

# CWNet: Causal Wavelet Network for Low-Light Image Enhancement

Tongshun Zhang<sup>1,2</sup> Pingping Liu<sup>1,2\*</sup> Yubing Lu<sup>1,2</sup> Mengen Cai<sup>1,2</sup> Zijian Zhang<sup>1,2</sup>  
Zhe Zhang<sup>1,2</sup> Qiuzhan Zhou<sup>3</sup>

<sup>1</sup>College of Computer Science and Technology, Jilin University

<sup>2</sup>Key Laboratory of Symbolic Computation and Knowledge Engineering of Ministry of Education

<sup>3</sup>College of Communication Engineering, Jilin University

{tszhang23, luyb24, caime24, zhezhang23}@mails.jlu.edu.cn,

{liupp, zhangzijian, zhouqz}@jlu.edu.cn

## Abstract

*Traditional Low-Light Image Enhancement (LLIE) methods primarily focus on uniform brightness adjustment, often neglecting instance-level semantic information and the inherent characteristics of different features. To address these limitations, we propose CWNet (Causal Wavelet Network), a novel architecture that leverages wavelet transforms for causal reasoning. Specifically, our approach comprises two key components: 1) Inspired by the concept of intervention in causality, we adopt a causal reasoning perspective to reveal the underlying causal relationships in low-light enhancement. From a global perspective, we employ a metric learning strategy to ensure causal embeddings adhere to causal principles, separating them from non-causal confounding factors while focusing on the invariance of causal factors. At the local level, we introduce an instance-level CLIP semantic loss to precisely maintain causal factor consistency. 2) Based on our causal analysis, we present a wavelet transform-based backbone network that effectively optimizes the recovery of frequency information, ensuring precise enhancement tailored to the specific attributes of wavelet transforms. Extensive experiments demonstrate that CWNet significantly outperforms current state-of-the-art methods across multiple datasets, showcasing its robust performance across diverse scenes. Code is available at [CWNet](#).*

## 1. Introduction

LLIE is essential in computer vision, addressing challenges like dimness and detail loss that degrade image quality. While traditional methods like gamma correction [29], Retinex theory [31], and histogram equalization [32] strug-

gle with non-uniform lighting and extreme darkness, deep learning approaches [3, 47, 48, 53, 54, 61] offer improved adaptability and performance. However, they often fail to fully exploit feature modeling and semantic information.

Frequency-based methods present a promising avenue for LLIE by effectively separating and enhancing high- and low-frequency information, improving detail and brightness while isolating noise. Nonetheless, existing methods [14, 19, 38, 53, 61] treat frequency features uniformly, which limits their potential. Additionally, maintaining color and semantic consistency is a significant challenge, as many advanced methods [2, 17, 39, 56] often overlook these aspects, leading to visually unnatural or semantically inaccurate results. This paper addresses these gaps by exploring two key questions:

**Firstly, how can we ensure consistency in color and semantic information while improving lighting conditions?** Current methods [43, 59] often rely on color histogram-based losses to maintain color consistency, while SCLLE [23] utilize downstream semantic segmentation consistency loss to enhance semantic brightness and color consistency. SKF [43] further improves semantic consistency at the feature level by extracting intermediate features through a semantic segmentation network. In contrast, the visual-language pre-trained model CLIP [33] demonstrates superior performance in maintaining color and semantic consistency. Many methods [17, 49, 55] leverage CLIP to learn diverse features, achieving semantically guided enhancement. However, these approaches primarily focus on global semantic and color consistency, lacking the ability to ensure instance-level consistency.

**Secondly, how can we establish a robust model that fully exploits frequency domain features?** The two commonly used frequency domain transformations are Fourier transform and wavelet transform. Fourier-based methods [14, 38] excels in capturing global information by ampli-

\*Corresponding author

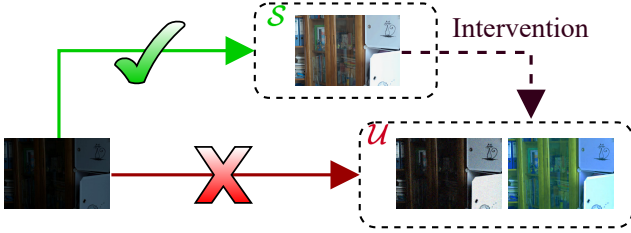


Figure 1. Structural causal model (SCM) for LLIE.

fying low-frequency components, enhancing overall brightness. However, its lack of spatial locality limits its ability to preserve high-frequency details like edges and textures, often resulting in brighter but less detailed images. Recent works [19, 53] have improved detail preservation by incorporating phase processing, but challenges remain in achieving fine-grained detail enhancement and spatial coherence. In contrast, wavelet transforms provide superior spatial locality, effectively separating image content from noise and enhancing edge and texture details. However, wavelet-based methods [17, 61] do not fully leverage the unique characteristics of the frequency domain, which limits their recovery potential.

To address these limitations, we propose a Structural Causal Model (SCM) [35] for Low-Light Image Enhancement (LLIE) based on causal inference, as shown in Fig. 1. Within this framework, we establish the core objective of the LLIE task: to maintain consistency in causal factors  $\mathcal{S}$  (semantic information). Meanwhile, non-causal factors  $\mathcal{U}$  (color and brightness anomalies) need to be filtered out. This causal perspective enables us to effectively distinguish meaningful semantic content from confounding degradations. Building on this causal analysis, we propose a two-level strategy. At the global level, we obtain non-causal factors  $\mathcal{U}$  through intervention procedures. Subsequently, we employ a causally-guided metric learning approach to filter out non-causal factors in the latent space. At the local level, we introduce an instance-level CLIP semantic loss to maintain fine-grained semantic consistency for each instance, achieving the objective of ensuring that causal factors  $\mathcal{S}$  remain consistent.

Based on the SCM, we meticulously propose a wavelet-based backbone network to support this causal concept, referred to as the Causal Wavelet Network (CWNNet). CWNNet incorporates a Hierarchical Feature Restoration Block (HFRB) after each sampling layer, consisting of three components: a Feature Extraction (FE), a High-Frequency Enhancement Block (HFEB), and a Low-Frequency Enhancement Block (LFEB). The FE adaptively extracts wavelet frequency domain features and compensates for missing information through interaction. HFEB, inspired by State Space Models (SSM) [9, 56, 61], employs a 2D Selective Scanning Module (2D-SSM) aligned with the

scanning order of wavelet high-frequency components, enabling accurate recovery of high-frequency details. For low-frequency information, we develop the LFEB module for comprehensive recovery.

In summary, our main contributions are as follows:

- We introduce a novel causal framework for LLIE, separating causal and non-causal factors to enhance image quality while preserving semantics.
- We propose a two-level consistency strategy, combining causally-guided metric learning for global consistency and instance-level CLIP loss for local semantic and color consistency.
- We develop the Causal Wavelet Network (CWNNet) with a Hierarchical Feature Restoration Block (HFRB) to model wavelet frequency features, enabling precise high-frequency recovery and robust low-frequency handling.
- Extensive experiments validate CWNNet’s state-of-the-art performance across diverse datasets, demonstrating its robustness and scalability.

## 2. Related Work

**Low-Light Image Enhancement (LLIE):** LLIE methods can be categorized into non-learning-based and learning-based approaches. Traditional techniques, such as histogram equalization (HE) [32] and Retinex theory [31], enhance images by improving contrast or adjusting illumination and reflectance maps. With the advent of deep learning, methods like LLNet [26] and Deep Retinex Decomposition [42] combined Retinex theory with neural networks, leading to significant advancements. Recent approaches, such as FourLLIE [38] and DMFourLLIE [53], leverage frequency domain features to enhance brightness, while Retinexformer [3] integrates transformers to address long-range dependencies. Advanced frameworks like UHD-former [39] and LightDiff [20] tackle ultra-high-definition restoration and unpaired low-light enhancement, respectively, expanding LLIE’s scope to more complex tasks.

**State Space Models (SSM):** State Space Models (SSMs) have emerged as efficient alternatives to CNNs and Transformers for handling long-range dependencies, with linear scalability [8]. Mamba, a structured SSM, has been applied to tasks like super-resolution [16], image classification [36], and restoration [9]. In LLIE, Retinexmamba [2] integrates SSMs into Retinexformer for faster processing, while Wave-Mamba [61] explores UHD low-light enhancement. Recent innovations include LocalMamba [15], which uses localized scanning for detail preservation, and LLE-Mamba [56], which employs bidirectional scanning to balance local and global focus. These advancements highlight the growing role of SSMs in LLIE.

**Causal Inference:** Causal inference focuses on identifying and quantifying causal relationships, with growing applications in computer vision. CIIM [46] removes modality

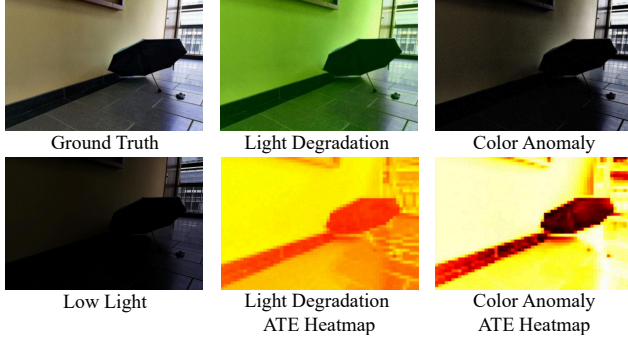


Figure 2. ATE Heatmap Analysis (PSNR). Top row: Ground truth, light degradation, and color anomaly examples. Bottom row: Low light input, ATE heatmaps for light degradation and color anomaly. Brighter regions indicate greater sensitivity to degradations.

bias through causal intervention, while DCIN [21] applies causal reasoning to reduce knowledge bias in image-text matching. MuCR [22] embeds semantic causal relationships for image synthesis. Despite its potential, causal inference remains underexplored in low-level image processing tasks, presenting an opportunity for innovation in LLIE.

### 3. Method

#### 3.1. Causal Inference Analysis for LLIE

##### 3.1.1. Meaningful and Harmless Causal Interventions

For effective causal analysis, interventions must be meaningful and harmless. To achieve these aims, we refer to [12] and design two types of interventions applied to ground truth (normal-light) images for synthetic degradation:

**Light Degradation Intervention:** Instead of simple ablation (which violates the harmless principle), we utilize a physics-based illumination degradation model. Given a normal-light image  $I$ , we generate a light-degraded version  $I_l$  as follows:

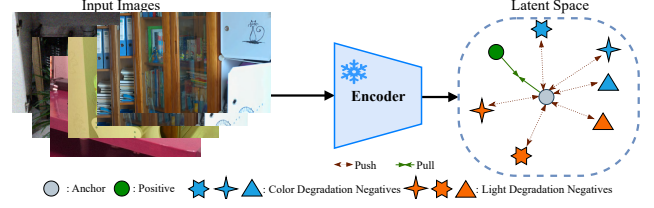
$$I_l = \frac{I}{L} L^\gamma + \varepsilon, \quad (1)$$

where  $L$  is the illumination map generated through LIME [11]. Here,  $\gamma \in [2, 5]$  controls the severity of the degradation, and  $\varepsilon$  is Gaussian noise with mean 0 and variance in  $[0.03, 0.08]$ . This approach offers substantial yet realistic lighting changes while preserving the original semantic content.

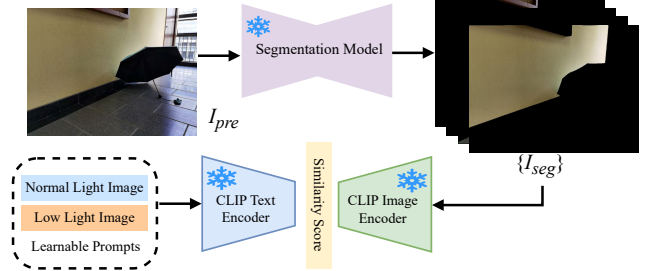
**Color Anomaly Intervention:** To assess the impact of color distortion, we apply the following transformation to the ground truth image:

$$I_c = \Delta H(I) + \Delta S(I) + \sum_{K=R,G,B} \Delta K(I) + \varepsilon, \quad (2)$$

where  $\Delta H \in [-30, 30]$  represents hue shift,  $\Delta S \in$



(a) Causality-Driven Metric Learning Strategy for Causal Inference. The latent space organization includes: Anchor (gray): The network-processed low-light image. Positive (green): The ground truth normal-light reference sharing the same semantic causal factors as the anchor. Color Degradation Negatives (blue): Counterfactual samples generated through color perturbation interventions on normal-light images from different scenes. Light Degradation Negatives (orange): Counterfactual samples generated through brightness perturbation interventions on normal-light images from different scenes.



(b) Instance-Level CLIP Semantic Loss. The enhanced image result  $I_{pre}$  is processed through a pre-trained segmentation network to obtain a series of segmented instance sub-images  $I_{seg}$ . Each sub-image is then iteratively aligned with corresponding textual prompts to ensure semantic consistency.

Figure 3. Global and local causal intervention methods. (a) Eliminate global non-causal interference from illumination and color. (b) Ensure causal semantic consistency.

$[-50, 50]$  represents saturation shift, and  $\Delta K \in [-50, 50]$  denotes RGB channel offsets.

##### 3.1.2. Average Treatment Effect Analysis

To quantitatively assess the impact of our interventions on different image regions, we employ Average Treatment Effect (ATE) analysis [1, 37].

For a feature or region  $p_i$ , the ATE is calculated as:

$$\phi_{\mathcal{F}}[p_i] = \mathbb{E}\{\mathcal{M}_R(I)\} - \mathbb{E}_{t \in \{1:T\}}\{\mathcal{M}_R(I|do(p_i = x_t))\}, \quad (3)$$

where  $I$  is the ground truth image,  $\mathcal{M}_R$  represents the quality metric (PSNR), and  $I|do(p_i = x_t)$  denotes the image with our interventions (light degradation, color anomaly, or noise) applied at intensity  $x_t$ .

To visualize region-specific sensitivity, we compute:

$$\phi_{\mathcal{F}}(I) = \{\phi_{\mathcal{F}}[p_i]\}_{i=1}^N, \quad (4)$$

creating attribution maps that highlight regions most affected by the interventions. As illustrated in Fig. 2, this analysis reveals that degradations affect different semantic regions with varying intensities, emphasizing the need

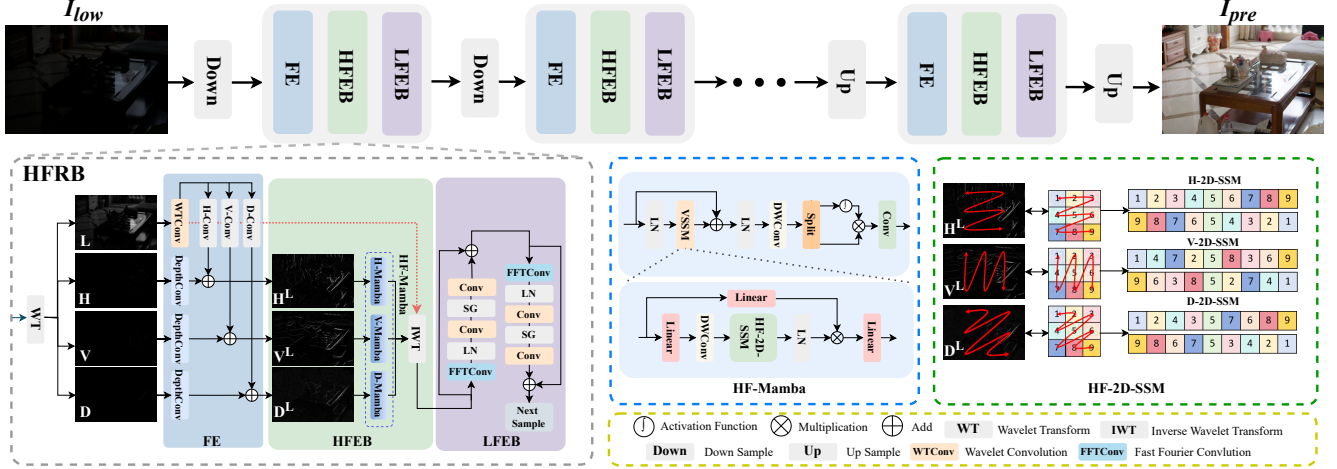


Figure 4. Overall Architecture of CWNet. The low-light image  $I_{low}$  is processed through the sampling and HFRB modules to generate the predicted image  $I_{pre}$ . The expanded structures of each module are shown below.

for both global causal consistency and instance-level causal preservation in our enhancement approach.

### 3.1.3. Causally-Guided Metric Learning

We first address global causal consistency through a causally-guided metric learning approach within a causal inference framework. As shown in Fig. 3(a), we introduce a sample mining strategy for metric learning to disentangle illumination-invariant semantic features (causal factors) from degradation-related factors (non-causal factors). The processed low-light image serves as the anchor, paired with its corresponding normal-light reference as the positive sample. The low-light image, as an extremely degraded instance of the reference in brightness and color (non-causal factors), forms an intrinsic hard positive pair, compelling the metric learning to focus on semantic invariance. Furthermore, we construct counterfactual negative samples by perturbing color and brightness in normal-light images from different scenes. This strategy deliberately excludes other low-light images (preventing confusion between causal and non-causal features), instead forcing the model to discriminate fundamental semantic differences (divergent causal factors) even when non-causal features may resemble the anchor. The metric loss is defined as:

$$\mathcal{L}_{ca}(F_p, \hat{F}, \{F_l\}, \{F_c\}) = \frac{\mathcal{L}_1(F_p, \hat{F})}{\xi(\sum_{l=1}^L \mathcal{L}_1(F_l, \hat{F}) + \sum_{c=1}^C \mathcal{L}_1(F_c, \hat{F}))}, \quad (5)$$

where  $F_p, \hat{F}, \{F_l\}, \{F_c\}$  are the feature representations of the positive sample, anchor, light negative samples, and color abnormal negative samples, respectively. The hyperparameter  $\xi = \frac{1}{L+C}$  normalizes the contributions from light and color negative samples, with  $L$  and  $C$  denoting their total counts.

### 3.1.4. Instance-Level Causal Consistency through CLIP

While our metric learning approach ensures global causal consistency, our ATE analysis identified significant region-specific variability in degradation sensitivity. To uphold local semantic integrity, we introduce an instance-level CLIP-based causal consistency module.

As depicted in Fig. 3(b), we employ HRNet [40], pre-trained on PASCAL-Context [30], to extract semantic instance maps. These maps, along with text prompts, are processed by CLIP encoders to assess semantic consistency:

$$\hat{y} = \frac{1}{K} \sum_{k=1}^K \frac{e^{\cos(\Phi_{\text{image}}(I_{seg}^k), \Phi_{\text{text}}(T_{low}))}}{\sum_{i \in \{low, normal\}} e^{\cos(\Phi_{\text{image}}(I_{seg}^k), \Phi_{\text{text}}(T_i))}}, \quad (6)$$

where  $K$  is the number of sub-instance maps and  $I_{seg}^k$  is a sub-instance map. We optimize using cross-entropy loss:

$$L_{sem} = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})), \quad (7)$$

where  $y$  is the label of the current image, 0 is for low light image and 1 is for normal light image.

## 3.2. Causal Wavelet Network (CWNet)

The architecture of CWNet is illustrated in Fig. 4, comprising upsampling and downsampling layers along with the HFRB. The HFRB consists of FE, HFE, and LFE. Below is a detailed introduction to each component:

### 3.2.1. Feature Extraction (FE)

Given a low-light image  $I_{low} \in \mathbb{R}^{H \times W \times C}$ , we use wavelet transform (WT) to decompose it into four different frequency sub-bands:

$$\{L, H, V, D\} = WT(I_{low}), \quad (8)$$



where  $L, H, V, D \in \mathbb{R}^{\frac{H}{2} \times \frac{W}{2} \times C}$  represent low-frequency component, horizontal, vertical and diagonal high-frequency components, respectively. The image can be reconstructed from these frequency sub-bands using inverse wavelet transform (IWT):

$$I_{low} = IWT\{L, H, V, D\}. \quad (9)$$

In the FE, we design the extraction process to leverage the wavelet domain's characteristics. High-frequency features are obtained using depthwise separable convolutions, while low-frequency features utilize WTConv [6]. The WTConv layer achieves a larger receptive field without additional parameter complexity, essential for capturing low-frequency information. Studies [17, 61] indicate that most information resides in low-frequency components, while high-frequency details are less sensitive in low-light scenarios. To enhance extraction, we employ three convolution kernels aligned with high-frequency directions. The extraction process is formalized as:

$$\begin{aligned} L' &= WTConv(L), \\ H^L &= DepthConv(H) + H-Conv(L'), \\ V^L &= DepthConv(V) + V-Conv(L'), \\ D^L &= DepthConv(D) + D-Conv(L'), \end{aligned} \quad (10)$$

where  $L', H^L, V^L, D^L$  denote the features extracted post-FE, with  $H-Conv$ ,  $V-Conv$ , and  $D-Conv$  extracting horizontal, vertical, and diagonal features, respectively. For detailed structures, please refer to our supplementary materials.

### 3.2.2. High-Frequency Enhancement Block (HFEB)

Recent advancements in state space models (SSM) have significantly improved image enhancement tasks [2, 16, 56, 61].

SSMs transform one-dimensional signals into outputs via latent state representations through linear ordinary differential equations. For a system with input  $x(t)$  and output  $y(t)$ , the model dynamics are described by:

$$h'(t) = Ah(t) + Bx(t), \quad y(t) = Ch(t) + Dx(t), \quad (11)$$

where  $A, B, C$ , and  $D$  are system parameters. The discrete versions, such as Mamba, utilize the zero-order hold (ZOH) discretization, represented as follows:

$$\begin{aligned} h'_t &= \bar{A}h_{t-1} + \bar{B}x_t, \quad y_t = Ch + Dx_t, \\ \bar{A} &= e^{\Delta A}, \quad \bar{B} = (\Delta A)^{-1}(e^{\Delta A} - I) \cdot \Delta B, \end{aligned} \quad (12)$$

where  $\bar{A}$  and  $\bar{B}$  are the discrete counterparts of  $A$  and  $B$ .

Inspired by SSM, we propose the High-Frequency Mamba (HF-Mamba) module, designed for wavelet high-frequency components. The HF-Mamba consists of three parts: D-Mamba, V-Mamba, and H-Mamba. As illustrated

in Fig.4, it applies Layer Normalization (LN), followed by a Visual State Space Module (VSSM) with residual connections, and culminates in a gated feedforward network to enhance channel information flow.

While many existing methods [2, 9, 16, 56, 61] directly adapt the 2D-SSM structure from VMamba [36], we propose distinct processing for wavelet high-frequency components. Horizontal scanning utilizes H-2D-SSM, vertical scanning employs V-2D-SSM, and diagonal scanning incorporates D-2D-SSM. This approach is formalized as:

$$\begin{aligned} \tilde{H}^L &= H-2D-SSM(H^L), \\ \tilde{V}^L &= V-2D-SSM(V^L), \\ \tilde{D}^L &= D-2D-SSM(D^L). \end{aligned} \quad (13)$$

This consistent scanning of high-frequency features extracted by the FE further enhances detail representation.

### 3.2.3. Low-Frequency Enhancement Block (LFEB)

Previous works [5, 45, 60] have employed structure-guided enhancement techniques to optimize image generation, demonstrating that refined high-frequency components can significantly aid in generation and restoration tasks, particularly in LLIE. Building on these insights, we enhance high-frequency components post-High-Frequency Enhancement Block (HFEB) and reconstruct the frequency domain using Inverse Wavelet Transform (IWT):

$$\bar{I}_{low} = IWT\{L', \tilde{H}^L, \tilde{V}^L, \tilde{D}^L\}, \quad (14)$$

where the reconstructed image  $\bar{I}_{low}$  serves as input to the LFEB. Fast Fourier Convolution (FFC) [4] integrates global context within early neural network layers using channel Fast Fourier Transform (FFT) for a broader receptive field.

Based on the above, we propose an LFEB, illustrated in Fig. 4, consisting of two residual blocks tailored for processing low-frequency components, which require large receptive fields. Both blocks employ Fast Fourier Convolution to enhance global features. The first block applies a 5×5 convolution with appropriate padding to expand the receptive field for local spatial context, while the SimpleGate mechanism ensures efficient information flow with minimal loss during activation. Finally, a 1×1 convolution restores the feature dimensions. The second block emphasizes inter-channel correlations and feature enhancement. It employs channel expansion through a 1×1 convolution that quadruples the channel dimensions. The SimpleGate mechanism selectively preserves important features, and another 1×1 convolution compresses the features back to the original channel size. Following the LFEB and subsequent modules, low-frequency components are refined under the guidance of high-frequency components, resulting in the predicted brightened output  $I_{pre}$ .

Category	Methods	LOL-v1			LOL-v2-Real			LOL-v2-Syn			LSRW-Huawei			#Param (M)	#Flops (G)
		PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$		
Traditional	NPE [41]	16.97	0.5928	0.2456	17.33	0.4642	0.2359	16.60	0.7781	0.1079	17.08	0.3905	0.2303	-	-
	LIME [11]	16.76	0.4440	0.2060	15.24	0.4190	0.2203	16.88	0.7578	0.1041	17.00	0.3816	0.2069	-	-
	SRIE [7]	11.80	0.5000	0.1862	14.45	0.5240	0.2160	14.50	0.6640	0.1484	13.42	0.4282	0.2166	-	-
CNN-based	Kind [57]	20.87	0.7995	0.2071	17.54	0.6695	0.3753	22.62	0.9041	0.0515	16.58	0.5690	0.2259	8.02	34.99
	MIRNet [52]	24.14	0.8305	0.2502	22.11	0.7942	0.1448	22.52	0.8997	0.0568	19.98	0.6085	0.2154	31.79	785.1
	Kind++ [58]	17.97	0.8042	0.1756	19.08	0.8176	0.1803	21.17	0.8814	0.0678	15.43	0.5695	0.2366	8.27	2970.5
Frequency-based	FourLLIE [38]	20.99	0.8071	0.0952	23.45	0.8450	0.0613	24.65	0.9192	0.0389	21.11	0.6256	0.1825	0.12	4.07
	UHDFour [19]	22.89	0.8147	0.0934	27.27	0.8579	0.0617	23.64	0.8998	0.0341	19.39	0.6006	0.2466	17.54	4.78
	DMFourLLIE [53]	22.98	0.8273	0.0792	26.40	0.8765	0.0526	25.74	0.9308	<u>0.0251</u>	21.09	<u>0.6328</u>	0.1804	0.41	1.70
Transformer-based	SNR-Aware [44]	<b>23.93</b>	<u>0.8460</u>	0.0813	21.48	0.8478	0.0740	24.13	0.9269	0.0318	20.67	0.5911	0.1923	39.12	26.35
	Retinexformer [3]	22.71	0.8177	0.0922	24.55	0.8434	0.0627	25.67	0.9295	0.0273	<u>21.23</u>	0.6309	0.1699	1.61	15.57
Mamba-based	Wave-Mamba [61]	22.76	0.8419	<u>0.0791</u>	<b>27.87</b>	<u>0.8935</u>	<u>0.0451</u>	24.69	0.9271	0.0584	21.19	0.6391	0.1818	1.26	7.22
	RetinexMamba [2]	23.15	0.8210	0.0876	27.31	0.8667	0.0551	<b>25.89</b>	<u>0.9346</u>	0.0389	20.88	0.6298	<u>0.1689</u>	3.59	34.76
	CWNet	<u>23.60</u>	<b>0.8496</b>	<b>0.0648</b>	<u>27.39</u>	<b>0.9005</b>	<b>0.0383</b>	<u>25.50</u>	<b>0.9362</b>	<b>0.0195</b>	<b>21.50</b>	<b>0.6397</b>	<b>0.1562</b>	1.23	11.3

Table 1. Quantitative comparison on LOL-v1 [51], LOL-v2-Real [51], LOL-v2-Syn [51], and LSRW-Huawei [13]. The **best** and second-best results are shown in bold and underlined respectively. Please note that we did not use the GT-Mean strategy.

### 3.3. Loss Function

Our total loss consists of five parts:

$$\mathcal{L}_{total} = \lambda_1 \mathcal{L}_2 + \lambda_2 \mathcal{L}_{ssim} + \lambda_3 \mathcal{L}_{per} + \lambda_4 \mathcal{L}_{ca} + \lambda_5 \mathcal{L}_{sem}, \quad (15)$$

where  $\lambda$  denotes the loss weights, we set  $\lambda_1, \lambda_2, \lambda_3, \lambda_4, \lambda_5 = [1.0, 0.3, 0.2, 0.01, 0.01]$ .  $\mathcal{L}_2$  represents the  $l_2$  loss.  $\mathcal{L}_{ssim}$  is the structure similarity loss.  $\mathcal{L}_{per}$  is the perceptual loss, which constrains the features extracted from VGG to obtain better visual results.  $\mathcal{L}_2$ ,  $\mathcal{L}_{ssim}$  and  $\mathcal{L}_{per}$  constrain the output  $I_{pre}$  and ground truth  $I_{gt}$  end-to-end.

## 4. Experiments

### 4.1. Datasets and Experimental Setting

CWNet is trained and evaluated on four LLIE datasets: LOL-v1 [51], LOL-v2-Real [51], LOL-v2-Synthesis [51], and LSRW-Huawei [13]. LOL-v1 contains 485 training and 15 testing pairs of real-world low-light/normal-light images. LOL-v2-Real provides 689 training and 100 testing pairs with more diverse real-world scenarios. For LOL-v2-Real evaluation, we use the model trained on LOL-v1 to demonstrate cross-dataset generalization. LOL-v2-Synthesis includes 900 training and 100 testing synthesized pairs. LSRW-Huawei comprises 2450 training and 30 testing pairs captured with different devices.

CWNet is implemented in PyTorch and trained end-to-end to jointly optimize all network parameters. The model employs a U-Net-like architecture with feature channels 16 and asymmetric block configurations of  $[1, 3, 4, 3, 1]$  and  $[1, 2, 2, 2, 1]$  for low and high-frequency branches respectively. During training, input images are randomly cropped to  $256 \times 256$  patches and augmented with random horizontal/vertical flips and rotations. We use the Adam optimizer with  $\beta_1 = 0.9$ ,  $\beta_2 = 0.99$ , and an initial learning rate of  $4.0 \times 10^{-4}$ . The total training is conducted for  $3.0 \times 10^5$  iterations with a batch size of 8.

### 4.2. Comparison with Current Methods

In this paper, our CWNet is compared to current state-of-the-art LLIE methods, including traditional approaches LIME [11], NPE [41] and SRIE [7], deep learning-based methods Kind [57], Kind++ [58], MIRNet [52], SGM [51], SNR-Aware [44], FourLLIE [38], FECNet [14], UHD-Four [19], Retinexformer [3], DMFourLLIE [53], Retinex-Mamba [2] and Wave-Mamba [61].

**Quantitative Results on LOL-v2-Real, LOL-v2-Synthesis, and LSRW-Huawei Datasets.** We comprehensively evaluate CWNet using PSNR, SSIM, and LPIPS metrics, where higher PSNR/SSIM and lower LPIPS indicate better image quality. As shown in Tab. 1, CWNet achieves superior performance across all benchmarks: PSNR of 23.60 dB on LOL-v1, 27.39 dB on LOL-v2-Real, 25.50 dB on LOL-v2-Synthesis, and 21.50 dB on LSRW-Huawei. Particularly noteworthy is the exceptional cross-dataset generalization from LOL-v1 training to LOL-v2-Real testing, achieving the best SSIM of 0.9005 and lowest LPIPS of 0.0383. Importantly, CWNet achieves this superior performance while maintaining computational efficiency with only 1.23M parameters and 11.3G FLOPs, significantly outperforming parameter-heavy methods like MIRNet [52] (31.79M) and SNR-Aware [44] (39.12M), demonstrating effective balance between performance and efficiency.

**Visualization Comparison on LOL-v2-Real and LSRW-Huawei Datasets.** We compare CWNet with state-of-the-art methods, including FECNet [14], FourLLIE [38], Wave-Mamba, Retinexformer [61], SKF-SNR [43], UHD-Former [39], UHDFour [19], and DMFourLLIE [53]. As shown in Fig. 5, while other methods improve brightness, they often fail to maintain color and semantic consistency or control noise. For example, FECNet, FourLLIE, and Wave-Mamba show color deviations and noise, while Retinexformer and SKF-SNR lack sufficient brightening. UHD-Former and UHDFour perform better but still exhibit noise

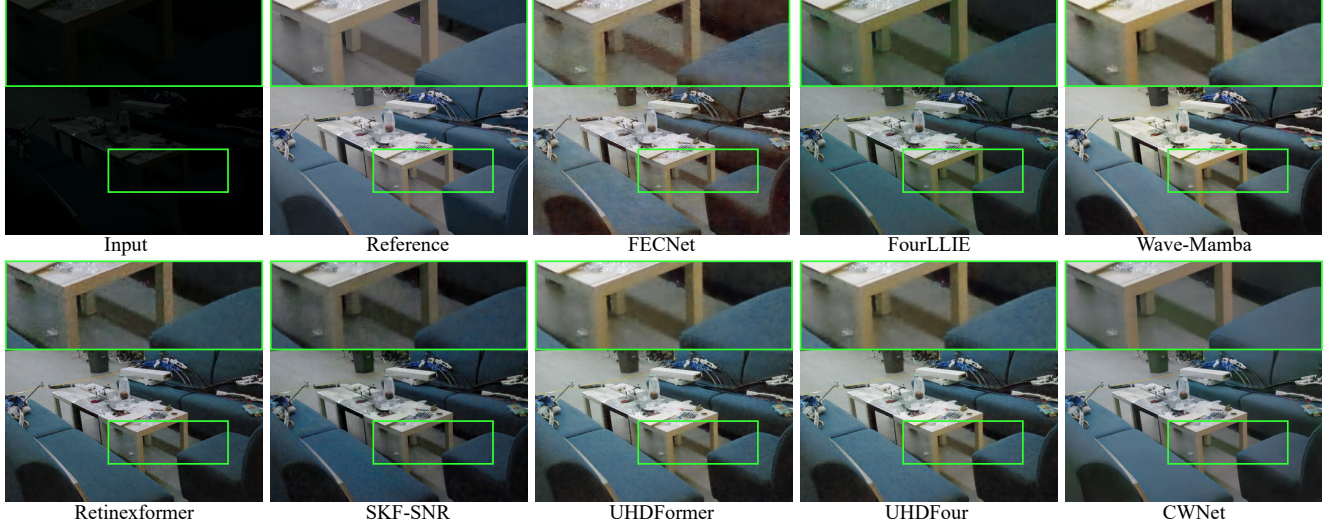


Figure 5. Visual comparison on LOL-v2-Real dataset.



Figure 6. Visual comparison on LSRW-Huawei dataset.

artifacts and lack smoothness. In contrast, CWNet produces clearer, more natural, and smoother results, ensuring semantic and color consistency, demonstrating its superior frequency domain modeling and structural fidelity.

**Visualization on LSRW-Huawei Dataset.** As shown in Fig. 6, CWNet outperforms others in clarity and detail preservation. Methods like FECNet, SKF-SNR, and UHDFormer show exposure and brightening deficiencies, while CWNet achieves balanced brightness and excellent detail retention. For additional comparisons, refer to supplementary materials.

### 4.3. Ablation Studies

We conduct comprehensive ablation studies on the LSRW-Huawei dataset to evaluate the effectiveness of CWNet’s design components. Tab. 2 presents the results across two experimental settings:

**Component Removal Analysis.** We systematically remove key components to assess their individual contributions. Removing the causal inference mechanism results in

a significant performance drop (PSNR: 20.87 dB vs. 21.53 dB). Similarly, ablating the feature extraction (FE) modules, HFEB, and LFEB leads to substantial performance degradation, with LFEB removal causing the most severe impact (PSNR drops to 20.41 dB), confirming the critical role of low-frequency processing in our dual-branch architecture.

**Component Replacement Analysis.** We validate the effectiveness of our proposed modules by replacing them with conventional alternatives. Substituting WTConv and FFTConv with standard convolutions reduces performance (PSNR: 21.42 dB and 21.36 dB respectively), highlighting the benefits of frequency-domain processing. Replacing HF-Mamba with standard VMamba’s 2D-SSM structure [36] also degrades performance (PSNR: 21.20 dB), demonstrating the superiority of our high-frequency Mamba design. Additionally, replacing semantic maps with global features shows performance reduction (PSNR: 21.48 dB), confirming the value of semantic guidance.

**Loss Weights Analysis.** We conducted comprehen-



Ablation Settings	PSNR↑	SSIM↑	LPIPS↓
Component Removal			
w/o Casual Inference	20.87	0.6375	0.1781
w/o FE	20.98	0.6387	0.1804
w/o HFEB	20.58	0.6317	0.1903
w/o LFEB	20.41	0.6302	0.1985
Component Replacement			
WTConv → Conv	21.42	0.6415	0.1690
FFTCConv → Conv	21.36	0.6396	0.1721
HF-Mamba → VMamba (2D-SSM)	21.20	0.6394	0.1735
Segmentic Maps → Global	21.48	0.6417	0.1652
CWNet	<b>21.53</b>	<b>0.6423</b>	<b>0.1631</b>

Table 2. Ablation experiment study of CWNet. By designing ablation experiments on component removal and replacement, the effectiveness of each component and composition of CWNet is fully verified.

Loss	$\mathcal{L}_1$	$\mathcal{L}_{ssim}$	$\mathcal{L}_{per}$	$\mathcal{L}_{ca}$	$\mathcal{L}_{sem}$	PSNR↑	SSIM↑	LPIPS↓
	1.0	0.3	0.2	0.01	0.01	<b>21.53</b>	<u>0.6423</u>	0.1631
A	1.0	0.4	0.2	0.01	0.01	21.39	<b>0.6433</b>	<b>0.1597</b>
B	1.0	0.3	0.3	0.01	0.01	21.34	0.6407	0.1601
C	1.0	0.3	0.2	0.05	0.01	<u>21.43</u>	0.6386	0.1614
D	1.0	0.3	0.2	0.001	0.01	21.17	0.6408	0.1651
E	1.0	0.3	0.2	0.01	0.05	20.89	0.6382	0.1701
F	1.0	0.3	0.2	0.01	0.001	20.94	0.6371	0.1652

Table 3. Ablation Study on Weight Configuration in Loss Functions. The loss function is defined as  $\mathcal{L}_{total} = \lambda_1 \mathcal{L}_2 + \lambda_2 \mathcal{L}_{ssim} + \lambda_3 \mathcal{L}_{per} + \lambda_4 \mathcal{L}_{ca} + \lambda_5 \mathcal{L}_{sem}$ . The **best** and second-best results are shown in bold and underlined respectively.

sive ablation experiments by systematically varying each weight parameter to analyze their individual impact on performance. As shown in Tab.3, different weight configurations lead to varying performance across metrics. The baseline configuration achieves the best PSNR (21.53) and the second-best SSIM (0.6423), demonstrating strong perceptual quality.

These ablation results collectively validate that each proposed component contributes meaningfully to CWNet’s superior performance. For additional ablation studies and downstream applications, please refer to our supplementary materials.

#### 4.4. Limitations

As shown in Fig.7, when facing its own degradation, such as blurring or haze, CWNet maintains the color and lighting intensity of the image as much as possible compared to other methods, yet the restoration quality is subpar. This opens up new avenues for us to explore how to ensure more effective recovery of low-light images that experience multiple degradation conditions simultaneously.

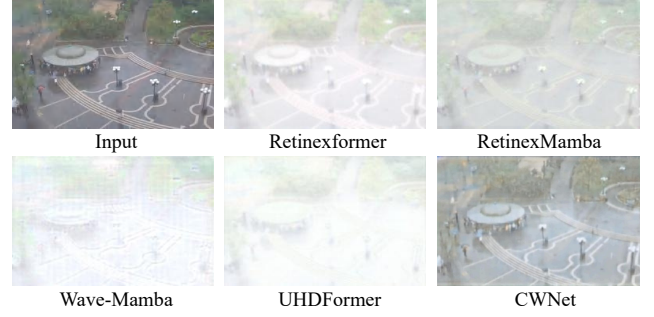


Figure 7. Failure cases in multiple degradation scenarios.

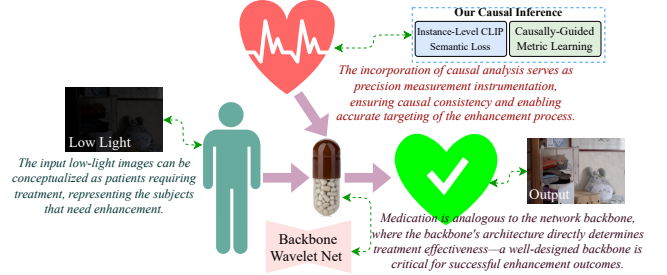


Figure 8. The connection between wavelet and causality.

## 5. Discussion and Conclusion

**Causality-Wavelet Connection.** The interpretability of CWNet and causality-wavelet connection: Wavelet structure and causal analysis are organically integrated (Fig.8) with causal treatment providing model interpretability. We clarify: **a) Causal Perspective:** Our focus centers on causal analysis (analogous to measurement instruments in causal treatment), treating scene information as causal factors in LLIE, ensuring causal consistency during enhancement. **b) Wavelet-based Backbone:** To achieve causal factor consistency (analogous to medication in treatment), our design leverages low-frequency enhancement for color and brightness consistency while high-frequency components incorporate Mamba consistency scanning to enhance detail modeling and promote structural consistency.

**Conclusion.** In this paper, we proposed the Causal Wavelet Network (CWNet), a novel architecture that integrates causal inference with wavelet transform to address low-light image enhancement. By leveraging causal analysis, we effectively separated causal factors from non-causal interference, ensuring both global and local semantic consistency. The HFRB further refined feature extraction by modeling wavelet frequency domain characteristics. Extensive experiments demonstrate that CWNet achieves superior performance over state-of-the-art methods, highlighting the effectiveness of causal reasoning in enhancing image quality. This work underscores the potential of causal inference as a powerful tool for advancing low-light image enhancement.



## Acknowledgements

This work was supported by Jilin Province Industrial Key Core Technology Tackling Project (20230201085GX).

## References

- [1] Joshua Angrist and Guido Imbens. Identification and estimation of local average treatment effects, 1995. [3](#)
- [2] Jiesong Bai, Yuhao Yin, and Qiyuan He. Retinexmamba: Retinex-based mamba for low-light image enhancement. *arXiv preprint arXiv:2405.03349*, 2024. [1](#), [2](#), [5](#), [6](#)
- [3] Yuanhao Cai, Hao Bian, Jing Lin, Haoqian Wang, Radu Timofte, and Yulun Zhang. Retinexformer: One-stage retinex-based transformer for low-light image enhancement. In *CVPR*, pages 12504–12513, 2023. [1](#), [2](#), [6](#)
- [4] Lu Chi, Borui Jiang, and Yadong Mu. Fast fourier convolution. *Advances in Neural Information Processing Systems*, 33:4479–4488, 2020. [5](#)
- [5] Qiaole Dong, Chenjie Cao, and Yanwei Fu. Incremental transformer structure enhanced image inpainting with masking positional encoding. In *ICCV*, pages 11358–11368, 2022. [5](#)
- [6] Shahaf E Finder, Roy Amoyal, Eran Treister, and Oren Freifeld. Wavelet convolutions for large receptive fields. In *ECCV*, pages 363–380. Springer, 2025. [5](#), [1](#)
- [7] Xueyang Fu, Delu Zeng, Yue Huang, Xiao-Ping Zhang, and Xinghao Ding. A weighted variational model for simultaneous reflectance and illumination estimation. In *CVPR*, pages 2782–2790, 2016. [6](#)
- [8] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. *arXiv preprint arXiv:2111.00396*, 2021. [2](#)
- [9] Hang Guo, Jinmin Li, Tao Dai, Zhihao Ouyang, Xudong Ren, and Shu-Tao Xia. Mambair: A simple baseline for image restoration with state-space model. In *ECCV*, pages 222–241. Springer, 2025. [2](#), [5](#)
- [10] Xiaojie Guo and Qiming Hu. Low-light image enhancement via breaking down the darkness. *IJCV*, 131(1):48–66, 2023. [1](#)
- [11] Xiaojie Guo, Yu Li, and Haibin Ling. Lime: Low-light image enhancement via illumination map estimation. *IEEE TIP*, 26(2):982–993, 2016. [3](#), [6](#), [2](#)
- [12] Yu Guo, Yuan Gao, Yuxu Lu, Huilin Zhu, Ryan Wen Liu, and Shengfeng He. Onerestore: A universal restoration framework for composite degradation. In *ECCV*, pages 255–272. Springer, 2024. [3](#)
- [13] Jiang Hai, Zhu Xuan, Ren Yang, Yutong Hao, Fengzhu Zou, Fang Lin, and Songchen Han. R2rnet: Low-light image enhancement via real-low to real-normal network. *Journal of Visual Communication and Image Representation*, 90: 103712, 2023. [6](#)
- [14] Jie Huang, Yajing Liu, Feng Zhao, Keyu Yan, Jinghao Zhang, Yukun Huang, Man Zhou, and Zhiwei Xiong. Deep fourier-based exposure correction network with spatial-frequency interaction. In *ECCV*, pages 163–180. Springer, 2022. [1](#), [6](#)
- [15] Tao Huang, Xiaohuan Pei, Shan You, Fei Wang, Chen Qian, and Chang Xu. Localmamba: Visual state space model with windowed selective scan. *arXiv preprint arXiv:2403.09338*, 2024. [2](#)
- [16] Yongsong Huang, Tomo Miyazaki, Xiaofeng Liu, and Shinichiro Omachi. Irsrmamba: Infrared image super-resolution via mamba-based wavelet transform feature modulation model. *arXiv preprint arXiv:2405.09873*, 2024. [2](#), [5](#)
- [17] Hai Jiang, Ao Luo, Haoqiang Fan, Songchen Han, and Shuaicheng Liu. Low-light image enhancement with wavelet-based diffusion models. *ACM Transactions on Graphics (TOG)*, 42(6):1–14, 2023. [1](#), [2](#), [5](#)
- [18] Chulwoo Lee, Chul Lee, and Chang-Su Kim. Contrast enhancement based on layered difference representation. In *ICIP*, pages 965–968. IEEE, 2012. [2](#)
- [19] Chongyi Li, Chun-Le Guo, Man Zhou, Zhixin Liang, Shangchen Zhou, Ruicheng Feng, and Chen Change Loy. Embedding fourier for ultra-high-definition low-light image enhancement. *arXiv preprint arXiv:2302.11831*, 2023. [1](#), [2](#), [6](#)
- [20] Jinlong Li, Baolu Li, Zhengzhong Tu, Xinyu Liu, Qing Guo, Felix Juefei-Xu, Runsheng Xu, and Hongkai Yu. Light the night: A multi-condition diffusion framework for unpaired low-light enhancement in autonomous driving. In *ICCV*, pages 15205–15215, 2024. [2](#)
- [21] Wenhui Li, Xinqi Su, Dan Song, Lanjun Wang, Kun Zhang, and An-An Liu. Towards deconfounded image-text matching with causal inference. In *ACM MM*, pages 6264–6273, 2023. [3](#)
- [22] Zhiyuan Li, Heng Wang, Dongnan Liu, Chaoyi Zhang, Ao Ma, Jieting Long, and Weidong Cai. Multimodal causal reasoning benchmark: Challenging vision large language models to infer causal links between siamese images. *arXiv preprint arXiv:2408.08105*, 2024. [3](#)
- [23] Dong Liang, Ling Li, Mingqiang Wei, Shuo Yang, Liyan Zhang, Wenhan Yang, Yun Du, and Huiyu Zhou. Semantically contrastive learning for low-light image enhancement. In *AAAI*, pages 1555–1563, 2022. [1](#)
- [24] Yudong Liang, Bin Wang, Wenqi Ren, Jiaying Liu, Wenjian Wang, and Wangmeng Zuo. Learning hierarchical dynamics with spatial adjacency for image enhancement. In *ACM MM*, pages 2767–2776, 2022. [1](#)
- [25] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, pages 740–755. Springer, 2014. [3](#)
- [26] Kin Gwn Lore, Adedotun Akintayo, and Soumik Sarkar. Ll-net: A deep autoencoder approach to natural low-light image enhancement. *Pattern Recognition*, pages 650–662, 2017. [2](#)
- [27] Kede Ma, Kai Zeng, and Zhou Wang. Perceptual quality assessment for multi-exposure image fusion. *IEEE TIP*, 24(11):3345–3356, 2015. [2](#)
- [28] Anish Mittal, Rajiv Soundararajan, and Alan C. Bovik. Making a “completely blind” image quality analyzer. *IEEE Signal Processing Letters*, 20(3):209–212, 2013. [2](#)

- [29] Nathan Moroney. Local color correction using non-linear masking. In *Color and Imaging conference*, pages 108–111. Society of Imaging Science and Technology, 2000. 1
- [30] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, pages 891–898, 2014. 4
- [31] Michael K Ng and Wei Wang. A total variation model for retinex. *SIAM Journal on Imaging Sciences*, 4(1):345–365, 2011. 1, 2
- [32] Etta D Pisano, Shuquan Zong, Bradley M Hemminger, Marla DeLuca, R Eugene Johnston, Keith Muller, M Patricia Braeuning, and Stephen M Pizer. Contrast limited adaptive histogram equalization image processing to improve the detection of simulated spiculations in dense mammograms. *Journal of Digital imaging*, 11:193–200, 1998. 1, 2
- [33] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PMLR, 2021. 1
- [34] Yoav Y Schechner and Nir Karpel. Recovery of underwater visibility and structure by polarization analysis. *IEEE Journal of oceanic engineering*, 30(3):570–587, 2005. 2
- [35] Bernhard Schölkopf. Causality for machine learning. In *Probabilistic and causal inference: The works of Judea Pearl*, pages 765–804. 2022. 2
- [36] Yuan Shi, Bin Xia, Xiaoyu Jin, Xing Wang, Tianyu Zhao, Xin Xia, Xuefeng Xiao, and Wenming Yang. Vmambair: Visual state space model for image restoration. *arXiv preprint arXiv:2403.11423*, 2024. 2, 5, 7
- [37] Tyler J VanderWeele and Miguel A Hernan. Causal inference under multiple versions of treatment. *Journal of causal inference*, pages 1–20, 2013. 3
- [38] Chenxi Wang, Hongjun Wu, and Zhi Jin. Fourllie: Boosting low-light image enhancement by fourier frequency information. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 7459–7469, 2023. 1, 2, 6
- [39] Cong Wang, Jinshan Pan, Wei Wang, Gang Fu, Siyuan Liang, Mengzhu Wang, Xiao-Ming Wu, and Jun Liu. Correlation matching transformation transformers for uhd image restoration. In *AAAI*, pages 5336–5344, 2024. 1, 2, 6
- [40] 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*, 43(10):3349–3364, 2020. 4
- [41] Shuhang Wang, Jin Zheng, Hai-Miao Hu, and Bo Li. Naturalness preserved enhancement algorithm for non-uniform illumination images. *IEEE TIP*, 22(9):3538–3548, 2013. 6, 2
- [42] C Wei, W Wang, W Yang, and J Liu. Deep retinex decomposition for low-light enhancement. *arXiv preprint arXiv:1808.04560*, 2018. 2
- [43] Yuhui Wu, Chen Pan, Guoqing Wang, Yang Yang, Jiwei Wei, Chongyi Li, and Heng Tao Shen. Learning semantic-aware knowledge guidance for low-light image enhancement. In *ICCV*, pages 1662–1671, 2023. 1, 6
- [44] Xiaogang Xu, Ruixing Wang, Chi-Wing Fu, and Jiaya Jia. Snr-aware low-light image enhancement. In *CVPR*, pages 17714–17724, 2022. 6
- [45] Xiaogang Xu, Ruixing Wang, and Jiangbo Lu. Low-light image enhancement via structure modeling and guidance. In *ICCV*, pages 9893–9903, 2023. 5
- [46] Jiexi Yan, Cheng Deng, Heng Huang, and Wei Liu. Causality-invariant interactive mining for cross-modal similarity learning. *IEEE TPAMI*, 2024. 2
- [47] Qirui Yang, Peng-Tao Jiang, Hao Zhang, Jinwei Chen, Bo Li, Huanjing Yue, and Jingyu Yang. Learning adaptive lighting via channel-aware guidance. *arXiv preprint arXiv:2412.01493*, 2024. 1
- [48] Qirui Yang, Qihua Cheng, Huanjing Yue, Le Zhang, Yihao Liu, and Jingyu Yang. Learning to see low-light images via feature domain adaptation. *IEEE Transactions on Image Processing*, 2025. 1
- [49] Shuzhou Yang, Moxuan Ding, Yanmin Wu, Zihan Li, and Jian Zhang. Implicit neural representation for cooperative low-light image enhancement. In *CVPR*, pages 12918–12927, 2023. 1
- [50] Wenhan Yang, Ye Yuan, Wenqi Ren, Jiaying Liu, Walter J Scheirer, Zhangyang Wang, Taiheng Zhang, Qiaoyong Zhong, Di Xie, Shiliang Pu, et al. Advancing image understanding in poor visibility environments: A collective benchmark study. *IEEE TIP*, 29:5737–5752, 2020. 2, 3
- [51] Wenhan Yang, Wenjing Wang, Haofeng Huang, Shiqi Wang, and Jiaying Liu. Sparse gradient regularized deep retinex network for robust low-light image enhancement. *IEEE TIP*, 30:2072–2086, 2021. 6, 1
- [52] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In *ECCV*, pages 492–511. Springer, 2020. 6, 1
- [53] Tongshun Zhang, Pingping Liu, Ming Zhao, and Haotian Lv. Dmfourllie: Dual-stage and multi-branch fourier network for low-light image enhancement. In *ACM MM*, pages 7434–7443, 2024. 1, 2, 6
- [54] Tongshun Zhang, Pingping Liu, Mengen Cai, Xiaoyi Wang, and Qiuzhan Zhou. Cross-modal guided and refinement-enhanced retinex network for robust low-light image enhancement. *Information Fusion*, page 103380, 2025. 1
- [55] Xiaozhe Zhang, Fengying Xie, Haidong Ding, Linpeng Pan, and Zhenwei Shi. Adapt clip as aggregation instructor for image dehazing. *arXiv preprint arXiv:2408.12317*, 2024. 1
- [56] Xuanqi Zhang, Haijin Zeng, Jinwang Pan, Qiangqiang Shen, and Yongyong Chen. Llemamba: Low-light enhancement via relighting-guided mamba with deep unfolding network. *arXiv preprint arXiv:2406.01028*, 2024. 1, 2, 5
- [57] Yonghua Zhang, Jiawan Zhang, and Xiaojie Guo. Kindling the darkness: A practical low-light image enhancer. In *ACM MM*, pages 1632–1640, 2019. 6, 1
- [58] Yonghua Zhang, Xiaojie Guo, Jiayi Ma, Wei Liu, and Jiawan Zhang. Beyond brightening low-light images. *IJCV*, 129:1013–1037, 2021. 6

- [59] Zhao Zhang, Huan Zheng, Richang Hong, Mingliang Xu, Shuicheng Yan, and Meng Wang. Deep color consistent network for low-light image enhancement. In *ICCV*, pages 1899–1908, 2022. [1](#)
- [60] Minfeng Zhu, Pingbo Pan, Wei Chen, and Yi Yang. Eemefn: Low-light image enhancement via edge-enhanced multi-exposure fusion network. In *AAAI*, pages 13106–13113, 2020. [5](#)
- [61] Wenbin Zou, Hongxia Gao, Weipeng Yang, and Tongtong Liu. Wave-mamba: Wavelet state space model for ultra-high-definition low-light image enhancement. In *ACM MM*, pages 1534–1543, 2024. [1](#), [2](#), [5](#), [6](#)

# CWNet: Causal Wavelet Network for Low-Light Image Enhancement

## Supplementary Material

### A. Network Structure

#### A.1. Detailed Structure of Feature Extraction (FE)

In the FE module, we utilize WTConv [6], H-Conv, V-Conv, and D-Conv. Each component is described in detail below:

WTConv [6] employs a convolution kernel of size 5 and is configured with 3 levels of wavelet downsampling to progressively reduce spatial resolution while preserving frequency information.

The **H-Conv**, **V-Conv**, and **D-Conv** modules utilize specifically designed convolution kernels to capture directional features in the input data. These kernels are as follows:

- Horizontal Convolution (H-Conv):

$$\text{Horizontal Kernel: } \begin{bmatrix} 1 & 0 & -1 \\ 1 & 0 & -1 \\ 1 & 0 & -1 \end{bmatrix}$$

This kernel is designed to emphasize horizontal edges in the input data.

- Vertical Convolution (V-Conv):

$$\text{Vertical Kernel: } \begin{bmatrix} 1 & 1 & 1 \\ 0 & 0 & 0 \\ -1 & -1 & -1 \end{bmatrix}$$

This kernel is designed to capture vertical edge features.

- Diagonal Convolution (D-Conv):

$$\text{Diagonal Kernel: } \begin{bmatrix} 0 & 1 & 0 \\ 1 & -4 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

This kernel is designed to detect diagonal structures.

### B. Experiments

#### B.1. Visualization Comparison on LOL-v2-Synthesis Dataset.

Fig.10 presents a comparison on the LOL-v2-Synthesis dataset between our CWNet and current state-of-the-art methods, including FECNet [14], FourLLIE [38], Retinexformer [3], UHDFormer [19], UHDFormer [39], and Wave-Mamba [61].

In the first row, other methods exhibit stripe noise in the sky region. Specifically, FECNet shows severe hue distortion, while other methods lack image sharpness. In the second row, similar artifacts are observed, with stripe-like noise appearing in the background areas.

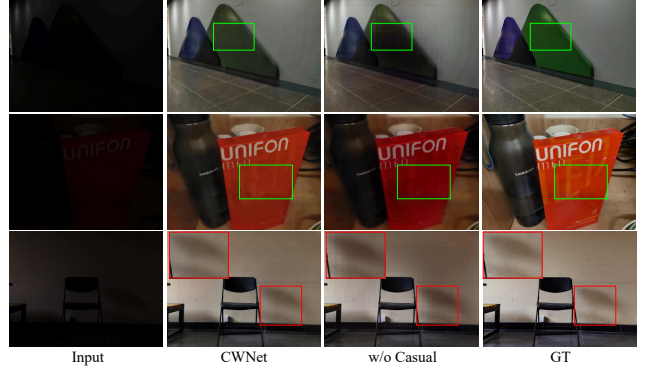


Figure 9. Visual comparison of results with and without the causal inference. Focus on the highlighted regions to observe differences in color and brightness, demonstrating the module’s effectiveness.

Methods	DICM	LIME	MEF	NPE	VV	AVG
Kind [57]	<b>3.61</b>	4.77	4.82	4.18	3.84	4.24
MIRNet [52]	4.04	6.45	5.50	5.24	4.74	5.19
SGM [51]	4.73	5.45	5.75	5.21	4.88	5.21
FECNet [14]	4.14	6.04	4.71	4.50	3.75	4.55
HDMNet [24]	4.77	6.40	5.99	5.11	4.46	5.35
Bread [10]	4.18	4.72	5.37	4.16	<b>3.30</b>	4.35
Retinexformer [3]	4.01	<b>3.44</b>	<u>3.73</u>	<u>3.89</u>	<u>3.71</u>	<u>3.76</u>
UHDFormer [39]	4.42	4.35	4.74	4.40	4.28	4.44
Wave-Mamba [61]	4.56	4.45	4.76	4.54	4.71	4.60
CWNet	<u>3.92</u>	<u>3.58</u>	<b>3.66</b>	<b>3.61</b>	3.74	<b>3.70</b>

Table 4. NIQE scores on DICM, LIME, MEF, NPE, and VV datasets. Lower NIQE scores indicate better perceptual quality. The **best** and second-best results in each column are shown in bold and underlined respectively. "AVG" denotes the average NIQE scores across the five datasets.

In the third and fifth rows, other methods exhibit color distortion when compared to the reference ground truth. In contrast, our results appear more natural and exhibit better visual quality. This demonstrates the effectiveness of our causal inference component in mitigating color distortion and preserving semantic structures by disentangling causal relationships in the feature space.

In the fourth and last rows, upon magnification, our results exhibit the most consistent and visually pleasing brightening effect. The comparison on the LOL-v2-Synthesis dataset further demonstrates that our CWNet achieves natural brightness restoration and excels in preserving semantic and color consistency.



## B.2. Quantitative and Qualitative Comparison on the DICM, LIME, MEF, NPE, and VV Datasets

**Quantitative Comparison.** We evaluated CWNet on five independent datasets: DICM [18] (69 images), LIME [11] (10 images), MEF [27] (17 images), NPE [41] (85 images), and VV (24 images). The evaluation was conducted using the no-reference image quality assessment metric NIQE [28], with lower scores indicating better perceptual quality and naturalness. Tab. 4 presents the NIQE results. As shown, CWNet outperforms several state-of-the-art (SOTA) methods, including SKF-SNR, UHDFormer, and Wave-Mamba. To ensure a fair comparison, all methods were pretrained on the LSRW-Huawei dataset.

**Qualitative Comparison.** We provide qualitative comparisons on the DICM, LIME, MEF, NPE, and VV datasets to visually demonstrate the effectiveness of CWNet.

1) Fig. 11 shows the qualitative comparison on the DICM dataset. In the first row, SKF-SNR produces distorted and overexposed results, which negatively impact the visual quality. UHDFormer generates relatively natural results but suffers from overexposure, particularly in the highlighted regions, where the flower colors are overly brightened. Wave-Mamba exhibits blurred edge details and insufficient exposure control. In contrast, CWNet produces clearer and more natural enhancement results. In the second row, CWNet produces sharper and more natural results, while SKF-SNR suffers from underexposure and noise artifacts. UHDFormer and Wave-Mamba produce relatively blurry results. In the third row, SKF-SNR introduces unavoidable noise, and both UHDFormer and Wave-Mamba generate unnatural sky colors. CWNet, however, produces results with consistent and natural sky colors, demonstrating its robustness in handling color distortions.

2) Fig. 12 presents the qualitative comparison on the LIME dataset. SKF-SNR fails to achieve sufficient brightness enhancement, while UHDFormer and Wave-Mamba produce reasonable brightness but suffer from blurriness and lack of detail. In contrast, CWNet generates sharper textures and richer details, further validating the robustness of our high- and low-frequency modeling in capturing fine-grained details and global structures.

3) Fig. 13 illustrates the qualitative comparison on the MEF dataset. SKF-SNR produces unnatural flame colors and insufficient brightness enhancement. UHDFormer and Wave-Mamba exhibit varying degrees of blurriness, while CWNet achieves the clearest and most visually pleasing enhancement results.

4) Fig. 14 shows the qualitative comparison on the NPE dataset. Similar to the MEF dataset, SKF-SNR fails to provide sufficient brightness enhancement, and both UHDFormer and Wave-Mamba suffer from blurriness. CWNet, on the other hand, produces visually superior results with

Metrics	UHDFour [19]	UHDFormer [39]	Wave-Mamba [61]	CWNet
UICM $\uparrow$	0.7464	0.9571	0.9082	<b>0.9663</b>
NIQE $\downarrow$	5.385	5.564	5.689	<b>4.494</b>

Table 5. Visualization quality comparison on DarkFace. The best results in each column are shown in bold.

richer detail information.

5) Fig. 15 demonstrates the qualitative comparison on the VV dataset. SKF-SNR generates distorted enhancement results with significant noise artifacts. UHDFormer suffers from overexposure, as observed in the highlighted facial regions, while Wave-Mamba produces blurry results. In contrast, CWNet produces natural enhancements with well-preserved details.

## B.3. Quantitative and Qualitative Comparison on the DarkFace Dataset

We conducted quantitative and qualitative comparisons against current state-of-the-art (SOTA) methods on the DarkFace dataset [50]. The DarkFace dataset, with 6,000 real-world low-light images, serves as a challenging benchmark for low-light image enhancement.

For quantitative evaluation, we employed two metrics: NIQE and the Underwater Image Colorfulness Measure (UICM)[34], which is commonly used to evaluate colorfulness and naturalness in enhanced images. Unlike NIQE, where lower scores indicate better perceptual quality, higher UICM scores reflect greater colorfulness and naturalness. As shown in Tab. 5, our method achieves the best performance across both metrics, highlighting its effectiveness in low-light image enhancement.

Fig. 17 illustrates the qualitative comparison on the DarkFace dataset. UHDFormer produces severe color distortions, particularly in bright regions, indicating its limited generalization ability on this dataset. In the first and last rows, the zoomed-in regions show that our method outperforms UHDFormer and Wave-Mamba in detail preservation and sharpness. Furthermore, in the second row, our method demonstrates superior exposure control, yielding natural and visually pleasing results. These observations validate the effectiveness of CWNet in low-frequency brightness control for natural exposure and high-frequency detail enhancement for sharper textures.

## C. Ablation Study and Downstream Application

### C.1. Ablation Study: Validating the Effectiveness of Core Modules

Fig. 9 demonstrates the impact of causal inference on CWNet’s performance. By focusing on the highlighted regions, it is evident that removing causal inference leads to

Number	PSNR $\uparrow$	SSIM $\uparrow$	LPIPS $\downarrow$
L=1, C=1	21.06	0.6381	0.1772
L=2, C=2	<b>21.50</b>	0.6397	<b>0.1562</b>
L=3, C=3	21.34	<b>0.6401</b>	0.1587

Table 6. Ablation Study on Negative Sample Quantity in the Causal Inference Module. The first column represents the number of  $I_l$  and  $I_c$ . The **best** results in each column are shown in bold.

inconsistencies in brightness and semantic coherence. In contrast, incorporating causal inference ensures both brightness restoration and semantic consistency, effectively preserving color fidelity and structural details. This validates the module’s ability to model causal relationships, enhancing overall enhancement quality.

## C.2. Ablation Study on Negative Sample Quantity in the Causal Inference

We conducted an ablation study within the Causal Inference module to determine the optimal number of negative samples for brightness degradation  $I_l$  (simulating underexposed regions) and color anomaly  $I_c$  (introducing color distortions), as summarized in Tab.6. The results indicate that the optimal performance is achieved when the number of samples for both brightness degradation and color anomaly is set to 3.

## C.3. Object Detection on DarkFace Dataset

To evaluate how enhanced images affect downstream tasks, we conducted object detection experiments on the DarkFace dataset [50]. We evaluated our method on 200 randomly selected images from the dataset using the official YOLOv5 model pretrained on the COCO dataset [25], as YOLOv5 is a widely used object detection framework and COCO provides a diverse set of object categories. Fig.18 presents the visual comparison results, showing that our method achieves superior detection accuracy compared to others.

In the first experiment, our method achieves higher confidence scores for detected pedestrians compared to other methods and uniquely detects the motorcycle on the right, likely due to its ability to enhance fine-grained details in low-light conditions. In the second experiment, our method detects more pedestrians with higher confidence and correctly identifies the traffic light, a task that is particularly challenging due to the small size and low contrast of the traffic light in the original image. In contrast, Wave-Mamba misclassifies building lights as traffic lights.

These results demonstrate that our enhancement method generates clearer images while preserving semantic structures effectively. By leveraging causal inference, CWNet emphasizes intrinsic image content, enhancing downstream

task performance such as object detection.

## C.4. Superior Edge Detection Performance

Fig. 16 compares edge detection performance across state-of-the-art methods, highlighting CWNet’s ability to restore details more precisely, particularly in highlighted regions. This superior performance validates the effectiveness of CWNet’s high-frequency module in preserving fine details and structural fidelity, further demonstrating the robustness of our architecture.



Figure 10. Visual comparison on LOL-v2-Synthesis dataset.



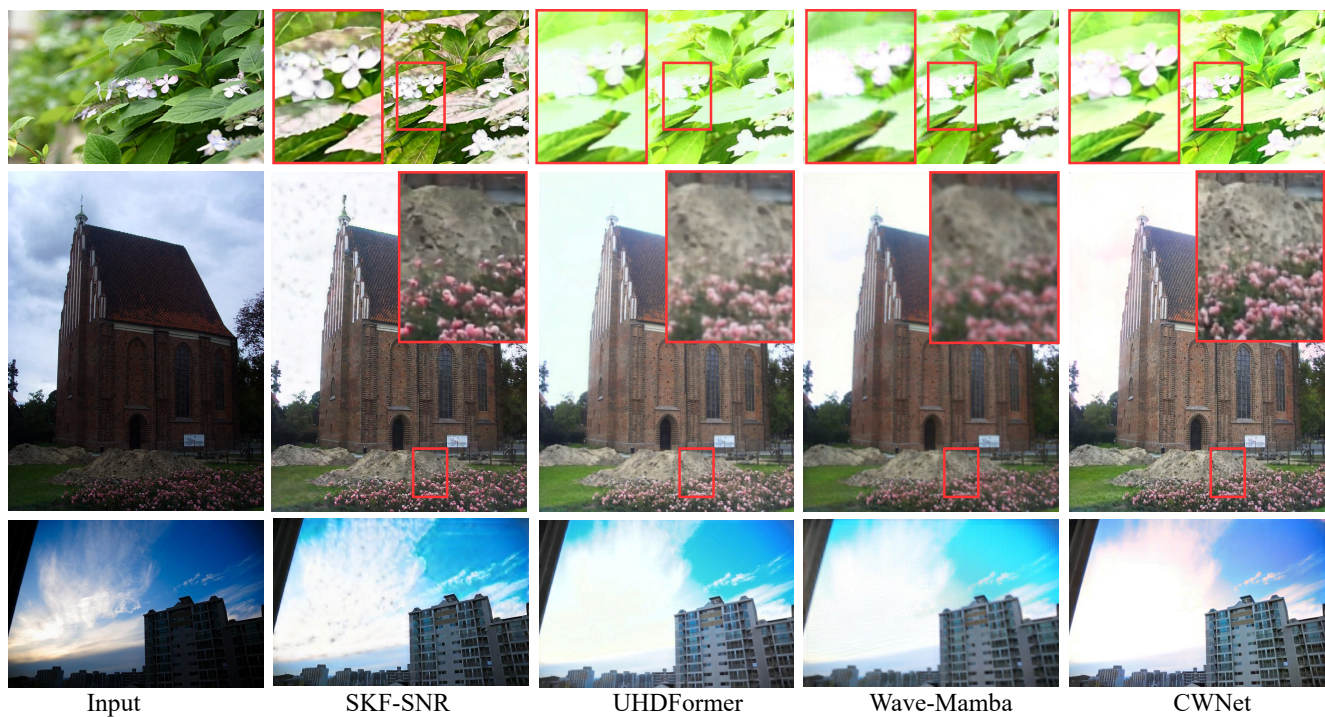


Figure 11. Visual comparison on DICM dataset.

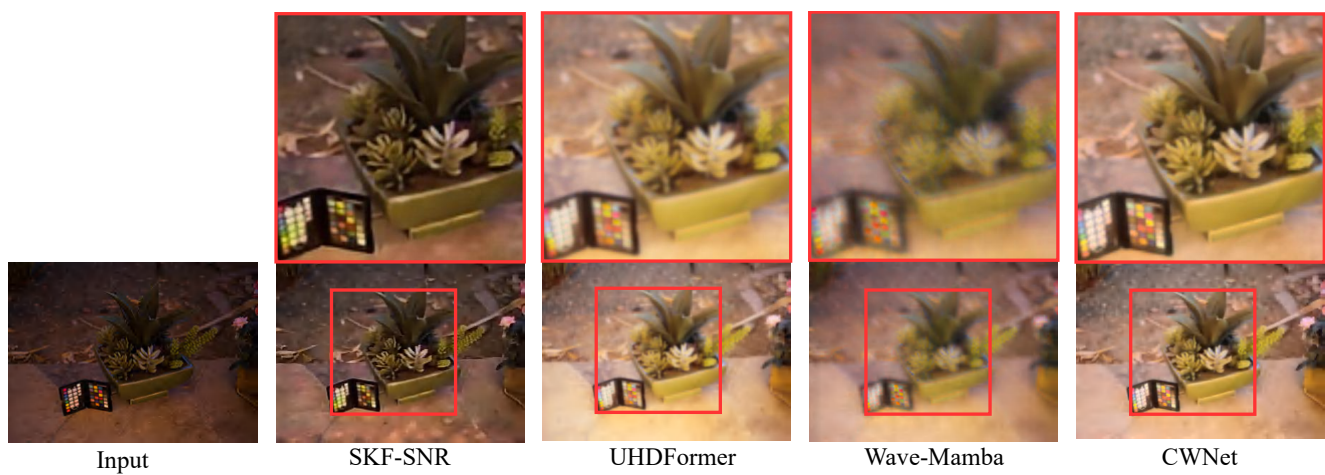


Figure 12. Visual comparison on LIME dataset.

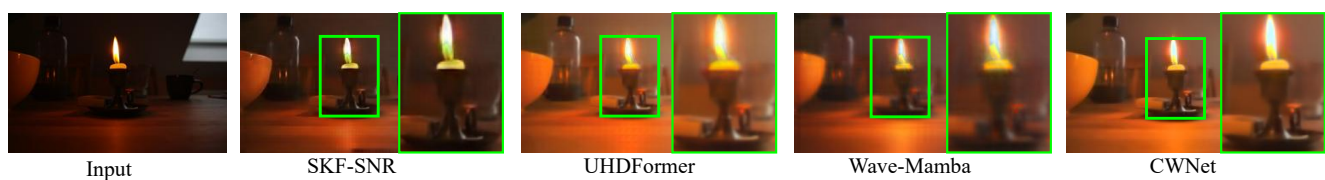


Figure 13. Visual comparison on MEF dataset.



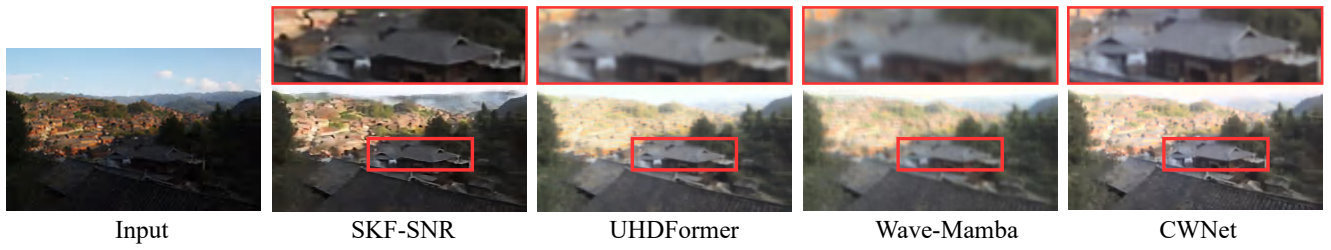


Figure 14. Visual comparison on NPE dataset.

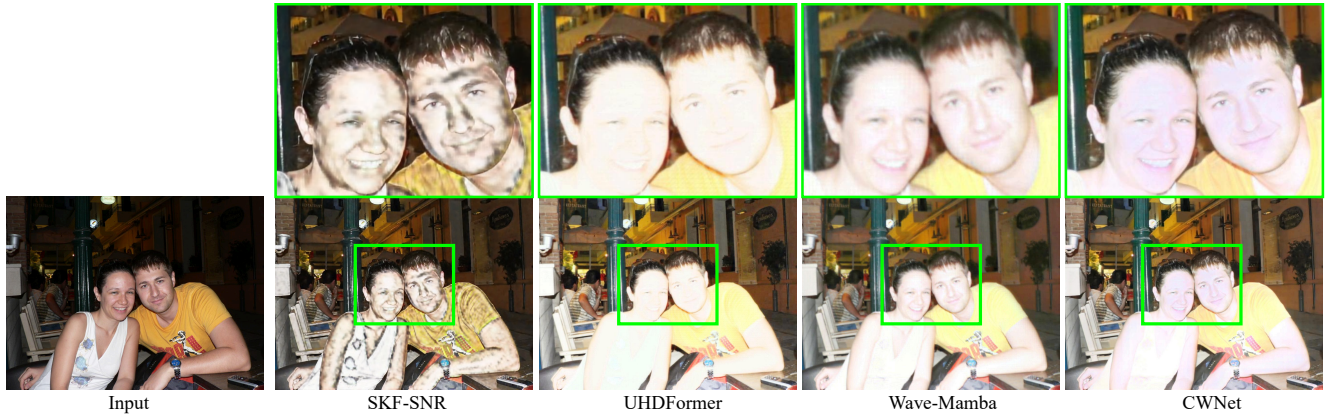


Figure 15. Visual comparison on VV dataset.

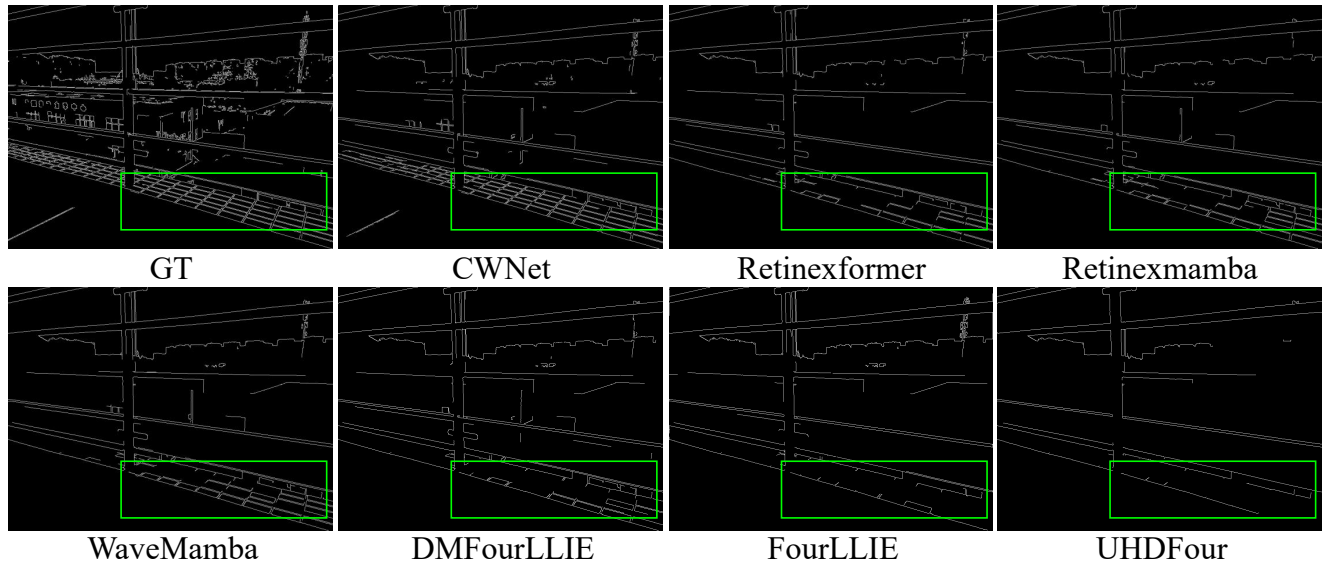


Figure 16. Edge detection comparison shows our method restores details more precisely, especially in highlighted regions, validating the high-frequency module.

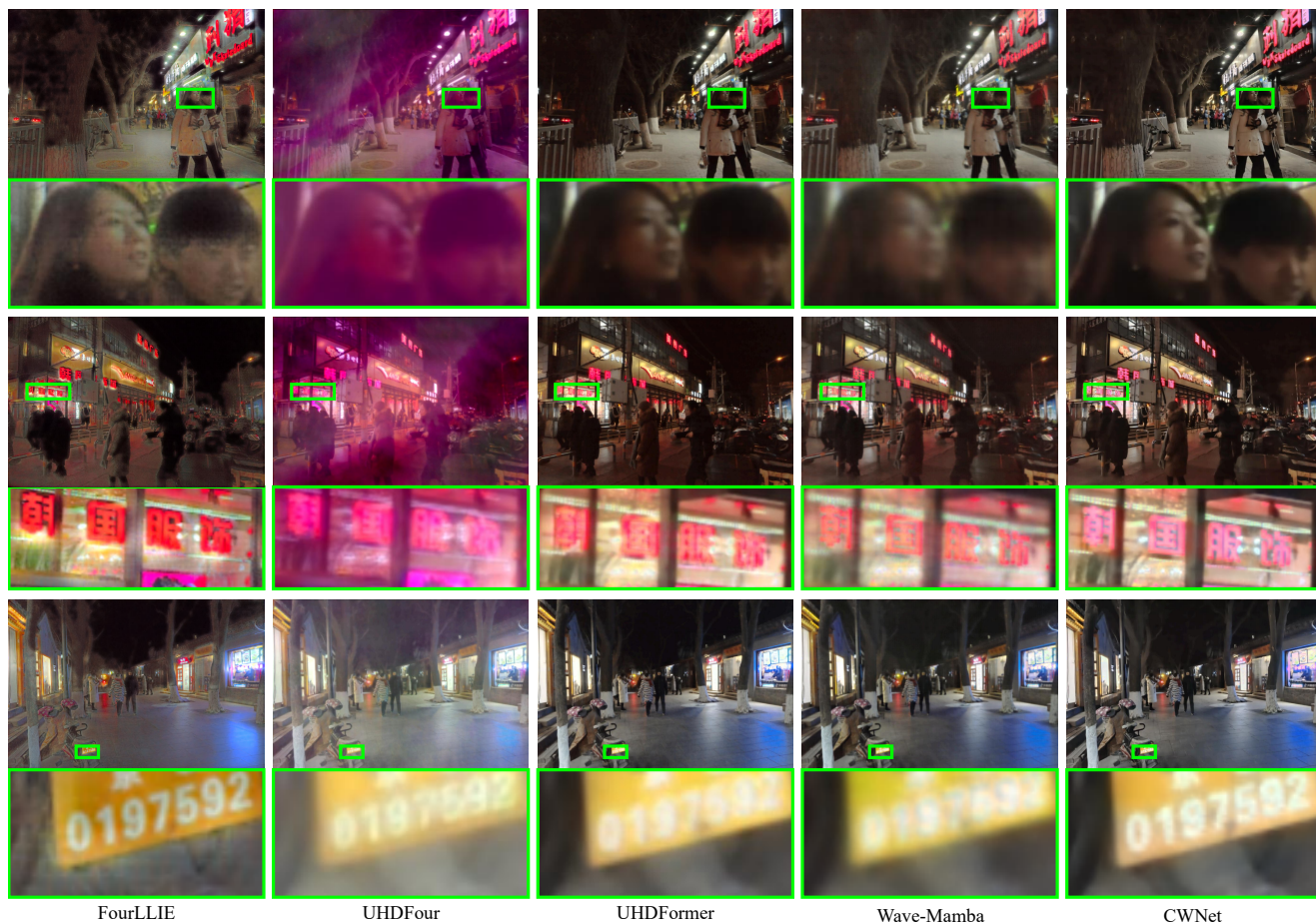
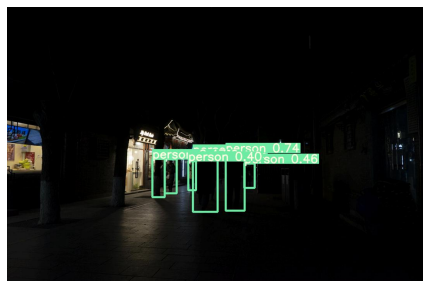


Figure 17. Visual comparison on DarkFace dataset.

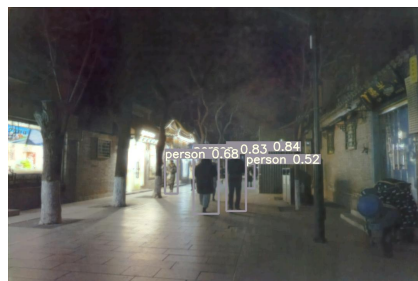




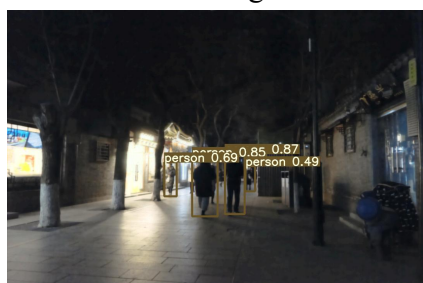
Low-Light



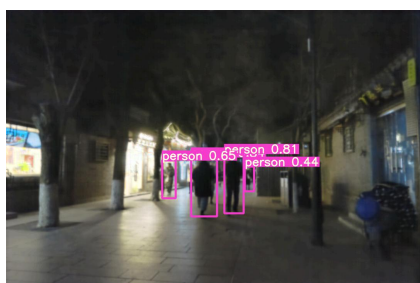
FourLLIE



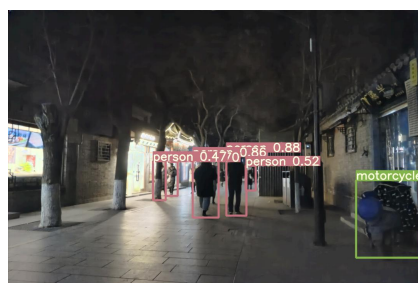
UHDFour



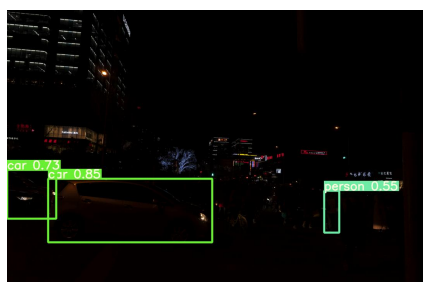
UHDFormer



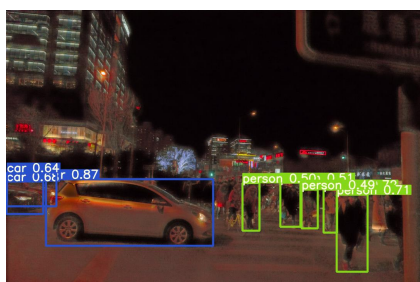
Wave-Mamba



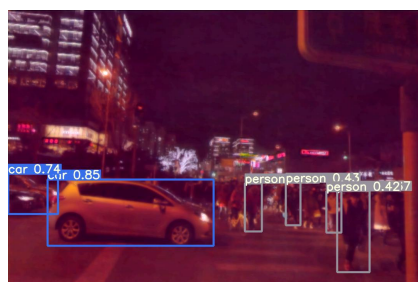
CWNet



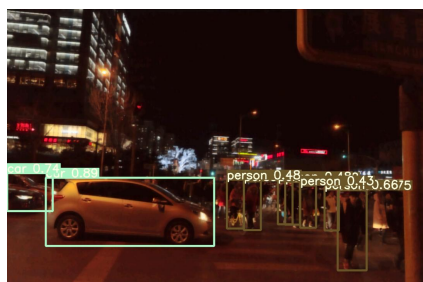
Low-Light



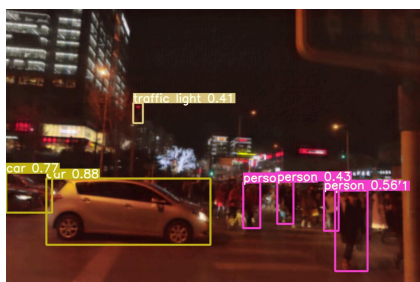
FourLLIE



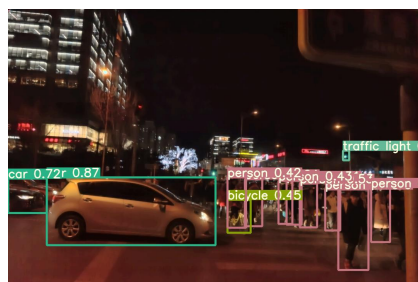
UHDFour



UHDFormer



Wave-Mamba



CWNet

Figure 18. Detection comparison results on DarkFace dataset.