
MambaLLIE: Implicit Retinex-Aware Low Light Enhancement with Global-then-Local State Space

Jiangwei Weng¹, Zhiqiang Yan¹, Ying Tai², Jianjun Qian¹, Jian Yang¹ and Jun Li¹

¹Nanjing University of Science and Technology, Nanjing, China

²Nanjing University, Nanjing, China

Abstract

Recent advances in low light image enhancement have been dominated by Retinex-based learning framework, leveraging convolutional neural networks (CNNs) and Transformers. However, the vanilla Retinex theory primarily addresses global illumination degradation and neglects local issues such as noise and blur in dark conditions. Moreover, CNNs and Transformers struggle to capture global degradation due to their limited receptive fields. While state space models (SSMs) have shown promise in the long-sequence modeling, they face challenges in combining local invariants and global context in visual data. In this paper, we introduce MambaLLIE, an implicit Retinex-aware low light enhancer featuring a global-then-local state space design. We first propose a Local-Enhanced State Space Module (LESSM) that incorporates an augmented local bias within a 2D selective scan mechanism, enhancing the original SSMs by preserving local 2D dependency. Additionally, an Implicit Retinex-aware Selective Kernel module (IRSK) dynamically selects features using spatially-varying operations, adapting to varying inputs through an adaptive kernel selection process. Our Global-then-Local State Space Block (GLSSB) integrates LESSM and IRSK with LayerNorm as its core. This design enables MambaLLIE to achieve comprehensive global long-range modeling and flexible local feature aggregation. Extensive experiments demonstrate that MambaLLIE significantly outperforms state-of-the-art CNN and Transformer-based methods. Project Page.

1 Introduction

low light image enhancement is a challenging task in computer vision due to insufficient lighting and sensor degradation. Consequently, images often suffer from poor global visibility and local issues such as color distortion and noise. These degraded images can adversely affect human perception and high-level vision tasks, such as object detection.

Traditional techniques, including histogram equalization [1] and gamma correction [5], enhance images through global mapping operations. However, these global operations often struggle to address local degradation effectively. In recent years, many methods based on CNNs and Transformers have gradually come to dominate this field [43, 52, 13, 31, 46, 3]. CNN-based methods [43, 52, 13, 31, 45] have achieved significant advancements by effectively aggregating local information, thus substantially improving performance in low light enhancement. Nevertheless, the limited receptive field and weight-sharing strategy of CNNs result in a local reductive bias, making the models less adaptive to varying inputs. On the other hand, Transformer-based methods [46, 3, 50] achieve a larger and adaptive receptive field by emphasizing long-term dependencies through the self-attention mechanism. However, the vanilla attention mechanism scales quadratically with input size, resulting in significant computational overhead.

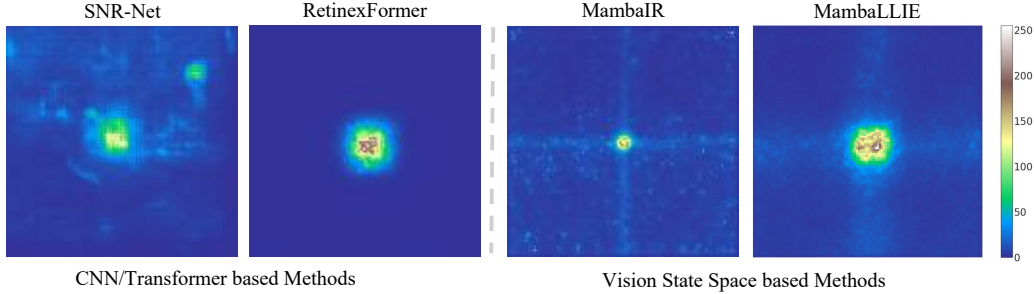


Figure 1: The Effective Receptive Field (ERF) visualization [28] for SNR-Net [46], RetinexFormer [3], MambaIR [14] and our MambaLLIE. A broader distribution of bright areas signifies a larger ERF. The receptive field of SNR-Net is large but messy, due to the SNR-aware mechanism, RetinexFormer achieves a larger receptive field of the central point, and MambaIR has the the global receptive field, but presents the limited local perception. Only our proposed MambaLLIE achieves a global perception ability outwards from central point and preserves the large local receptive field.

Recently, Mamba[8, 25, 22] have garnered significant attention in the field of computer vision. These internal state space models (SSMs) demonstrate great potential for global information modeling with linear complexity. However, a straightforward implementation of vision state space models for low light image enhancement is inadequate. This is because SSMs are primarily designed for long-range modeling and lack the flexibility to capture local information effectively [54]. For instance, as illustrated in Fig. 1, the receptive field of MambaIR [14], a simple yet effective vision state space model, achieves longer-range dependencies compared to CNN and Transformer-based methods. but it falls short in refining local interactions.

In this work, we introduce MambaLLIE, a novel framework for enhancing low light images that integrates an implicit Retinex-aware approach within a global-then-local state space model. MambaLLIE not only explores the capabilities of state space models in low light image enhancement but also incorporates a Retinex-aware structure providing both explicit and implicit guidance. Our framework introduces a unique global-then-local state space block, enhancing global long-range degradation modeling and local feature aggregation through an augmented state space. Additionally, we incorporate a Retinex-aware selective kernel mechanism in the enhancement process, enabling adaptive modulation of illumination strength through specific spatial operations.

Our contributions and main findings can be summarized as follows: 1) We introduce a novel global-then-local state space block that integrat a local-enhanced state space module and an implicit Retinex-aware selective kernel module. This design effectively captures intricate global and local dependencies. 2) We devise an implicit Retinex-aware selective kernel mechanism to guide deeper neural representations, eliminating the need for complex structural design and constraints to estimate physical priors, the prior feature tends to segregate into independent positive and negative illumination components before integrating them, a capability lacking in explicit methods. 3) Experimental results on benchmark datasets and real-world evaluations consistently demonstrate the superior performance of our proposed method compared to the state-of-the-art approaches.

2 Related work

Low Light Image Enhancement. Nowadays, the exciting deep learning-based methods have mainly been categorized into end-to-end and Retinex-based methods [21]. To the best our knowledge, LLNet [27] firstly introduced a deep neural network for low light image enhancement by supervised learning. LightenNet [2] adopted the CNN for single image contrast enhancement. MBLLEN [29] proposed the multi-branch fusion within CNN to extract rich features. Besides, SNR-Net [46], Restormer [50] LLFormer [18] and [30] adopted the self-attention mechanism to achieve excellent performance. However, all these end-to-end models mainly depend on the distribution of training dataset and ignore the inherent illumination prior. As contrast, ZeroDCE [13], RUAS [24], and subsequent works [31, 7, 41] represent impressive solutions for image enhancement, as ones precisely using physical priors to enhance the images. However, due to the absence of an ideal reference for guidance, these methods usually exhibit a certain gap compared to the supervised learning models.

As for the supervised Retinex-based models, these methods aim to decompose the image into illumination and reflectance maps, and then enhance the image by optimizing these maps. For instance, Retinex-Net [43] divided image enhancement into decomposition, adjustment and reconstruction stages, which providing a good representation of image enhancement process. KinD [52] and URetinex-Net [45] further introduced the novel multi-branch and multi-stage frameworks, respectively. However, striking a balance between complexity and efficiency remains challenging for these methods. Recently, RetinexFormer [3] simplified a one-stage Retinex-based low light enhancer with a efficient Transformer. Diff-Retinex [49] designed a transformer-based decomposition network and adopted generative diffusion networks to reconstruct the results. Overall, they typically applied the Retinex theory in a direct way, which may be limited for low light enhancement problem.

Vision State Space Model. State Space Model (SSMs) [11, 10, 9] are burgeoning new sequence models for deep learning, which first swept the natural language processing (NLP) community such as language understanding [35], content-based reasoning [54]. Recently, SSMs have also garnered considerable attention in computer vision (CV) tasks. To our knowledge, S4ND [32] first explored state space mechanism into CV tasks by swapping Conv2D and self-attention layers with S4ND layers in existing models. VMamba [25] bridged the gap between ordered sequences and non-causal visual images, enabling the extension of vision selective state space model with global receptive fields. Vim [53] proposed the bidirectional state space modeling with positional awareness, achieved the global visual perception. Furthermore, LocalMamba [15] was focused on the local scanning strategy, preservation of local context dependencies. EfficientVMamba [34] designed a light-weight SSMs with an additional convolution branch to learn both global and local representational features. MambaIR [14] employed convolution and channel attention to enhance the capabilities of the Mamba. But existing vision state space model does not pay enough attention on capturing local information, as vanilla SSMs are designed for long sequence and the invariant of local vision data is not taken into account in the existing vision state space models.

3 Methodology

This work aims to introduce a novel implicit Retinex-aware low light enhancer with global-then-local state space. In this section, we revisit the Retinex theory and the state space model, offering a concise overview. Following that, the details of our proposed MambaLLIE are introduced.

3.1 Preliminaries

Retinex Theory. The ideal Retinex theory [20] for low light enhancement assumes that the captured images can be decomposed into reflectance and illumination maps. Following [31, 37], explicit Retinex-based methods emphasize estimating either an illumination map while regarding the reflectance map as the enhanced result, or estimating concrete reflectance and illumination maps and then restoring the well-exposed images. Specifically, given a low light image $\mathbf{L} \in \mathbb{R}^{H \times W \times 3}$, where H and W represent height and width respectively, the derived maps can be denoted as:

$$\mathbf{L} = \mathbf{R} \cdot \mathbf{I}, \quad \mathbf{N} = \mathbf{L} / \tilde{\mathbf{I}}, \quad \mathbf{N} = \tilde{\mathbf{R}} \cdot \tilde{\mathbf{I}}. \quad (1)$$

where \cdot denotes the element-wise multiplication, $\mathbf{R} \in \mathbb{R}^{H \times W \times 3}$ denotes reflectance map, a static property of captured objects; $\mathbf{I} \in \mathbb{R}^{H \times W}$ denotes illumination map; $\mathbf{N} \in \mathbb{R}^{H \times W \times 3}$ denotes normal images; $\tilde{\mathbf{R}}, \tilde{\mathbf{I}} \in \mathbb{R}^{H \times W \times 3}$ denotes the estimated reflectance and illumination maps, respectively.

Consequently, the former assumption ignore the noise and artifacts resulting from sensor degradation in the captured images, while pixel-wise adjustments for the illumination is inadequate. The latter aims to restore the reflectance and illumination maps to enhance the images. However, this requires the design of multiple branches and constraints to guide the training [52].

State Space Model. The SSMs, such as structured state space sequence models (S4) [10] and Mamba [8], can be regarded as the continuous linear time-invariant (LTI) systems [44]. Given an one-dimension sequence $x(t) \in \mathbb{R}$, it projects into a new one-dimension sequence $y(t) \in \mathbb{R}$ through the hidden state $h(t) \in \mathbb{R}^m$, the whole system can be defined as a linear ordinary differential equation (ODE):

$$\begin{aligned} h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t), \\ y(t) &= \mathbf{C}h(t) + \mathbf{D}x(t). \end{aligned} \quad (2)$$

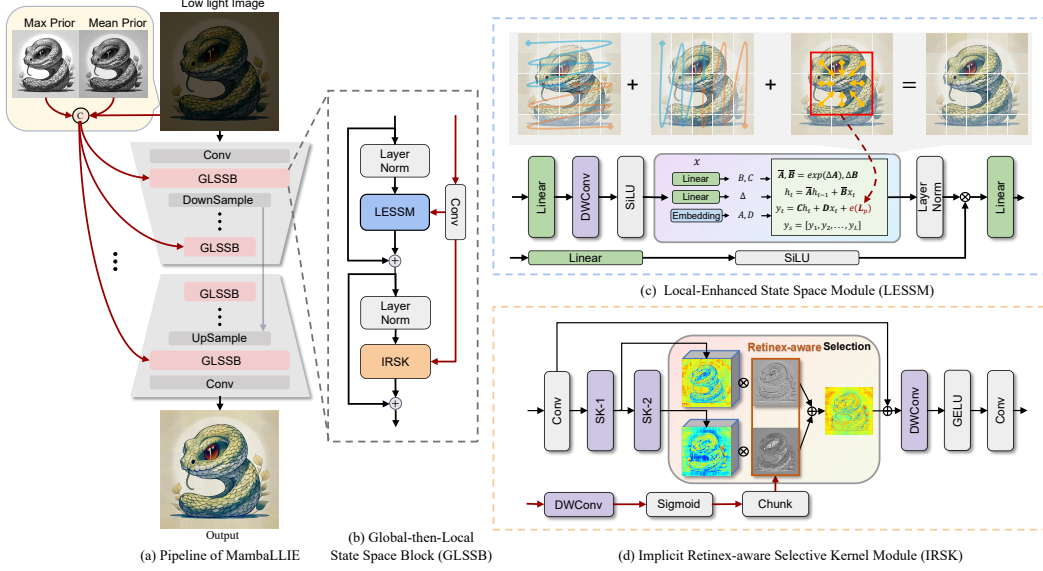


Figure 2: The overall pipeline of the proposed MambaLLIE. Our Global-then-Local State Space Block (GLSSB) integrates Local-enhanced state space module (LESSM) and implicit Retinex-aware selective kernel module (IRSK) with layer normalization as its core.

where m denotes the state size, $\mathbf{A} \in \mathbb{R}^{m \times m}$, $\mathbf{B} \in \mathbb{R}^{m \times 1}$, $\mathbf{C} \in \mathbb{R}^{1 \times m}$ and $\mathbf{D} \in \mathbb{R}$ denotes state, input projection, output projection, and feedthrough parameters.

As the raw state-space models are continuous, the systems adopt the discrete versions before feeding the computer, in which the zero-order hold (ZOH) is used to transform the continuous parameters \mathbf{A} and \mathbf{B} to discrete parameters $\bar{\mathbf{A}}$ and $\bar{\mathbf{B}}$ as follows

$$\bar{\mathbf{A}} = \exp(\Delta \mathbf{A}), \quad \bar{\mathbf{B}} = (\exp(\Delta \mathbf{A}) - \mathbf{I}) \cdot \Delta \mathbf{B}. \quad (3)$$

where Δ denotes a step size. Overall, the discretized version can be rewritten as:

$$h_t = \bar{\mathbf{A}} h_{t-1} + \bar{\mathbf{B}} x_t, \quad y_t = \mathbf{C} h_t + \mathbf{D} x_t. \quad (4)$$

However, the current system remains static for varying inputs. To address this limitation, Mamba [8] introduces selective state space models, allowing parameters to adapt with the input, thereby enhancing selective information processing across sequences. This parameter selection mechanism can be expressed as:

$$\bar{\mathbf{B}} = f_{\mathbf{B}}(x_t), \quad \bar{\mathbf{C}} = f_{\mathbf{C}}(x_t), \quad \Delta = \vartheta_{\mathbf{A}}(\mathbf{P} + f_{\mathbf{A}}(x_t)). \quad (5)$$

where $f_{\mathbf{B}}(x_t)$, $f_{\mathbf{C}}(x_t)$ and $f_{\mathbf{A}}(x_t)$ are linear functions that broadens feature to the hidden state dimensions. As SSMs are tailored for long sequences, it is limited in capturing complicated local information. For visual data, VMamba [25] and Vim [53] proposed the specific location-aware scan strategies to maintains the integrity of 2D image structures. However, the specific directed sequences overlook the vision information of pixels neighborhood structure. Inspired by [54], we explore a global-then-local state space, which receives the global perception before the details, supplementing the lack of local information.

3.2 Overall Pipeline

We first present the overall pipeline of our MambaLLIE, an U-shaped architecture as shown in Fig. 2(a), which includes encode and decode parts with the convolutional downsampling and upsampling layers. The encoder features are concatenated with the decoder features via skip connections. Next, We propose a global-then-local state space block (GLSSB) as the basic core of MambaLLIE, the max and mean priors concatenated with low light image are projected into GLSSB by convolutional layers as augmented input. Therein, GLSSB is composed of the local-enhanced state space module (LESSM) and the implicit Retinex-aware selective kernel module (IRSK), interleaved with layer norm layer.

Specifically, given a low light image $\mathbf{L} \in \mathbb{R}^{H \times W \times 3}$, we employ a 3×3 convolution layer to project the neural features $\mathbf{F} \in \mathbb{R}^{H \times W \times C}$ from input feature space, and then project features into each GLSSB, which will be described in Section 3.3. Besides, IRSK integrates original input, maximum prior $\mathbf{L}_{\max} \in \mathbb{R}^{H \times W}$ and mean prior $\mathbf{L}_{\text{mean}} \in \mathbb{R}^{H \times W}$ as augmented input $\mathbf{L}_p \in \mathbb{R}^{H \times W \times 5}$,

$$\mathbf{L}_p = \text{Concat}(\mathbf{L}, \text{mean}(\mathbf{L}), \text{max}(\mathbf{L})). \quad (6)$$

We first define \mathbf{F}_g is the output of GLSSB. Subsequently, the downsampling layer and following GLSSB achieve the feature extraction to acquire the deep feature, which can be denoted as $\mathbf{F}_g \in \mathbb{R}^{\frac{H}{2^i} \times \frac{W}{2^i} \times 2^i C}$, where $i = 0, 1, 2$. Moreover, the feature is later concatenated with the upsampling layer with a symmetrical structure. Finally, using a 3×3 convolution layer projects into $\mathbf{F}_{\text{out}} \in \mathbb{R}^{H \times W \times 3}$ and the enhanced image can be expressed as $\mathbf{N} = \mathbf{F}_{\text{out}} + \mathbf{L}$.

3.3 Global-then-Local State Space Block

As illustrated in Fig. 2(b), GLSSB follows the LayerNorm, LESSM, LayerNorm and IRSK flow, motivated by Transformer [38] and Mamba[8] usage of similar structures in a basic block. Given the input feature, it first undergoes the LayerNorm and LESSM to capture the local-enhanced global information. the above process can be denoted as:

$$\mathbf{M} = \text{LESSM}(\text{LN}(\mathbf{F}_g^{i-1})) + \mathbf{F}_g^{i-1}. \quad (7)$$

And then, another LayerNorm and our proposed IRSK are used for Retinex-aware guidance. The above process can be formulated as:

$$\mathbf{F}_g^i = \text{IRSK}(\text{LN}(\mathbf{M})) + \mathbf{M}. \quad (8)$$

Overall, at the prior module of GLSSB, we capture global dependencies using a local-enhanced SSM. Because the SSM is better at learning global information, the subsequent module aims to handle more refined and complicated local dependencies.

Local-Enhanced State Space Module. Existing state space models [6, 10, 8] excels at capturing the causal processing of input data in long range dependencies. However, the unidirectional scan manner encounters difficulties in vision data to modeling non-causal relationships. To accommodate vision data, [53, 25, 34] process the input data from different 2D scan directions. However, these methods ignore the local invariants of vision data. In other word, the fixed scanning methods widen the distance between neighborhood data and snarl the causal relationships.

The most SSMs [25] can be regarded as the continuous linear time-invarian systems, we further introduce the a $e(\mathbf{L}_p)$ augmented local bias, enhancing the original SSMs by preserving local 2D dependency as shown in Fig. 2(c). Following [47, 16], we propose a global-then-local state space:

$$\begin{aligned} h_t &= \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t, \\ y_t &= \mathbf{C}h_t + \mathbf{D}x_t + e(\mathbf{L}_p). \end{aligned} \quad (9)$$

where $e(\mathbf{L}_p)$ is independent of the hidden state space. Hence, the model can be computed in a simple way, given a feature $\mathbf{F} \in \mathbb{R}^{H \times W \times 5}$ and illumination feature $\mathbf{L}_p \in \mathbb{R}^{H \times W \times C}$, we adopts the Layernorm followed by our proposed LESSM to integrate the spatial long-term dependency. Following [8], the input feature are chunk into $\tilde{\mathbf{F}}_1$ and $\tilde{\mathbf{F}}_2$ in two branches. The first branch projects the feature into a linear layer, followed by a depth-wise convolution, SiLU activation function, accompanied by our proposed augmented local bias and Layernorm. In the second branch, the features is also projected to a linear layer followed by the SiLU activation function. Finally, features from the two branches are aggregated with the element-wise product and then are projected back to input feature space by linear layer. The entire process can be delineated as

$$\begin{aligned} \tilde{\mathbf{F}}_1 &= \text{LN}(2\text{DSSM}(\text{SiLU}(\text{DWConv}(\text{Linear}(\mathbf{F}_1)))) + \text{Conv}(\mathbf{L}_p)) \\ \tilde{\mathbf{F}}_2 &= \text{SiLU}(\text{Linear}(\mathbf{F}_2)) \\ \tilde{\mathbf{F}}_{\text{out}} &= \text{Linear}(\tilde{\mathbf{F}}_1 \odot \tilde{\mathbf{F}}_2) \end{aligned} \quad (10)$$

Implicit Retinex-Aware Selective Kernel Module. We further employ a implicit Retinex-aware selective kernel network to enhance the capabilities of the information integration ability. IRSK

Table 1: Quantitative comparisons on LOL-V2-real, LOL-V2-syn, SMID, SDSD-indoor and SDSD-outdoor datasets. The best result is in red color while the second best result is in blue color.

Methods	Ref.	LOL-V2-real		LOL-V2-syn		SMID		SDSD-indoor		SDSD-outdoor		Complexity	
		PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	PSNR	SSIM	FLOPS	Param
RetinexNet	BMVC 2018	15.47	0.567	17.13	0.798	22.83	0.684	20.84	0.617	20.96	0.629	587.47	0.84
DeepUPE	CVPR 2019	13.27	0.452	15.08	0.623	23.91	0.690	21.70	0.662	21.94	0.698	21.10	1.02
SID	ICCV 2019	13.24	0.442	15.04	0.610	24.78	0.718	23.29	0.703	24.90	0.693	13.73	7.76
KinD	MM 2019	14.74	0.641	13.29	0.578	22.18	0.634	21.95	0.672	21.97	0.654	34.99	8.02
MIRNet	ECCV 2020	20.02	0.820	21.94	0.876	25.66	0.762	24.38	0.864	27.13	0.837	785.00	31.76
EnGAN	TIP 2021	18.23	0.617	16.57	0.734	22.62	0.718	23.29	0.703	24.90	0.693	61.01	114.35
Restormer	CVPR 2022	19.94	0.827	21.41	0.830	26.97	0.758	25.67	0.827	24.79	0.802	144.25	26.13
SNR-Net	CVPR 2022	21.48	0.849	24.14	0.928	28.49	0.805	29.44	0.894	28.66	0.866	26.35	4.01
QuadPrior	CVPR 2024	20.48	0.811	16.11	0.758	15.50	0.604	22.22	0.783	18.26	0.662	/	/
MambaIR	Arxiv 2024	21.25	0.831	25.55	0.929	27.07	0.774	28.97	0.884	29.75	0.861	60.66	4.30
RetinexFormer	ICCV 2023	22.80	0.840	25.67	0.930	29.15	0.815	29.77	0.896	29.84	0.877	15.57	1.61
MambaLLIE	/	22.95	0.847	25.87	0.940	29.26	0.818	30.12	0.900	30.00	0.869	20.85	2.28

constructs a sequence of depth-wise convolutions with an alterable kernel to select the feature with different receptive field, using a spatial selection mechanism by illumination prior. Inspired by LSKNet [23], for each of the feature maps from different selective kernel, a Sigmoid activation function is applied to obtain the individual illumination maps from illumination prior. Fig. 2(d) shows a detailed conceptual illustration of IRSK module where we intuitively demonstrate how the implicit Retinex-aware module works. The above process can be formulated as:

$$\tilde{\mathbf{F}}_k = \tilde{\mathbf{F}}_{\text{out}}, \quad \tilde{\mathbf{F}}_{k+1} = f_{\text{DWconv}}^k(\tilde{\mathbf{F}}_k). \quad (11)$$

The output of the Retinex-aware maps are concatenated with the input features via residual connections, followed by a depth-wise convolution, GELU activation function and convolution layer.

$$\{\mathbf{S}_1, \mathbf{S}_2\} = \text{Chunk}(\text{Sigmoid}(\text{Conv}(\mathbf{L}_p))) \quad (12)$$

$$\mathbf{F}_g = \text{Conv}(\text{GELU}(\text{DWConv}(\sum_{k=1}^K \tilde{\mathbf{F}}_k \mathbf{S}_k + \tilde{\mathbf{F}}_{\text{out}}))) \quad (13)$$

4 Experiments

4.1 Benchmark Datasets and Implementation Details

Datasets. We employ five paired low light image datasets for evaluation, including LOL-V2-real [48], LOL-v2-syn [48], SMID [4], SDSD-indoor [39] and SDSD-outdoor [39] datasets. Therein, LOL-V2-real contains 689 low-normal light paired images for training and 100 pairs for testing; LOL-V2-syn includes 900 paired images for training and the 100 pairs for testing; Besides, SMID is composed of the 15763 short-long exposure paired images for training and the remaining images for testing; SDSD-indoor and SDSD-outdoor are all subsets of SDSD dataset (the static version), which extract the paired images from 62 and 116 pairs for training, and the left 6 and 10 pairs for testing.

Implementation Details. We implement MambaLLIE in PyTorch [33] on a server with the 4090GPUs. Random cropping the image pairs into 128×128 patches as training samples, data augmentation is performed on the training samples such as rotation and flipping. The batch size is 8. In terms of optimization procedure, Adam [19] is adopted as the optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$; The training iterations is set to 1.5×10^5 . The initial learning rate is set to 2×10^{-4} and steadily decreased by the cosine annealing scheme. The loss criterion is mean absolute error (MAE), thus peak signal-to-noise ratio (PSNR) and structural similarity (SSIM) [42] is selected as the evaluation metrics for the paired datasets.

4.2 Main Results on Benchmarks.

Quantitative Comparison. As shown in Tab. 1, we evaluated the performance of our MambaLLIE against 11 SOTA image enhancement methods, including RetinexNet [43], DeepUPE [40], SID [4], KinD [52], MIRNet [51], EnGan [17], Restormer [50], SNR-Net [46], QuadPrior [41], MambaIR [14] and Retinexformer [3]. Our MambaLLIE demonstrates superior performance than SOTA methods on

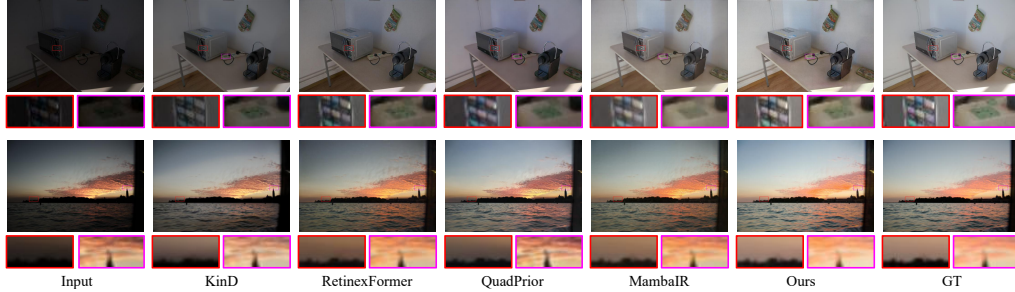


Figure 3: Qualitative comparison with previous methods on LOL-V2-real and LOL-V2-syn datasets. Our MambaLLIE effectively enhances the illumination and preserves the color.

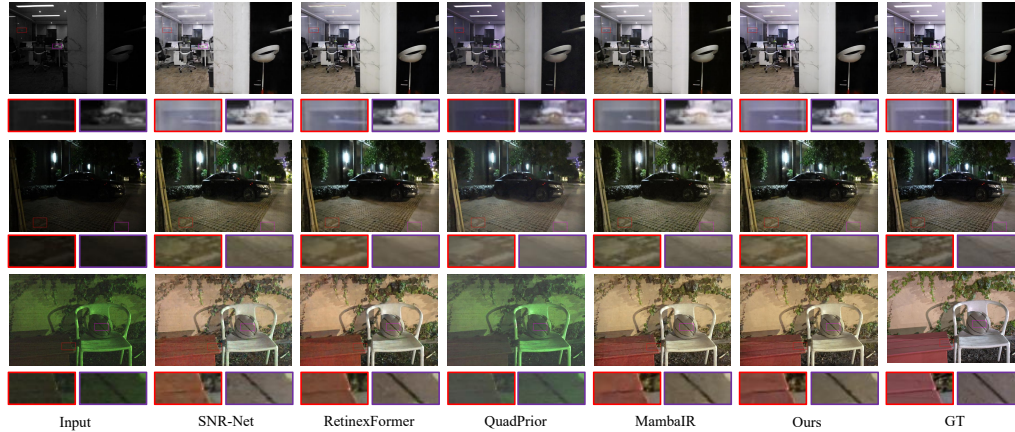


Figure 4: Qualitative comparison with previous methods on SMID, SDDS-indoor and SDDS-outdoor datasets. Our MambaLLIE restore the texture and color under challenging degradation, such as the wooden bench and reflective glass.

the adopted benchmark datasets in terms of PSNR and SSIM, while achieves comparable results of SSIM with the SOTA methods in LOL-V2-real and SDDS-outdoor. Therein, when the parameters are roughly similar, our MambaLLIE achieves an average improvement of 0.2 dB on benchmark datasets compared to the previous Transformer based SOTA method, *i.e.* RetinexFormer. Compared with the earlier Transformer based SNR-Net, MambaLLIE outperforms it by average 1dB PSNR on the all datasets. When compared to the MambaIR, MambaLLIE achieves 1.70, 0.32, 2.19, 1.15 and 0.25 dB PSNR improvements on the adopted datasets, respectively. Besides, Our MambaLLIE gains the improvements over 7 dB on all datasets than traditional Retinex-based models, such as RetinexNet, DeepUPE and KinD.

Qualitative Comparison. Figs. 3 & 4 report the vision results for comparing our method with latest the SOTA methods and traditional Retinex-based models. Existing methods suffer from insufficient illumination and fail to restore the details as shown in Fig. 3. As we can see, color distortion and image degradation also affect the enhanced results of previous methods in Fig. 4, yet our MambaLLIE not only enhances brightness but also faithfully preserves colors with reference to ground truth images, all while restoring the details.

4.3 Real World Experimental Evaluation

Enhancing low light images in real-world scenarios is exceptionally challenging because not only can provide benefits for downstream tasks, such as dark object detection, but must the enhanced images be pleasing to the human perception.

Low Light Object Detection. We utilized ExDark dataset [26] to compare the enhancement of preprocessing methods for high-level vision tasks. There are 7363 challenging low light images annotated with 12 bounding box classes, of which 5,890 for training and 1,473 for testing. Note that all supervised methods were pretrained on the LOL-V2-syn dataset, the low light image underwent different enhancement methods and then finetuned YOLOv3 [36] as the object detector.



Figure 5: Vision comparison of our MambaLLIE with recent SOTA methods.(a) Qualitative comparison on object detection, (b) A toy example of user study.

Table 2: Low light object detection results on the ExDark dataset. The best result is in red color while the second best result is in blue color.

Methods	Bicycle	Boat	Bottle	Bus	Car	Cat	Chair	Cup	Dog	Motor	People	Table	Mean
RetinexNet	0.790	0.741	0.743	0.908	0.820	0.665	0.651	0.750	0.721	0.703	0.784	0.556	0.736
EnGAN	0.733	0.710	0.687	0.892	0.786	0.675	0.656	0.650	0.741	0.657	0.731	0.528	0.704
KinD	0.800	0.721	0.788	0.919	0.822	0.718	0.672	0.771	0.775	0.736	0.803	0.555	0.757
ZeroDCE	0.806	0.750	0.762	0.914	0.837	0.681	0.677	0.769	0.788	0.728	0.801	0.535	0.754
SCI	0.821	0.742	0.749	0.916	0.846	0.695	0.690	0.784	0.756	0.758	0.810	0.555	0.760
SNR-Net	0.802	0.721	0.750	0.932	0.840	0.694	0.677	0.758	0.763	0.755	0.789	0.559	0.753
Retinexformer	0.809	0.769	0.753	0.914	0.814	0.688	0.689	0.763	0.766	0.769	0.805	0.543	0.757
MambaIR	0.803	0.763	0.752	0.903	0.830	0.687	0.684	0.761	0.721	0.738	0.813	0.556	0.751
MambaLLIE	0.802	0.764	0.779	0.926	0.846	0.701	0.692	0.800	0.781	0.751	0.812	0.560	0.768

As shown in Tab. 2, our methods achieved the best average result compared with all adopted models, and yielded the best results on Car, Chair, Cup and Table classes. Fig. 5(a) further reported the visual comparison, compared with suboptimal preprocessing method SCI, detector through our MambaLLIE can detect the objects in extreme dark regions including two persons and a chair, while other methods failed.

User Study. We conducted a user study to evaluate the human visual perception quality of the enhanced results in challenging scenarios. Due to the lack of the ideal reference for training, we selected the pretrained model from the benchmarks to enhance the photos. There are 7 random selected low light images from the benchmarks and ExDark datasets under different lighting conditions. Human perception primarily focuses on the presence of global visual effect, local detail, color distortion (noise), which significantly reflect the quality of the enhanced images. Thus, We assigned ratings on a scale of 1 (worst) to 5 (best), evaluating the quality of the enhancements in terms of overall rating, local detail and color distortion(noise), respectively. Overall, 70 participants were invited to assess the visual quality. The average scores are reported in Tab. 3, our MambaLLIE achieves the best score in the involved voting aspects. Fig. 5(b) shows the toy example of user study, which display the input and the random enhanced results and local details by different algorithms.

Table 3: User study on the challenging low light image enhancement.

Methods	RetinexNet	EnGAN	SCI	QuadPrior	SNR-Net	Retinexformer	MambaIR	MambaLLIE
Overall Rating	3.093	3.314	3.943	3.014	3.821	4.100	3.857	4.243
Local Detail	2.871	3.143	3.686	3.129	3.779	3.950	3.629	4.129
Color distortion(noise)	2.914	3.164	3.776	2.929	3.657	3.971	3.750	4.100

4.4 Ablation Study

Implicit Retinex-Aware Framework. We compare the improvement of using an implicit Retinex-aware model with the end-to-end and explicit Retinex-aware models. Specifically, Baseline-1 is a simple variant of our MambnaLLIE by removing Retinex-aware guidance, namely estimates enhanced result directly from input without any prior. Baseline-2 is designed to estimate the illumination map and then light up the low light image by element-wise multiplication. Tab. 4 reveals our implicit Retinex-Aware framework significantly outperforms Baseline-1 with the improvement of 1.25 dB in PSNR, while achieving a PSNR enhancement of 1.00 dB compared to Baseline-2.

Global-then-Local State Space. As the core component, our GLSSB comprises the LESSM and IRSK. We demonstrate the effect of each component through ablation study. The results, presented at the bottom of Tab. 4, indicate that our LESSM achieves improvements of 0.33 dB and 0.08 dB in PSNR compared to Baseline-1 and Baseline-2, respectively, which utilize vanilla state space blocks. Additionally, our IRSK produces PSNR enhancements of 0.96, 0.74, and 0.63 dB compared

Table 4: Effects of design choices.

Methods	Params	FLOPS	PSNR	SSIM
Baseline-1	2.14	18.39	28.87	0.865
Baseline-2	2.14	18.39	29.12	0.862
Ours w/o LESSM	2.26	20.64	29.83	0.889
Ours w/o IRSK	2.19	19.94	29.20	0.887
Ours	2.28	20.85	30.12	0.900

Table 5: Effects of different selective kernel.

Kernel Sizes	Params	FLOPS	PSNR	SSIM
3*3	2.25	20.47	29.55	0.899
5*5	2.31	21.23	29.48	0.896
5*7	2.35	21.79	28.88	0.892
5*3	2.38	20.85	29.31	0.892
3*5 (Ours)	2.28	20.85	30.12	0.900

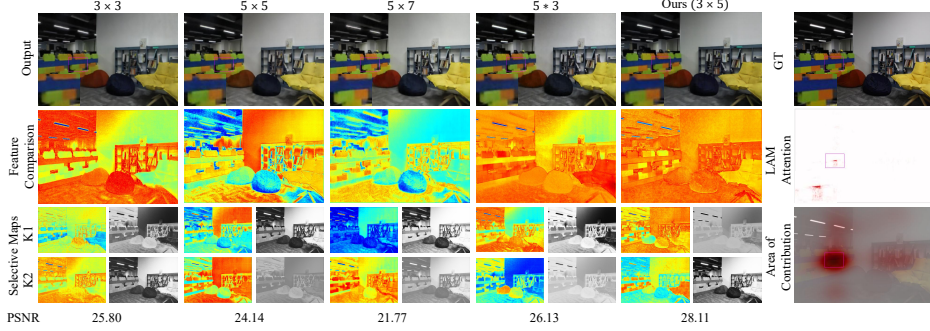


Figure 6: The details of selective kernel behaviour, the LAM visualization [12] demonstrates influence of similar local information is higher than that of global dependence, our local-enhanced strategy underscores the feature. Besides, the larger receptive fields can provide globally consistent results.

to Baselines and when only applying LESSM. Our full version indicating that although LESSM improves the vanilla SSM with local enhanced modeling ability, IRSK should be considered for further improvements, when GLSSB integrates LESSM and IRSK, our MambaLLIE achieves the highest PSNR and SSIM.

Selective Kernel Behaviour. We further investigate the kernel selection behaviour in our MambaLLIE as shown in Fig. 6. We find the implicit Retinex-aware selection pattern tend to learn two independent positive and negative illumination, resulting in complementary features. Compared with explicit Retinex-based methods, our IRSK can guide from a flexible deeper neural representation. The quantitative results are reported in Tab. 5. Different with LSKNet [23], we put small kernels in front and larger kernels in higher levels. This is because object detection needs larger receptive field, thus adopts a sequence of depth-wise convolutions with growing kernel and increasing dilation, while has to introduce a lots of padding. But image enhancement may suffers from padding operation at the edge of the image, especially upsampling further expands the padding values. Thus, the the former small kernels can quickly focus on local information and the the latter kernels contain larger receptive fields for better feature fusion.

5 Limitation and Discussion

We adopt an implicit Retinex-Aware guidance within a global-then-local state space framework to address global insufficient illumination and local degradation for low light enhancement. However, our approach has several limitations. 1) Unlike end-to-end methods, our technique requires the design of a reasonable illumination prior, which relies on prior experience. 2) Most existing enhancement models, including ours, primarily focus on mean square error and use PSNR and SSIM to evaluate image quality. To mitigate inherent biases in these metrics, we conducted additional real-world experimental evaluations to reconcile the bias and further validate the effectiveness of our approach.

6 Conclusion

In this paper, we introduced a novel state space-based model, MambaLLIE. Our proposed core of GLSSB effectively combines global and local information by implicit Retinex-aware selective kernel into global-then-local state space. Extensive experiments on benchmarks, low light object detection and user study demonstrate that our framework consistently achieves the best performance. Our future work is to address the dual challenges of local redundancy and global dependencies in low light video enhancement via efficient state space modeling.

References

- [1] Mohammad Abdullah-Al-Wadud, Md. Hasanul Kabir, M. Ali Akber Dewan, and Oksam Chae. A dynamic histogram equalization for image contrast enhancement. *IEEE Trans. Consumer Electron.*, 53(2):593–600, 2007.
- [2] Jianrui Cai, Shuhang Gu, and Lei Zhang. Learning a deep single image contrast enhancer from multi-exposure images. *IEEE Trans. Image Process.*, 27(4):2049–2062, 2018.
- [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 *ICCV*, pages 12504–12513, October 2023.
- [4] Chen Chen, Qifeng Chen, Minh N. Do, and Vladlen Koltun. Seeing motion in the dark. In *ICCV*, pages 3184–3193. IEEE, 2019.
- [5] Hany Farid. Blind inverse gamma correction. *IEEE Trans. Image Process.*, 10(10):1428–1433, 2001.
- [6] Daniel Y. Fu, Tri Dao, Khaled Kamal Saab, Armin W. Thomas, Atri Rudra, and Christopher Ré. Hungry hungry hippos: Towards language modeling with state space models. In *ICLR*. OpenReview.net, 2023.
- [7] Zhenqi Fu, Yan Yang, Xiaotong Tu, Yue Huang, Xinghao Ding, and Kai-Kuang Ma. Learning a simple low-light image enhancer from paired low-light instances. In *CVPR*, pages 22252–22261, 2023.
- [8] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *ArXiv*, abs/2312.00752, 2023.
- [9] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, *NeurIPS*, 2022.
- [10] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *ICLR*, 2022.
- [11] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. In *NeurIPS*, pages 572–585, 2021.
- [12] Jinjin Gu and Chao Dong. Interpreting super-resolution networks with local attribution maps. In *CVPR*, pages 9199–9208. Computer Vision Foundation / IEEE, 2021.
- [13] Chunle Guo, Chongyi Li, Jichang Guo, Chen Change Loy, Junhui Hou, Sam Kwong, and Runmin Cong. Zero-reference deep curve estimation for low-light image enhancement. In *CVPR*, pages 1777–1786, 2020.
- [14] 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. *arXiv preprint arXiv:2402.15648*, 2024.
- [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.
- [16] Salim Ibrir and Sette Diopt. Novel lmi conditions for observer-based stabilization of lipschitzian nonlinear systems and uncertain linear systems in discrete-time. *Applied Mathematics and Computation*, 206(2):579–588, 2008.
- [17] Yifan Jiang, Xinyu Gong, Ding Liu, Yu Cheng, Chen Fang, Xiaohui Shen, Jianchao Yang, Pan Zhou, and Zhangyang Wang. Enlightengan: Deep light enhancement without paired supervision. *IEEE Trans. Image Process.*, 30:2340–2349, 2021.
- [18] Haoxiang Jie, Xinyi Zuo, Jian Gao, Wei Liu, Jun Hu, and Shuai Cheng. Llformer: An efficient and real-time lidar lane detection method based on transformer. In Wenbing Zhao and Xinguo Yu, editors, *PRIS*, pages 18–23, 2023.
- [19] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR (Poster)*, 2015.
- [20] Edwin H Land and John J McCann. Lightness and retinex theory. *Josa*, 61(1):1–11, 1971.
- [21] Chongyi Li, Chunle Guo, Linghao Han, Jun Jiang, Ming-Ming Cheng, Jinwei Gu, and Chen Change Loy. Low-light image and video enhancement using deep learning: A survey. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(12):9396–9416, 2022.
- [22] Kunchang Li, Xinhao Li, Yi Wang, Yanan He, Yali Wang, Limin Wang, and Yu Qiao. Videomamba: State space model for efficient video understanding. *arXiv preprint arXiv:2403.06977*, 2024.
- [23] Yuxuan Li, Qibin Hou, Zhaohui Zheng, Ming-Ming Cheng, Jian Yang, and Xiang Li. Large selective kernel network for remote sensing object detection. In *ICCV*, pages 16748–16759. IEEE, 2023.
- [24] Risheng Liu, Long Ma, Jiaao Zhang, Xin Fan, and Zhongxuan Luo. Retinex-inspired unrolling with cooperative prior architecture search for low-light image enhancement. In *CVPR*, pages 10561–10570. Computer Vision Foundation / IEEE, 2021.

- [25] Yue Liu, Yunjie Tian, Yuzhong Zhao, Hongtian Yu, Lingxi Xie, Yaowei Wang, Qixiang Ye, and Yunfan Liu. Vmamba: Visual state space model. *arXiv preprint arXiv:2401.10166*, 2024.
- [26] Yuen Peng Loh and Chee Seng Chan. Getting to know low-light images with the exclusively dark dataset. *Comput. Vis. Image Underst.*, 178:30–42, 2019.
- [27] Kin Gwn Lore, Adedotun Akintayo, and Soumik Sarkar. LLNet: A deep autoencoder approach to natural low-light image enhancement. *Pattern Recognit.*, 61:650–662, 2017.
- [28] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. *NeurIPS*, 29, 2016.
- [29] Feifan Lv, Feng Lu, Jianhua Wu, and Chongsoon Lim. MBLLEN: low-light image/video enhancement using cnns. In *BMVC*, page 220. BMVA Press, 2018.
- [30] Xiaoqian Lv, Shengping Zhang, Chenyang Wang, Weigang Zhang, Hongxun Yao, and Qingming Huang. Unsupervised low-light video enhancement with spatial-temporal co-attention transformer. *IEEE Trans. Image Process.*, 32:4701–4715, 2023.
- [31] Long Ma, Tengyu Ma, Risheng Liu, Xin Fan, and Zhongxuan Luo. Toward fast, flexible, and robust low-light image enhancement. In *CVPR*, pages 5627–5636, 2022.
- [32] Eric Nguyen, Karan Goel, Albert Gu, Gordon W. Downs, Preey Shah, Tri Dao, Stephen Baccus, and Christopher Ré. S4ND: modeling images and videos as multidimensional signals with state spaces. In *NeurIPS*, 2022.
- [33] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Köpf, Edward Z. Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, pages 8024–8035, 2019.
- [34] Xiaohuan Pei, Tao Huang, and Chang Xu. Efficientvmamba: Atrous selective scan for light weight visual mamba. *arXiv preprint arXiv:2403.09977*, 2024.
- [35] Biqing Qi, Junqi Gao, Dong Li, Kaiyan Zhang, Jianxing Liu, Ligang Wu, and Bowen Zhou. S4++: Elevating long sequence modeling with state memory reply, 2024.
- [36] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. *arXiv preprint arXiv:1804.02767*, 2018.
- [37] Shangquan Sun, Wenqi Ren, Jingyang Peng, Fenglong Song, and Xiaochun Cao. Di-retinex: Digital-imaging retinex theory for low-light image enhancement. *arXiv preprint arXiv:2404.03327*, 2024.
- [38] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017.
- [39] Ruixing Wang, Xiaogang Xu, Chi-Wing Fu, Jiangbo Lu, Bei Yu, and Jiaya Jia. Seeing dynamic scene in the dark: A high-quality video dataset with mechatronic alignment. In *ICCV*, pages 9680–9689, 2021.
- [40] Ruixing Wang, Qing Zhang, Chi-Wing Fu, Xiaoyong Shen, Wei-Shi Zheng, and Jiaya Jia. Underexposed photo enhancement using deep illumination estimation. In *CVPR*, pages 6849–6857, 2019.
- [41] Wenjing Wang, Huan Yang, Jianlong Fu, and Jiaying Liu. Zero-reference low-light enhancement via physical quadruple priors, 2024.
- [42] Zhou Wang, Alan C. Bovik, Hamid R. Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE Trans. Image Process.*, 13(4):600–612, 2004.
- [43] Chen Wei, Wenjing Wang, Wenhan Yang, and Jiaying Liu. Deep retinex decomposition for low-light enhancement. In *BMVC*, page 155, 2018.
- [44] Jan C Willems. From time series to linear system—part i. finite dimensional linear time invariant systems. *Automatica*, 22:561–580, 1986.
- [45] Wenhui Wu, Jian Weng, Pingping Zhang, Xu Wang, Wenhan Yang, and Jianmin Jiang. Uretinex-net: Retinex-based deep unfolding network for low-light image enhancement. In *CVPR*, pages 5891–5900. IEEE, 2022.
- [46] Xiaogang Xu, Ruixing Wang, Chi-Wing Fu, and Jiaya Jia. Snr-aware low-light image enhancement. In *CVPR*, pages 17693–17703. IEEE, 2022.
- [47] Shyi-Kae Yang and Chieh-Li Chen. Observer-based robust controller design for a linear system with time-varying perturbations. *Journal of Mathematical Analysis and applications*, 213(2):642–661, 1997.
- [48] Wenhan Yang, Wenjing Wang, Haofeng Huang, Shiqi Wang, and Jiaying Liu. Sparse gradient regularized deep retinex network for robust low-light image enhancement. *IEEE Trans. Image Process.*, 30:2072–2086, 2021.

- [49] Xunpeng Yi, Han Xu, Hao Zhang, Linfeng Tang, and Jiayi Ma. Diff-retinex: Rethinking low-light image enhancement with A generative diffusion model. In ICCV, pages 12268–12277. IEEE, 2023.
- [50] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In CVPR, pages 5718–5729. IEEE, 2022.
- [51] Syed Waqas Zamir, Aditya Arora, Salman H. Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Learning enriched features for real image restoration and enhancement. In ECCV.
- [52] Yonghua Zhang, Jiawan Zhang, and Xiaojie Guo. Kindling the darkness: A practical low-light image enhancer. In ACM Multimedia, pages 1632–1640. ACM, 2019.
- [53] Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. arXiv preprint arXiv:2401.09417, 2024.
- [54] Simiao Zuo, Xiaodong Liu, Jian Jiao, Denis Charles, Eren Manavoglu, Tuo Zhao, and Jianfeng Gao. Efficient long sequence modeling via state space augmented transformer. arXiv preprint arXiv:2212.08136, 2022.

A Broader Impact.

Low light image enhancement is the classical task that improves the quality of degraded images, exhibiting the promising value of research and application. Our proposed global-then-local state space enhances the feature extraction ability by integrating implicit Retinex-aware strategy. We believe our method has the potential to advance other low-level tasks and may inspire future research in state space models. However, there could be negative effects brought by the proposed method. For example, the inevitable deviations of training data distribution, the generated results for the real world scenarios may exist the color deviation.

B More Results.



Figure 7: More qualitative comparisons with SOTAs.(Zoom in for best view)

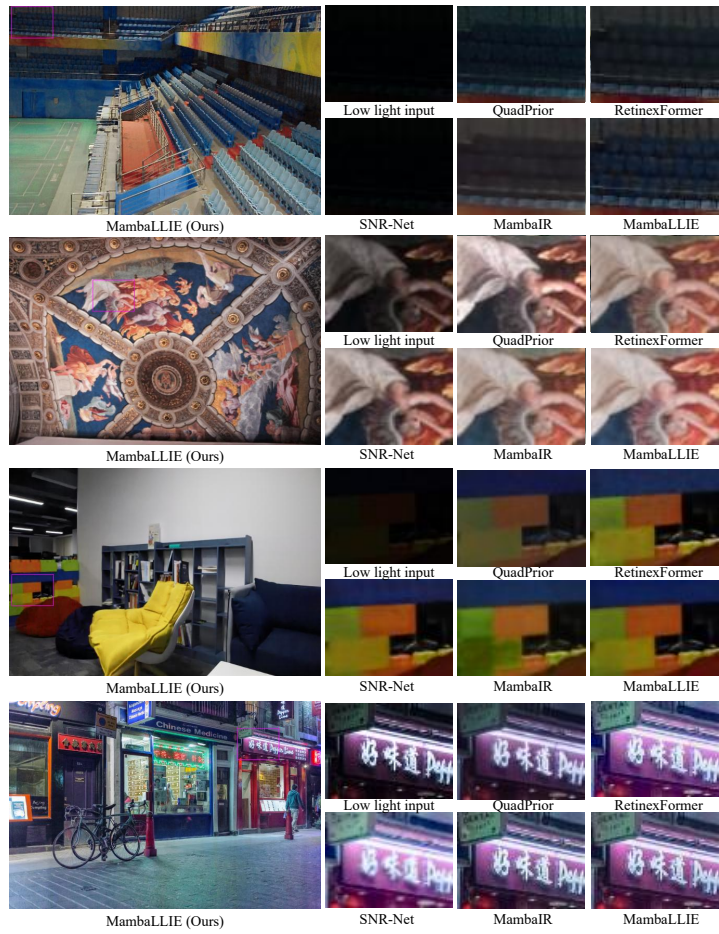


Figure 8: More qualitative comparisons with SOTAs.(Zoom in for best view)

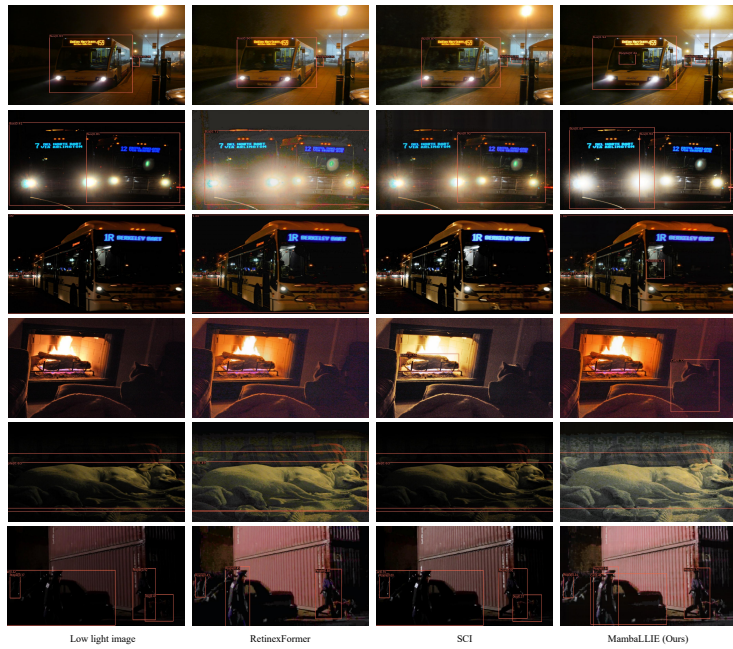


Figure 9: Object detection qualitative comparisons with SOTAs.(Zoom in for best view)