0.839	0.858
$L-v2-s$	PSNR	25.38	26.49	25.92	26.80	27.18	27.53	$L-v2-s$	PSNR	26.25	26.62	27.02	27.53
	SSIM	0.918	0.925	0.913	0.944	0.947	0.951		SSIM	0.939	0.945	0.941	0.951

(a) Break down ablation.

(b) Effect of our Retinex prior.

**Table 6:** Ablation study on the low-light image enhancement task.

Datasets	Metrics	Res [74]		Ret [3]		Res+RLDM	Ret+RLDM	PSNR	SSIM	Gain
		Res	Res+RLDM	Ret	Ret+RLDM					
$L-v2-s$	PSNR	21.41	24.15	25.67	26.81	—	—	21.41	0.830	—
	SSIM	0.830	0.862	0.930	0.942					
	Gain	—	8.33%	—	2.87%					
$L-v2-r$	PSNR	19.94	21.56	22.80	23.16	—	—	19.94	0.827	—
	SSIM	0.827	0.837	0.840	0.849					
	Gain	—	4.67%	—	1.33%					

**Table 7:** Generalization of our RLDM. “Res” and “Ret” are Restormer and Retformer.**Fig. 6:** Ablation study of the number of iterations in RLDM on  $LOL-v2-syn$ .

*i.e.*, *DICM* [36], *LIME* [19], *MEF* [61], *NPE* [50], and *VV*<sup>7</sup>, which only have the low-quality images without paired high-quality ground-truth. Therefore, akin to [13], we leverage models pretrained on the *LOL-v2-syn* dataset for inference and select PI [2] and NIQE [52] as evaluation metrics. In both metrics, lower scores indicate better results. As presented in Tab. 4, our method achieves optimal results and surpasses the second-based method (DCC-Net [80]) by 13.39%. This verifies the generalizability of our Reti-Diff in addressing unknown degradation scenarios. Note that all the methods abandon the GT-mean strategy.

### 4.3 Ablation Study

We conduct ablation studies on the low-light image enhancement task with the  $L-v2-r$  and  $L-v2-s$  datasets, which are short for *LOL-v2-real* and *LOL-v2-syn*.

**Effect of RLDM.** As illustrated in Tab. 6a, we ablate RLDM by directly removing RLDM or retraining RLDM in the RGB domain, *i.e.*, w/o Retinex, rather than in the reflectance and illumination domain (RGformer is guided by one RGB prior instead of the Retinex priors in this time). The two modifications result in significant drops in performance. This outcome underscores the critical role of RLDM in enhancing the restoration process. Furthermore, to assess the generalizability of RLDM, we conducted additional experiments by replacing our RGformer with two transformer-based frameworks, namely Res (Restormer [74]) and Ret (Retformer [3]). Note that the training settings are kept consistent with our Reti-Diff. The results are presented in Tab. 7. Tab. 7 reveals that RLDM significantly improves the performance of both frameworks, where “Gain” is the average gain of PSNR and SSIM. This demonstrates that our RLDM serves as a plug-and-play module with strong generalization capabilities.

**Effect of RGformer.** We conduct an analysis to assess the impact of our RGformer, and the results are presented in Tab. 6a. In this study, we systematically removed critical components, such as DFA, RG-MCA, and the auxiliary decoder  $D_a(\cdot)$ , from the model architecture. The outcomes of this ablation study

<sup>7</sup> <https://sites.google.com/site/vonikakis/datasets>

Methods	$L-v1$	$L-v2-r$	$L-v2-s$	$SID$	Mean	Methods (AP)	Bicycle	Boat	Bottle	Bus	Car	Cat	Chair	Cup	Dog	Motor	People	Table	Mean
KinD [79]	2.31	2.25	2.46	2.33	2.34	Baseline	74.7	64.9	70.7	84.2	79.7	47.3	58.6	67.1	64.1	66.2	73.9	45.7	66.4
EnGAN [29]	2.63	1.69	2.23	1.24	1.95	RetinexNet [64]	72.8	66.4	67.3	87.5	80.6	52.8	60.0	67.8	68.5	69.3	71.3	46.2	67.5
RUAS [42]	3.57	3.06	3.01	2.23	2.97	KinD [79]	73.2	67.1	64.6	86.8	79.5	58.7	63.4	67.5	67.4	62.3	75.5	51.4	68.1
Restormer [74]	3.26	3.32	3.41	2.53	3.13	MIRNet [75]	74.9	69.7	68.3	89.7	77.6	57.8	56.9	66.4	69.7	64.6	74.6	53.4	68.6
Uretinex [65]	3.82	3.98	3.70	3.28	3.70	RUAS [42]	75.7	71.2	73.5	90.7	80.1	59.3	67.0	66.3	68.3	66.9	72.6	50.6	70.2
SNR-Net [69]	3.76	4.12	3.58	3.42	3.72	Restormer [74]	77.0	71.0	68.8	91.6	77.1	62.5	57.3	68.0	69.6	69.2	74.6	49.7	69.7
CUE [81]	3.62	3.81	3.28	3.09	3.45	SCI [51]	73.4	68.0	69.5	86.2	74.5	63.1	59.5	61.0	67.3	63.9	73.2	47.3	67.2
Retformer [3]	3.35	4.02	3.71	3.35	3.61	SNR-Net [69]	78.3	74.2	74.5	89.6	82.7	66.8	66.3	62.5	74.7	63.1	73.3	57.2	71.9
Ours	4.05	4.33	3.92	3.75	4.01	Retformer [3]	78.1	74.5	74.2	91.2	82.2	65.0	63.3	67.0	75.4	68.6	75.3	55.6	72.5
						Ours	82.0	77.9	76.4	92.2	83.3	69.6	67.4	74.4	75.5	74.3	78.3	57.9	75.8

Table 8: User study.

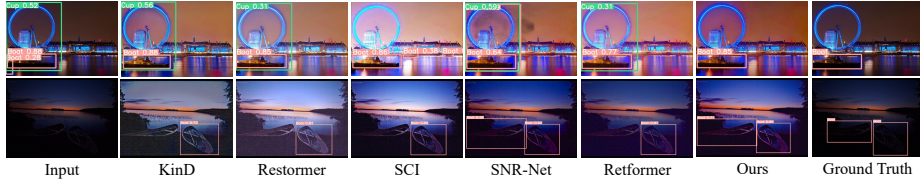
Table 9: Low-light image detection on *ExDark* [45].

Fig. 7: Results on the low-light object detection task.

clearly indicate that the performance deteriorates when these components are removed, highlighting their essential role in the system. Additionally, in Tab. 6a, we conduct an evaluation to affirm the significance of joint training in our approach. This analysis reinforces the importance of the joint training process.

**Ablations on iteration number.** The number of iterations in the diffusion model plays a crucial role in determining the method’s efficiency. To explore this, we conducted experiments with different iteration numbers for Reti-Diff, specifically  $T$  values selected from the set  $\{1, 2, 4, 8, 16, 32\}$ . We adjusted  $\beta^t$  as defined in Eq. (9) accordingly. The results in terms of PSNR for different iterations, as shown in Fig. 6, illustrate that Reti-Diff exhibits rapid convergence and generates stable guidance priors with just 4 iterations. This efficiency is attributed to our application of the diffusion model within the compact latent space.

#### 4.4 User Study and Downstream Tasks

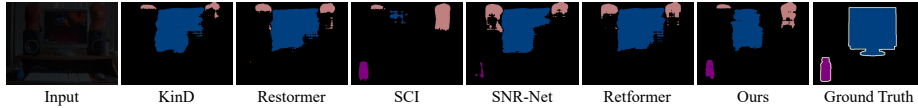
**User Study.** We conduct a user study to assess the subjective visual perception of low-light image enhancement. In this study, 29 human subjects are invited to assign scores to the enhanced results based on four criteria: (1) The presence of underexposed or overexposed regions. (2) The existence of color distortion. (3) The occurrence of undesired noise or artifacts. (4) The inclusion of essential structural details. Participants rate the results on a scale from 1 (worst) to 5 (best). Each low-light image is presented alongside its enhanced results, with the names of the enhancement methods concealed. The scores are reported in Tab. 8, where our method receives the highest scores across all four datasets. This highlights our effectiveness in generating visually appealing results.

**Low-light Object Detection.** The enhanced images are expected to have better downstream performance than the original ones. We first verify this on low-light object detection. Following [3], all compared methods are performed on *ExDark* [45] with YOLOv3, which is retrained from scratch with their corresponding enhanced results. As shown in Tab. 9, our Reti-Diff exhibits a substantial advantage over existing methods and the performance of our method

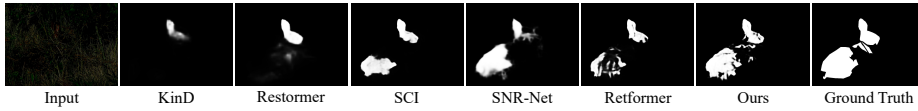


Methods (IoU)	Bicycle	Boat	Bottle	Bus	Car	Cat	Chair	Dog	Horse	People	Mean	Methods	COD10K				NC4K			
	$M \downarrow$	$F_\beta \uparrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$	$M \downarrow$	$F_\beta \uparrow$	$E_\phi \uparrow$	$S_\alpha \uparrow$												
Baseline	43.5	36.3	48.6	70.5	67.3	46.6	11.2	42.4	56.7	57.8	48.1	Baseline	0.050	0.625	0.812	0.756	0.071	0.733	0.816	0.763
RetinexNet [64]	48.6	41.7	51.7	77.6	68.3	52.7	15.8	46.3	60.2	62.3	52.5	RetinexNet [64]	0.041	0.667	0.845	0.789	0.055	0.750	0.842	0.819
KinD [79]	51.3	40.2	53.2	76.8	69.4	50.8	14.6	47.3	60.3	60.9	52.5	KinD [79]	0.039	0.673	0.849	0.792	0.052	0.762	0.875	0.822
MIRNet [75]	50.3	42.9	47.4	73.6	62.7	50.4	15.8	46.3	61.0	63.3	51.4	MIRNet [75]	0.037	0.697	0.857	0.799	<b>0.049</b>	<b>0.802</b>	0.888	0.833
RUAS [42]	53.0	37.3	50.4	71.3	72.3	47.6	15.9	50.8	63.6	60.8	52.3	RUAS [42]	<b>0.036</b>	0.705	0.861	0.803	0.051	0.795	0.883	0.827
Restormer [74]	53.8	43.8	51.4	68.7	66.8	52.6	21.6	<b>54.8</b>	59.8	63.3	53.7	Restormer [74]	<b>0.036</b>	0.700	0.859	0.800	0.050	0.792	0.880	0.830
SCI [51]	54.5	46.3	57.2	78.4	73.3	49.1	22.8	49.0	62.1	66.9	56.0	SCI [51]	0.037	<b>0.710</b>	0.863	0.805	0.051	0.782	0.880	0.836
SNR-Net [69]	<b>57.7</b>	<b>48.6</b>	<b>59.5</b>	<b>81.3</b>	<b>74.8</b>	50.2	<b>24.4</b>	50.7	<b>64.3</b>	68.7	<b>58.0</b>	SNR-Net [69]	<b>0.036</b>	0.703	<b>0.865</b>	0.803	<b>0.049</b>	0.801	<b>0.892</b>	<b>0.838</b>
Retformer [3]	50.9	47.7	58.6	77.2	68.1	<b>53.2</b>	17.4	52.0	61.3	<b>71.5</b>	55.8	Retformer [3]	0.037	0.682	0.861	<b>0.806</b>	0.052	0.766	0.881	0.832
Ours	<b>59.8</b>	<b>51.5</b>	<b>62.1</b>	<b>85.5</b>	<b>76.6</b>	<b>57.7</b>	<b>28.9</b>	<b>56.3</b>	<b>66.2</b>	<b>73.4</b>	<b>61.8</b>	Ours	<b>0.034</b>	<b>0.725</b>	<b>0.880</b>	<b>0.813</b>	<b>0.047</b>	<b>0.804</b>	<b>0.897</b>	<b>0.841</b>

**Table 10:** Low-light semantic segmentation, where images are darkened by [76]. **Table 11:** Low-light concealed object segmentation.



**Fig. 8:** Results on the low-light semantic segmentation task.



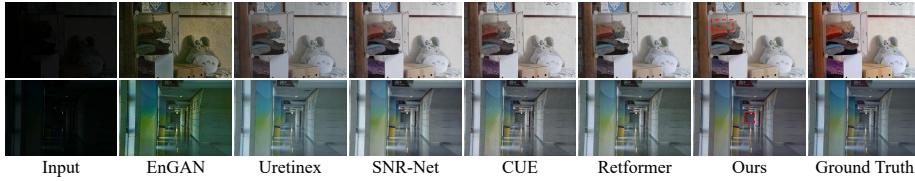
**Fig. 9:** Results on the low-light concealed object segmentation task.

surpasses that of the second-best method, Retformer [3], by 4.72%, which verifies our efficacy in facilitating high-level vision understanding.

**Low-light Image Segmentation.** We extend our experiment to segmentation tasks, *i.e.*, semantic segmentation and concealed object segmentation. Following the practice in detection, we also retrain the segmentor for each method. This means that each method’s enhanced results are segmented by the corresponding segmentor with specific weights. We argue this could better exploit the potential of image enhancement methods as a degraded data restoration module.

For semantic segmentation, following [31], we apply image darkening to samples from the *VOC* [9] dataset according to [76]. We then employ Mask2Former [7] to perform segmentation on the enhanced results of these darkened images. We select Intersection over Union (IoU) for evaluation, and the results are presented in Tab. 10. As shown in Tab. 10, our method achieves the highest performance across all classes, surpassing the second-best method by 7.53%.

We further venture into concealed object segmentation (COS) on two widely-used datasets, *COD10K* [10] and *NC4K* [49], which represents a challenging segmentation task aimed at delineating objects with inherent background similarity. We also apply image darkening [76] and enlist the cutting-edge COS segmentor, FEDER [22], to perform segmentation on the enhanced results. We evaluate the results using four metrics: mean absolute error ( $M$ ), adaptive F-measure ( $F_\beta$ ), mean E-measure ( $E_\phi$ ), and structure measure ( $S_\alpha$ ), which are presented in Tab. 11. As depicted in Tab. 11, our method exhibits superior performance compared to the second-best method, SNR-Net, with a margin of 2.16% on average. Note that it is a notable improvement in COS. Collectively, the exceptional results achieved in these two segmentation tasks substantiate our proficiency in recovering image-level illumination degraded information.



**Fig. 10:** Failure cases. Our results show blurred texture details in the dashed boxes.

## 5 Discussions

Our Reti-Diff is the first LDM-based solution specifically designed for the IDIR task, setting it apart from existing LDM-based methods applied in other tasks. To illustrate the distinctions, we compare it with a general enhancement method, DiffIR [66]: **(1) Motivation.** Reti-Diff targets enhancing details and correcting degraded illumination. Thus, we enable RLDM to learn Retinex knowledge and generate Retinex priors from the low-quality input. We contend that relying solely on priors extracted from the RGB domain struggles to fully represent valuable texture details and correct illumination cues, leading to suboptimal restoration performance. To verify this, we substitute our RLDM for the LDM structure used in DiffIR. In *LOL-v2-syn*, we observe that the PSNR rises from 24.76 to 26.14 and the SSIM increases from 0.921 to 0.933. **(2) Implementation.** Apart from proposing RLDM to extract Retinex priors, we further modify the structure of RGformer to implicitly model the Retinex theory at the feature level and introduce an auxiliary decoder to reconstruct the decomposed Retinex components to the RGB domain. **(3) Performance.** As shown in Tab. 1, our Reti-Diff significantly outperforms DiffIR [66] by 20.6% on average.

## 6 Limitations and Future Work

As shown in Fig. 10, our Reti-Diff encounters challenges in simultaneously recovering illumination and restoring texture details when the low-quality inputs exhibit severe illumination degradation. This issue persists across existing methods and remains unresolved. We attribute this to the loss of texture information during illumination recovery. To address this limitation in future research, we propose excavating texture priors from other domains, *e.g.*, the frequency domain. These priors can complement the reflectance priors extracted from the RGB domain, enhancing the preservation of critical texture features. Additionally, we consider the use of multimodal data [11] to aid in improving image reconstruction performance, such as using infrared images [21, 32, 68] to aid in low-light visible image enhancement. Besides, We will explore whether our approach is downstream task-friendly with more segmentation algorithms [23, 25, 67]. We also aim to extend our approach to tackle IDIR problems afflicted by other types of degradation, such as haze and motion blur, using some domain adaptation strategies [57, 58]. These endeavors will further advance the capabilities and applicability of Reti-Diff in real-world scenarios.

## 7 Conclusions

To balance generation capability and computational efficiency, our approach adopts DM within a compact latent space to generate guidance priors. Specifically, we introduce RLDM to extract Retinex priors, which are subsequently supplied to RGformer for feature decomposition. This process ensures precise detailed reconstruction and effective illumination correction. RGformer then refines and aggregates the decomposed features, enhancing the robustness in handling complex degradation scenarios. Our approach is extensively validated through experiments, establishing the clear superiority of the proposed Reti-Diff.

## References

1. Abdullah-Al-Wadud, M., Kabir, M.H., Dewan, M.A.A., Chae, O.: A dynamic histogram equalization for image contrast enhancement. *IEEE transactions on consumer electronics* **53**(2), 593–600 (2007)
2. Blau, Y., Mechrez, R., Timofte, R., Michaeli, T., Zelnik-Manor, L.: The 2018 pirm challenge on perceptual image super-resolution. In: *ECCV*. pp. 0–0 (2018)
3. Cai, Y., Bian, H., Lin, J., Wang, H., Timofte, R., Zhang, Y.: Retinexformer: One-stage retinex-based transformer for low-light image enhancement. In: *ICCV*. pp. 12504–12513 (2023)
4. Chen, C., Chen, Q., Do, M.N., Koltun, V.: Seeing motion in the dark. In: *ICCV*. pp. 3185–3194 (2019)
5. Chen, H., Wang, Y., Guo, T., Xu, C., Deng, Y., Liu, Z., Ma, S., Xu, C., Xu, C., Gao, W.: Pre-trained image processing transformer. In: *CVPR*. pp. 12299–12310 (2021)
6. Chen, Z., Zhang, Y., Liu, D., Xia, B., Gu, J., Kong, L., Yuan, X.: Hierarchical integration diffusion model for realistic image deblurring. In: *NeurIPS* (2023)
7. Cheng, B., Misra, I., Schwing, A.G., Kirillov, A., Girdhar, R.: Masked-attention mask transformer for universal image segmentation. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 1290–1299 (2022)
8. Cong, R., Yang, W., Zhang, W., Li, C., Guo, C.L., Huang, Q., Kwong, S.: Pagan: Physical model-guided underwater image enhancement using gan with dual-discriminators. *IEEE Transactions on Image Processing* (2023)
9. Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. *International journal of computer vision* **88**, 303–338 (2010)
10. Fan, D.P., Ji, G.P., Cheng, M.M., Shao, L.: Concealed object detection. *IEEE transactions on pattern analysis and machine intelligence* **44**(10), 6024–6042 (2021)
11. Fang, C.Y., Han, X.F.: Joint geometric-semantic driven character line drawing generation. In: *ICMR*. pp. 226–233 (2023)
12. Fei, B., Lyu, Z., Pan, L., Zhang, J., Yang, W., Luo, T., Zhang, B., Dai, B.: Generative diffusion prior for unified image restoration and enhancement. In: *CVPR*. pp. 9935–9946 (2023)
13. Feng, Y., Zhang, C., Wang, P., Wu, P., Yan, Q., Zhang, Y.: You only need one color space: An efficient network for low-light image enhancement. *arXiv preprint arXiv:2402.05809* (2024)
14. Fu, X., Zeng, D., Huang, Y., Liao, Y., Ding, X., Paisley, J.: A fusion-based enhancing method for weakly illuminated images. *Signal Processing* **129**, 82–96 (2016)
15. Fu, X., Zeng, D., Huang, Y., Zhang, X.P., Ding, X.: A weighted variational model for simultaneous reflectance and illumination estimation. In: *CVPR*. pp. 2782–2790 (2016)
16. Fu, Z., Wang, W., Huang, Y., Ding, X., Ma, K.K.: Uncertainty inspired underwater image enhancement. In: *ECCV*. pp. 465–482. Springer (2022)
17. Fu, Z., Yang, Y., Tu, X., Huang, Y., Ding, X., Ma, K.K.: Learning a simple low-light image enhancer from paired low-light instances. In: *CVPR*. pp. 22252–22261 (2023)
18. Guo, C., Wu, R., Jin, X., Han, L., Zhang, W., Chai, Z., Li, C.: Underwater ranker: Learn which is better and how to be better. In: *AAAI*. vol. 37, pp. 702–709 (2023)

19. Guo, X., Li, Y., Ling, H.: Lime: Low-light image enhancement via illumination map estimation. *IEEE Trans. Image Process.* **26**(2), 982–993 (2016)
20. He, C., Li, K., Xu, G., Yan, J., Tang, L., Zhang, Y., Wang, Y., Li, X.: Hqg-net: Unpaired medical image enhancement with high-quality guidance. *IEEE Transactions on Neural Networks and Learning Systems* (2023)
21. He, C., Li, K., Xu, G., Zhang, Y., Hu, R., Guo, Z., Li, X.: Degradation-resistant unfolding network for heterogeneous image fusion. In: *ICCV*. pp. 12611–12621 (2023)
22. He, C., Li, K., Zhang, Y., Tang, L., Zhang, Y., Guo, Z., Li, X.: Camouflaged object detection with feature decomposition and edge reconstruction. In: *CVPR*. pp. 22046–22055 (2023)
23. He, C., Li, K., Zhang, Y., Xu, G., Tang, L., Zhang, Y., Guo, Z., Li, X.: Weakly-supervised concealed object segmentation with sam-based pseudo labeling and multi-scale feature grouping. *NIPS* **36** (2024)
24. He, C., Li, K., Zhang, Y., Zhang, Y., Guo, Z., Li, X., Danelljan, M., Yu, F.: Strategic preys make acute predators: Enhancing camouflaged object detectors by generating camouflaged objects (2024)
25. He, C., Wang, X., Deng, L., Xu, G.: Image threshold segmentation based on glle histogram. In: *CPSCoM*. pp. 410–415. *IEEE* (2019)
26. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. *NeurIPS* **30** (2017)
27. Huang, S.C., Cheng, F.C., Chiu, Y.S.: Efficient contrast enhancement using adaptive gamma correction with weighting distribution. *IEEE transactions on image processing* **22**(3), 1032–1041 (2012)
28. Islam, M.J., Xia, Y., Sattar, J.: Fast underwater image enhancement for improved visual perception. *IEEE Robotics and Automation Letters* **5**(2), 3227–3234 (2020)
29. Jiang, Y., Gong, X., Liu, D., Cheng, Y., Fang, C., Shen, X., Yang, J., Zhou, P., Wang, Z.: Enlightengan: Deep light enhancement without paired supervision. *IEEE transactions on image processing* **30**, 2340–2349 (2021)
30. Jinhui, H., Zhu, Z., Hou, J., Hui, L., Zeng, H., Yuan, H.: Global structure-aware diffusion process for low-light image enhancement. In: *NeurIPS* (2023)
31. Ju, M., Guo, C.A., Chen, C., Pan, J., Tang, J., Tao, D.: Sllen: Semantic-aware low-light image enhancement network. *arXiv preprint arXiv:2211.11571* (2022)
32. Ju, M., He, C., Liu, J., Kang, B., Su, J., Zhang, D.: Ivf-net: An infrared and visible data fusion deep network for traffic object enhancement in intelligent transportation systems. *IEEE Transactions on Intelligent Transportation Systems* **24**(1), 1220–1234 (2022)
33. Kingma, D., Salimans, T., Poole, B., Ho, J.: Variational diffusion models. *NeurIPS* **34**, 21696–21707 (2021)
34. Kingma, D.P., Welling, M.: Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114* (2013)
35. Land, E.H.: The retinex theory of color vision. *Scientific american* **237**(6), 108–129 (1977)
36. Lee, C., Lee, C., Kim, C.S.: Contrast enhancement based on layered difference representation of 2d histograms. *IEEE Trans. Image Process.* **22**(12), 5372–5384 (2013)
37. Li, C., Anwar, S., Hou, J., Cong, R., Guo, C., Ren, W.: Underwater image enhancement via medium transmission-guided multi-color space embedding. *IEEE Transactions on Image Processing* **30**, 4985–5000 (2021)

38. Li, C., Guo, C.L., Zhou, M., Liang, Z., Zhou, S., Feng, R., Loy, C.C.: Embedding-fourier for ultra-high-definition low-light image enhancement. In: ICLR (2023)
39. Li, C., Guo, C., Ren, W., Cong, R., Hou, J., Kwong, S., Tao, D.: An underwater image enhancement benchmark dataset and beyond. *IEEE Transactions on Image Processing* **29**, 4376–4389 (2019)
40. Li, M., Liu, J., Yang, W., Sun, X., Guo, Z.: Structure-revealing low-light image enhancement via robust retinex model. *IEEE Transactions on Image Processing* **27**(6), 2828–2841 (2018)
41. Liang, Z., Li, C., Zhou, S., Feng, R., Loy, C.C.: Iterative prompt learning for unsupervised backlit image enhancement. In: ICCV. pp. 8094–8103 (2023)
42. Liu, R., Ma, L., Zhang, J., Fan, X., Luo, Z.: Retinex-inspired unrolling with cooperative prior architecture search for low-light image enhancement. In: CVPR. pp. 10561–10570 (2021)
43. Liu, Y., Huang, T., Dong, W., Wu, F., Li, X., Shi, G.: Low-light image enhancement with multi-stage residue quantization and brightness-aware attention. In: ICCV. pp. 12140–12149 (2023)
44. Locatello, F., Weissenborn, D., Unterthiner, T., Mahendran, A., Heigold, G., Uszkoreit, J., Dosovitskiy, A., Kipf, T.: Object-centric learning with slot attention. *NeurIPS* **33**, 11525–11538 (2020)
45. Loh, Y.P., Chan, C.S.: Getting to know low-light images with the exclusively dark dataset. *Computer Vision and Image Understanding* **178**, 30–42 (2019)
46. Lore, K.G., Akintayo, A., Sarkar, S.: Llnet: A deep autoencoder approach to natural low-light image enhancement. *Pattern Recognition* **61**, 650–662 (2017)
47. Loshchilov, I., Hutter, F.: Stochastic gradient descent with warm restarts. In: ICLR. pp. 1–16 (2016)
48. Lv, X., Zhang, S., Liu, Q., Xie, H., Zhong, B., Zhou, H.: Backlitnet: A dataset and network for backlit image enhancement. *Computer Vision and Image Understanding* **218**, 103403 (2022)
49. Lv, Y., Zhang, J., Dai, Y., Li, A., Liu, B., Barnes, N., Fan, D.P.: Simultaneously localize, segment and rank the camouflaged objects. In: CVPR. pp. 11591–11601 (2021)
50. Ma, K., Zeng, K., Wang, Z.: Perceptual quality assessment for multi-exposure image fusion. *IEEE Trans. Image Process.* **24**(11), 3345–3356 (2015)
51. Ma, L., Ma, T., Liu, R., Fan, X., Luo, Z.: Toward fast, flexible, and robust low-light image enhancement. In: CVPR. pp. 5637–5646 (2022)
52. Mittal, A., Soundararajan, R., Bovik, A.C.: Making a “completely blind” image quality analyzer. *IEEE Signal Processing Lett.* **20**(3), 209–212 (2012)
53. Moorthy, A.K., Bovik, A.C.: A two-step framework for constructing blind image quality indices. *IEEE Signal processing letters* **17**(5), 513–516 (2010)
54. Naik, A., Swarnakar, A., Mittal, K.: Shallow-wnet: Compressed model for underwater image enhancement (student abstract). In: AAAI. vol. 35, pp. 15853–15854 (2021)
55. Panetta, K., Gao, C., Agaian, S.: Human-visual-system-inspired underwater image quality measures. *IEEE Journal of Oceanic Engineering* **41**(3), 541–551 (2015)
56. Peng, L., Zhu, C., Bian, L.: U-shape transformer for underwater image enhancement. *IEEE Transactions on Image Processing* (2023)
57. Tang, L., Li, K., He, C., Zhang, Y., Li, X.: Consistency regularization for generalizable source-free domain adaptation. In: ICCV. pp. 4323–4333 (2023)
58. Tang, L., Li, K., He, C., Zhang, Y., Li, X.: Source-free domain adaptive fundus image segmentation with class-balanced mean teacher. In: MICCAI. pp. 684–694. Springer (2023)

59. Ueng, N.T., Scharf, L.L.: The gamma transform: A local time-frequency analysis method. In: ACSSC. vol. 2, pp. 920–924. IEEE (1995)
60. Wang, R., Zhang, Q., Fu, C.W., Shen, X., Zheng, W.S., Jia, J.: Underexposed photo enhancement using deep illumination estimation. In: CVPR. pp. 6849–6857 (2019)
61. Wang, S., Zheng, J., Hu, H.M., Li, B.: Naturalness preserved enhancement algorithm for non-uniform illumination images. *IEEE Trans. Image Process.* **22**(9), 3538–3548 (2013)
62. Wang, Y., Liu, Z., Liu, J., Xu, S., Liu, S.: Low-light image enhancement with illumination-aware gamma correction and complete image modelling network. In: ICCV. pp. 13128–13137 (2023)
63. Wang, Z., Cun, X., Bao, J., Zhou, W., Liu, J., Li, H.: Uformer: A general u-shaped transformer for image restoration. In: CVPR. pp. 17683–17693 (2022)
64. Wei, C., Wang, W., Yang, W., Liu, J.: Deep retinex decomposition for low-light enhancement. arXiv preprint arXiv:1808.04560 (2018)
65. Wu, W., Weng, J., Zhang, P., Wang, X., Yang, W., Jiang, J.: Uretinex-net: Retinex-based deep unfolding network for low-light image enhancement. In: CVPR. pp. 5901–5910 (2022)
66. Xia, B., Zhang, Y., Wang, S., Wang, Y., Wu, X., Tian, Y., Yang, W., Van Gool, L.: Diffir: Efficient diffusion model for image restoration. In: ICCV (2023)
67. Xiao, F., Zhang, P., He, C., Hu, R., Liu, Y.: Concealed object segmentation with hierarchical coherence modeling. In: CAAI. pp. 16–27. Springer (2023)
68. Xu, G., He, C., Wang, H., Zhu, H., Ding, W.: Dm-fusion: Deep model-driven network for heterogeneous image fusion. *IEEE Transactions on Neural Networks and Learning Systems* (2023)
69. Xu, X., Wang, R., Fu, C.W., Jia, J.: Snr-aware low-light image enhancement. In: CVPR. pp. 17714–17724 (2022)
70. Xu, X., Wang, R., Lu, J.: Low-light image enhancement via structure modeling and guidance. In: CVPR. pp. 9893–9903 (2023)
71. Yang, M., Sowmya, A.: An underwater color image quality evaluation metric. *IEEE Transactions on Image Processing* **24**(12), 6062–6071 (2015)
72. Yang, W., Wang, W., Huang, H., Wang, S., Liu, J.: Sparse gradient regularized deep retinex network for robust low-light image enhancement. *IEEE Transactions on Image Processing* **30**, 2072–2086 (2021)
73. Yi, X., Xu, H., Zhang, H., Tang, L., Ma, J.: Diff-retinex: Rethinking low-light image enhancement with a generative diffusion model. In: ICCV. pp. 12302–12311 (2023)
74. Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H.: Restormer: Efficient transformer for high-resolution image restoration. In: CVPR. pp. 5728–5739 (2022)
75. Zamir, S.W., Arora, A., Khan, S., Hayat, M., Khan, F.S., Yang, M.H., Shao, L.: Learning enriched features for real image restoration and enhancement. In: ECCV. pp. 492–511. Springer (2020)
76. Zhang, F., Li, Y., You, S., Fu, Y.: Learning temporal consistency for low light video enhancement from single images. In: CVPR. pp. 4967–4976 (2021)
77. Zhang, R., Isola, P., Efros, A.A., Shechtman, E., Wang, O.: The unreasonable effectiveness of deep features as a perceptual metric. In: CVPR. pp. 586–595 (2018)
78. Zhang, Y., Guo, X., Ma, J., Liu, W., Zhang, J.: Beyond brightening low-light images. *Int. J. Comput. Vision* **129**, 1013–1037 (2021)
79. Zhang, Y., Zhang, J., Guo, X.: Kindling the darkness: A practical low-light image enhancer. In: ACM MM. pp. 1632–1640 (2019)

80. Zhang, Z., Zheng, H., Hong, R., Xu, M., Yan, S., Wang, M.: Deep color consistent network for low-light image enhancement. In: CVPR. pp. 1899–1908 (2022)
81. Zheng, N., Zhou, M., Dong, Y., Rui, X., Huang, J., Li, C., Zhao, F.: Empowering low-light image enhancer through customized learnable priors. In: ICCV. pp. 12559–12569 (2023)
82. Zhou, D., Yang, Z., Yang, Y.: Pyramid diffusion models for low-light image enhancement. arXiv preprint arXiv:2305.10028 (2023)
83. Zhou, J., Liu, Q., Jiang, Q., Ren, W., Lam, K.M., Zhang, W.: Underwater camera: Improving visual perception via adaptive dark pixel prior and color correction. *International Journal of Computer Vision* pp. 1–19 (2023)