# SNNSIR: A Simple Spiking Neural Network for Stereo Image Restoration

Ronghua Xu, Jin Xie, Jing Nie, Jiale Cao, Yanwei Pang

*Abstract*—Spiking Neural Networks (SNNs), characterized by discrete binary activations, offer high computational efficiency and low energy consumption, making them well-suited for computation-intensive tasks such as stereo image restoration. In this work, we propose SNNSIR, a simple yet effective Spiking Neural Network for Stereo Image Restoration, specifically designed under the spike-driven paradigm where neurons transmit information through sparse, event-based binary spikes. In contrast to existing hybrid SNN-ANN models that still rely on operations such as floating-point matrix division or exponentiation, which are incompatible with the binary and event-driven nature of SNNs, our proposed SNNSIR adopts a fully spike-driven architecture to achieve low-power and hardware-friendly computation. To address the expressiveness limitations of binary spiking neurons, we first introduce a lightweight Spike Residual Basic Block (SRBB) to enhance information flow via spike-compatible residual learning. Building on this, the Spike Stereo Convolutional Modulation (SSCM) module introduces simplified nonlinearity through element-wise multiplication and highlights noise-sensitive regions via cross-view-aware modulation. Complementing this, the Spike Stereo Cross-Attention (SSCA) module further improves stereo correspondence by enabling efficient bidirectional feature interaction across views within a spike-compatible framework. Extensive experiments on diverse stereo image restoration tasks, including rain streak removal, raindrop removal, low-light enhancement, and super-resolution demonstrate that our model achieves competitive restoration performance while significantly reducing computational overhead. These results highlight the potential for real-time, low-power stereo vision applications. The code will be available after the article is accepted.

*Index Terms*—Spiking neural network, Spike-driven, Stereo image restoration.

## I. INTRODUCTION

IMAGE restoration aims to reconstruct high-quality images from degraded or low-quality inputs, including those affected by adverse weather (e.g., rain), insufficient lighting (e.g., low-light environments), or limited resolution. Stereo image restoration extends this task by leveraging a pair of degraded left and right images to produce more detailed and geometrically consistent results. One of its key advantages lies in the interactive information exchange between the two views, which allows the model to recover scene details that are degraded or missing in one image using cues from its counterpart. By exploiting such complementary information,
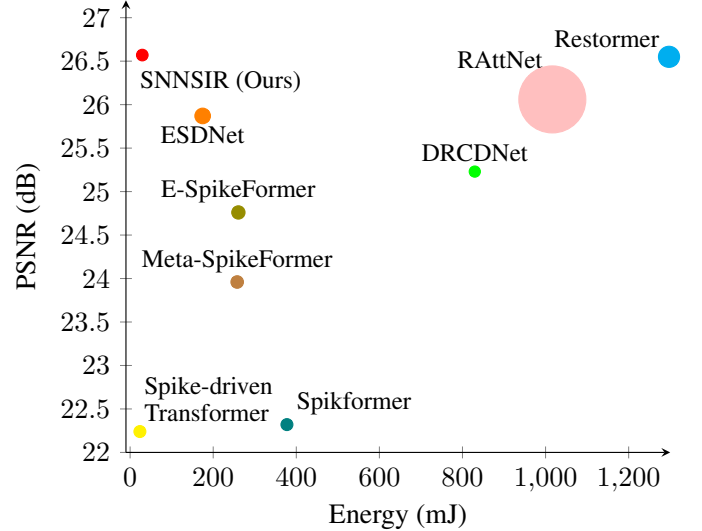


Fig. 1. Comparison of PSNR (dB) vs. Energy Consumption (mJ) on the SeteroWaterdrop test set. Circle size denotes model size. Our SNNSIR achieves competitive restoration performance with significantly lower energy consumption and model size.

stereo restoration often surpasses single-image approaches in both visual fidelity and structural coherence.

Additionally, this capability is especially critical in applications such as stereo matching [1], [2], depth estimation [3], and 3D perception tasks like object tracking [4] and detection [1], [2]. Moreover, high-quality stereo image restoration serves as a fundamental component in robot vision and embodied intelligence systems, where accurate spatial perception and robust visual input are essential for decision making, navigation, and interaction with dynamic environments.

Cross-view interaction lies at the core of stereo restoration tasks. Recent stereo restoration methods [5]–[11] have commonly enhanced such interaction through guided alignment, parallax attention, and semantic fusion, achieving notable improvements in restoration quality. However, these approaches typically rely on dual-branch encoders, and computationally intensive interaction modules, leading to considerable computational overhead. This complexity hinders real-time deployment on resource-limited platforms such as robots, UAVs [12], [13], and edge devices. Consequently, there is an increasing demand for lightweight stereo restoration frameworks that preserve the benefits of cross-view modeling while ensuring computational efficiency.

To address the growing need for computational efficiency in vision tasks, recent research has explored Spiking Neural Networks (SNNs) [14], often referred to as the third genera-
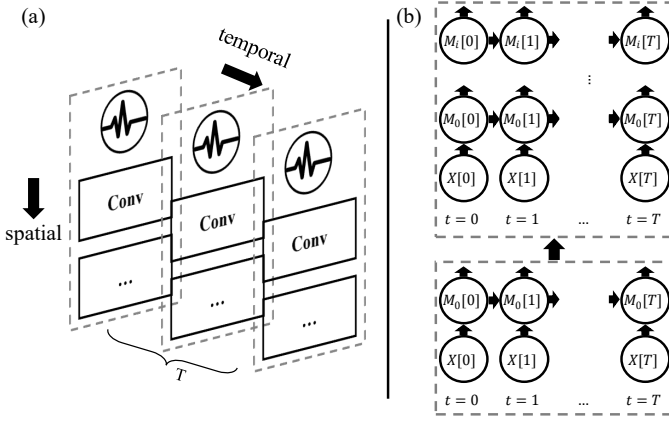
Fig. 2. (a) and (b) respectively illustrate the spatiotemporal characteristics of SNN in terms of network architecture and data propagation.

tion of neural networks. Unlike ANNs, SNNs operate using discrete spikes triggered by threshold-based membrane potentials, offering sparse, event-driven computation and inherent temporal modeling. These characteristics enable significant reductions in energy consumption and computational cost. SNNs have shown promising results in high-level tasks such as image classification [15]–[19], object detection [20], and semantic segmentation [21]. Notably, ESDNet [22] recently demonstrated the potential of SNNs in rain removal, achieving effective restoration with low energy consumption. These developments point to the potential of SNNs as a lightweight solution for stereo image restoration, particularly on resource-constrained platforms. Despite their strengths in temporal modeling and computational efficiency, SNNs have yet to be fully explored in this domain. Moreover, existing SNN-based single-image restoration methods, such as ESDNet [22], adopt hybrid SNN-ANN architectures that rely on non-spike-compatible operations, including floating-point division and sigmoid activations, thereby limiting their energy efficiency. For example, the Mixed Attention Unit (MAU) in ESDNet employs sigmoid functions, underscoring the lack of exploration into achieving nonlinear representation within fully spike-driven frameworks.

In this paper, we propose a simple spiking neural network for stereo image restoration, called SNNSIR. The architecture adopts a coarse-to-fine framework, enabling hierarchical refinement while maintaining a fully spike-driven processing pipeline. Unlike ANNs that depend on floating-point operations and nonlinear activations, SNNSIR achieves nonlinear representation and cross-view reasoning using spike-compatible modules.

To realize this design, we introduce several SNN-based modules. To enhance the representational capacity of binary spike neurons, we deisgn the Spike Residual Basic Block (SRBB), which integrates a residual learning scheme tailored for SNNs. Two SRBBs form a Feature Extraction Block (FEB) that benefits from the discrete computation nature to avoid excessive computational burden. To introduce nonlinearity into the SNN model, we draw inspiration from prior findings [23] showing that element-wise multiplication can serve as an

effective substitute for traditional activations. Based on this, we design the Spike Stereo Convolutional Modulation (SSCM) module, which employs multiplication as a spike-compatible nonlinear activation. Additionally, its channel- and spatial-wise modulation enhances feature representation and highlights degraded regions to facilitate subsequent restoration. To capture long-range dependencies and enhance cross-view interaction, we introduce the Spike Stereo Cross-Attention (SSCA) module, which facilitates efficient inter-view information exchange within a spike-driven framework. In addition, to compensate for the loss of fine-grained details in deeper layers, a coarse-to-fine architecture is adopted, where lightweight Spike Stereo Refinement Blocks (SSRBs) are introduced in the final stage for local restoration. Furthermore, to enable direct training of the SNN-based model, we adopt surrogate gradient techniques [24], which effectively address the non-differentiability of spike functions. Our model achieves comparable restoration performance to existing ANN- and SNN-based methods, while significantly reducing energy consumption.

The main contributions of this paper are as follows:

- We propose a simple spiking neural network for stereo image restoration, termed SNNSIR, offering a novel and energy-efficient perspective for adapting stereo restoration to real-world, resource-constrained devices.
- We design the Spike Stereo Cross-Attention (SSCA) module, Spike Stereo Convolutional Modulation (SSCM) module, and Spike Stereo Refinement Block (SSRB). SSCM introduces nonlinearity and guides the network to focus on degraded regions. SSCA enables efficient long-range cross-view interaction. SSRB refines local details in a lightweight manner.
- Extensive experiments on stereo image raindrop removal, rain streak removal, low-light enhancement, and super-resolution validate the effectiveness of our method. SNNSIR achieves comparable performance to representative ANN-based methods while reducing energy consumption by up to 97.73%, and consistently outperforms existing SNN-based approaches. To our knowledge, this work establishes the first dedicated and high-performing baseline for stereo image restoration using spiking neural networks.

## II. PRELIMINARY

**Spatiotemporal Nature.** SNNs are spatiotemporal. In Figure 2(a), we follow [25], where the spatial dimension is represented by the vertical stacking of convolutional layers (Conv) within each time step, capturing spatial dependencies in the input, and the temporal dimension is depicted by the sequential propagation of neural activities across different time steps $T$, emphasizing the dynamic nature of SNNs. Specifically, the feature propagates layer by layer along the temporal dimension which is represented in Figure 2(b), as described by [26], where $\mathbf{X}$ and $\mathbf{M}$ denote input features and intermediate result in each layer and time step with total $T$, separately Therefore, for static images, we first temporally replicate them over $T$ time steps to generate a sequence $\mathbf{X}[T]$, as shown in Figure 3. Considering the biological characteristics of Leaky Integrate-and-Fire (LIF) neuron [14] that can transform input features
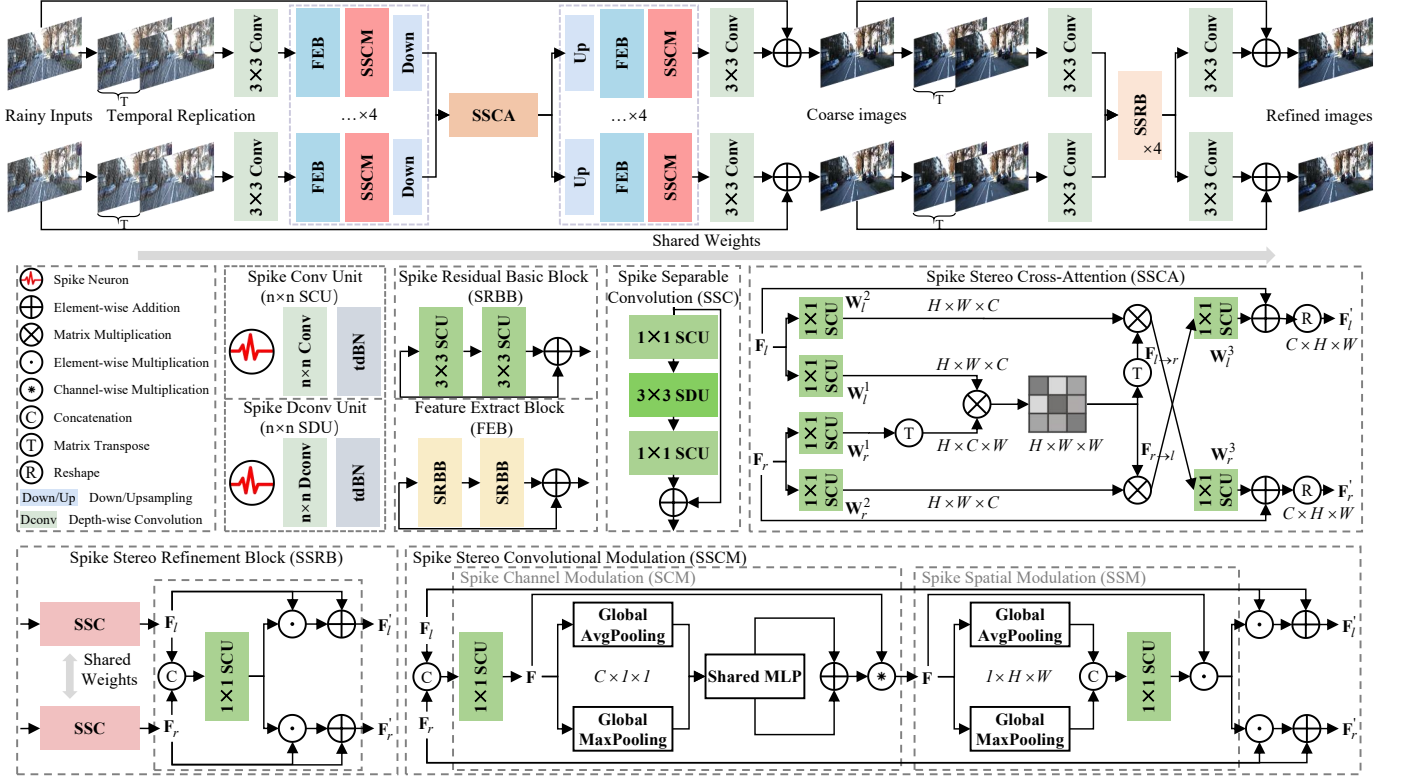
Fig. 3. The architecture of the proposed two-stage spike-driven stereo image restoration framework. In the first stage, a coarse restored image is generated, followed by a lightweight refinement stage that removes artifacts to produce the final output. The framework incorporates several core components: a Feature Extraction Block (FEB) with two Spike Residual Basic Blocks (SRBBs) for efficient feature extraction, a Spike Stereo Convolutional Modulation (SSCM) module to enhance nonlinearity, and a Spike Stereo Cross Attention (SSCA) module to facilitate cross-view interaction.

into binary (0/1) spike sequences and its ease of simulation on computers, we select it to implement the propagation of spike signals in the network. Its dynamic equation can be described as:

$$\mathbf{U}_i[t] = \mathbf{V}_i[t-1] + \frac{1}{\tau}(\mathbf{X}_i[t] - (\mathbf{V}_i[t-1] - u_{\text{rest}})), \quad (1)$$

$$\mathbf{S}_i[t] = \theta(\mathbf{U}_i[t] - u_{\text{th}}), \quad (2)$$

$$\theta(x) = \begin{cases} 1, & x \geq 0 \\ 0, & x < 0 \end{cases}, \quad (3)$$

$$\mathbf{V}_i[t] = (1 - \mathbf{S}_i[t])\mathbf{U}_i[t] + \mathbf{S}_i[t]u_{\text{rest}}, \quad (4)$$

where $t$ and $i$ represents the $t$-th time step and the $i$-th spike neuron. The whole procedure encompasses three distinct phases: charging, firing, and resetting. Equation (1) represents the charging progression, where the interaction between spatial input $\mathbf{X}_i[t]$ and the membrane potential $\mathbf{V}_i[t-1]$ at the last moment generates the current membrane potential $\mathbf{U}_i[t]$. $\tau$ is the membrane time constant. Equation (2) represents the firing progression. If the membrane potential after charging exceeds the threshold potential $u_{\text{th}}$, a spike is emitted through the step function in Equation (3). Otherwise, the output spike is 0. After firing, reset the membrane potential which is described in Equation (4). Here, $u_{\text{rest}}$ and $\mathbf{V}_i[t]$ represent the reset potential and the final potential of this neuron.

**Energy Consumption.** In ANNs, floating-point operations (FLOPs) are commonly used to evaluate the computational cost, where most of the FLOPs are multiply-accumulate operations (MACs). However, in SNNs, which handle binary matrices, the operations reduce to only accumulation operations (ACs), referred to as synaptic operations (SOPs) [16], [17]. We assume that the implementation of MAC and AC operations is carried out on 45nm hardware [27], with the energy consumption $E_{\text{MAC}} = 4.6pJ$ and $E_{\text{AC}} = 0.9pJ$. For an input $\mathbf{X}$, We follow [16], defining SOPs as follows:

$$\text{SOPs}(\mathbf{X}) = T \times fr \times \text{FLOPs}(\mathbf{X}), \quad (5)$$

where $T$ is the time step and $fr$ denotes the firing rate. Here, the product of $fr$ and $\text{FLOPs}(\mathbf{X})$ corresponds to the number of 1s in the feature matrix after processing by the spike neurons.

For our proposed network, to be more precise, we describe its energy consumption as follows:

$$E = 0.9pJ \times \sum \text{SOPs}(\mathbf{X}_s) + 4.6pJ \times \sum \text{FLOPs}(\mathbf{X}_a) \quad (6)$$

Where $\mathbf{X}_s$ is the binary spike matrix processed by the spike neurons. $\mathbf{X}_a$ is the floating-point matrix.

At the code level, calculating SOPs is equivalent to counting the number of 1s in the output from the spike neurons shown in Figure 2(a), which corresponds to the number of spikes emitted. Since the spike firing rate varies for each dataset, we first calculate the total number of spikes in the test set for each dataset, and then compute the ratio of the total spikes

to the number of samples. This provides us the SOPs per sample for that dataset, from which the energy consumption can be derived. As we all know, FLOPs are fixed and only need to account for the computation that occurs in non-spike sequences. Taking our proposed SNNSIR as an example, the only non-spike sequences are the $3 \times 3$ convolutional layers at the beginning and end of the degradation removal and refinement stages, while the rest of the operations are spike sequences. In this paper, SOPs and FLOPs are measured on $256 \times 256$, except for super-resolution, which uses a size of $64 \times 64$. The results of FLOPs and SOPs for different tasks on different datasets are shown in Table VIII. For tasks involving multiple datasets, the computational complexity is taken as the average across all datasets for that task.

## III. METHOD

The overall architecture of the proposed SNNSIR is illustrated in Figure 3. SNNSIR restores a degraded stereo image pair through a coarse-to-fine framework composed of a U-shaped encoder-decoder and a lightweight refinement stage.

The first stage adopts a 5-layer U-shaped encoder-decoder with channel dimensions [32, 64, 96, 128, 160], enabling multi-scale feature representation with a simple structure. Given a static stereo input of size $3 \times H \times W$, temporal replication produces inputs $\mathbf{X}_l, \mathbf{X}_r \in \mathbb{R}^{T \times 3 \times H \times W}$ for the left and right views. These replicated inputs are then passed through a $3 \times 3$ convolution to extract shallow features $\mathbf{F}_l, \mathbf{F}_r \in \mathbb{R}^{T \times C \times H \times W}$, serving as the initial representation for subsequent processing.

These features are processed by the Feature Extraction Block (FEB) for single-view feature extraction, followed by the Spike Stereo Convolutional Modulation (SSCM) module, which introduces nonlinearity and highlight noise-sensitive region. Feature maps are progressively downsampled, reducing spatial resolution while increasing channel depth. The encoded features are then fed into the Spike Stereo Cross-Attention (SSCA) module to model cross-view long-range dependencies. The decoder mirrors the encoder and includes skip connections via element-wise addition. After decoding, a $3 \times 3$ convolution and temporal average pooling generate the left and right residual maps $\mathbf{O}_l, \mathbf{O}_r \in \mathbb{R}^{3 \times H \times W}$, which are added to the original inputs to obtain coarse predictions.

Finally, to recover fine details, a refinement stage composed of four Spike Stereo Refinement Blocks (SSRBs) operates at a fixed resolution with 32 channels, avoiding information loss from further downsampling.

Next, we provide detailed descriptions of the core modules in SNNSIR.

**Spike Convolution Unit (SCU) and Spike Depth-wise Convolution Unit(SDU).** Given that SNNs process binary spike sequences, their foundational structure [15], [16], [18]–[20], [22], [28] typically consists of a LIF neuron [14] followed by a convolutional or linear layer. tdBN [28] performs normalization simultaneously in both spatial and temporal dimensions that can alleviate the gradient and firing rate issues in deep SNN training. [22] proposed a Spike Convolution Unit (SCU), but it is not specific enough. We design $n \times n$ SCU and $n \times n$

Spike Depth-wise Convolution Unit ($n \times n$ SDU) based on the kernel size and convolution type. These two structures are described as:

$$
\begin{aligned}
\mathbf{X}_{i+1}^{SCU}[t] &= \text{tdBN}(\text{Conv}(\mathbf{S}_i[t])), \\
\mathbf{X}_{i+1}^{SDU}[t] &= \text{tdBN}(\text{Dconv}(\mathbf{S}_i[t])),
\end{aligned}
\tag{7}
$$

where the definition of $\mathbf{S}$, $i$ and $t$ can be found in Section II.
**Spike Residual Basic Block (SRBB).** To ensure sufficient network depth while preserving the ability to extract rich local features, we adopt the residual basic block [29] as the fundamental building unit of our network. Unlike vanilla residual basic block, we replace the standard convolutional layers with $3 \times 3$ SCUs to better accommodate the characteristics of SNN. Compared to traditional ANNs, implementing residual connections in SNNs requires specialized designs. Among existing strategies [15], [25], [28], we adopt the Membrane Shortcut (MS) [25] as it best preserves identity mapping within the spike-driven paradigm, forming our Spike Residual Basic Block (SRBB). We construct each FEB by stacking two SRBBs to learn richer, hierarchical spatiotemporal representations. FEBs' weights are shared across the left and right stereo views to reduce the parameter count and enforce consistent feature extraction.

**Spike Stereo Convolutional Modulation (SSCM) module.** To address the limited non-linearity and lack of cross-view interaction in spiking neural networks, the proposed Spike Stereo Convolutional Modulation (SSCM) module introduces spike-compatible non-linear activation while enabling stereo-aware feature refinement through convolutional modulation. Inspired by [23], SSCM employs element-wise multiplication to introduce lightweight non-linearity without relying on conventional activation functions such as sigmoid, which are incompatible with the spike-driven paradigm.

As illustrated in Figure 3, SSCM takes the concatenated left and right view features as input to form a fused representation $\mathbf{F}$, which encodes coarse inter-view correlations and enables the network to leverage global context for compensating view-specific degradations caused by noise or occlusion. It then sequentially applies Spike Channel Modulation (SCM) and Spike Spatial Modulation (SSM), which modulate feature responses along the channel and spatial dimensions, respectively. Both submodules rely on element-wise multiplication to introduce non-linearity. The detailed computational process is defined as:

$$
\begin{aligned}
\text{SCM}(\mathbf{F}) &= \mathbf{F} \circledast (\mathbf{W}(\text{GAP}(\mathbf{F})) + \mathbf{W}(\text{GMP}(\mathbf{F}))), \\
\text{SSM}(\mathbf{F}) &= \mathbf{F} \odot (\text{SCU}([\text{GAP}(\mathbf{F}), \text{GMP}(\mathbf{F})])), \\
\mathbf{F}_{l/r}^{'} &= \text{SSM}(\mathbf{F}) \odot \mathbf{F}_{l/r} + \mathbf{F}_{l/r},
\end{aligned}
\tag{8}
$$

where $\text{GAP}(\cdot)$ and $\text{GMP}(\cdot)$ denote global average pooling and global max pooling, respectively; $\mathbf{W}(\cdot)$ is a shared linear operation; $\circledast$ and $\odot$ represent channel- and spatial-wise element-wise multiplication, respectively. The final outputs, $\mathbf{F}_l^{'}$ and $\mathbf{F}_r^{'}$, are obtained through residual modulation based on spike-compatible operations. In summary, SSCM acts as a unified component that simultaneously enhances non-linearity and facilitates cross-view information integration, playing a

TABLE I

COMPARISON FOR STEREO IMAGE RAINDROP REMOVAL ON THE STEREO WATERDROP DATASET. P AND E INDICATE THE NUMBER OF PARAMETERS AND ENERGY CONSUMPTION (mJ), RESPECTIVELY.

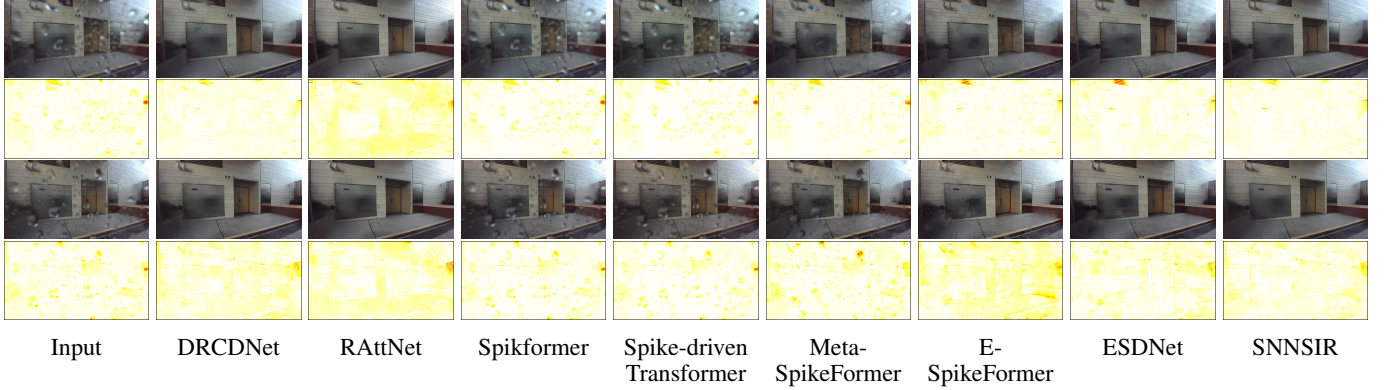| Methods | Type | PSNR (dB) ↑ | MS-SSIM ↑ | LPIPS ↓ | P (M) ↓ | E (mJ) ↓ | Times (s) ↓ |
|---|---|---|---|---|---|---|---|
| Pix2Pix | | 22.76 | 0.895 | 0.217 | 54.41 | - | **0.045** |
| Qian et al. | | 24.47 | 0.900 | 0.163 | 6.24 | 822.39 | 0.083 |
| Liu et al. | | 22.70 | 0.833 | 0.247 | 16.69 | - | 1.241 |
| Quan et al. | ANN | 24.97 | 0.913 | 0.153 | 7.27 | - | 0.115 |
| DRCDNet | | 25.23 | 0.906 | 0.176 | **2.25** | 829.64 | 0.173 |
| Restormer | | **26.55** | 0.945 | 0.124 | 26.13 | 1297.11 | 0.439 |
| RAttNet | | 26.06 | **0.950** | 0.096 | 136.55 | 1016.30 | 0.210 |
| Spikformer | | 22.32 | 0.828 | 0.186 | 4.04 | 377.45 | 0.680 |
| Spike-driven Transformer | | 22.24 | 0.830 | 0.182 | 4.04 | **23.53** | 0.642 |
| Meta-SpikeFormer | SNN | 23.96 | 0.874 | 0.151 | 5.54 | 257.74 | 0.663 |
| E-SpikeFormer | | 24.76 | 0.918 | 0.102 | 7.30 | 260.61 | **0.032** |
| ESDNet | | 25.87 | 0.941 | **0.072** | 12.81 | 174.63 | 16.148 |
| SNNSIR(Ours) | | **26.57** | **0.949** | **0.062** | **3.26** | **29.32** | 0.411 |



Fig. 4. Visual comparison for the stereo image raindrop removal task on the Stereowaterdrop dataset. The top and bottom rows show the left and right views, respectively. Best viewed by zooming in for details.

key role in high-quality stereo restoration under the constraints of biologically inspired spiking computation.

**Spike Stereo Cross-Attention (SSCA) module.** The Spike Stereo Cross-Attention (SSCA) module, a SNN variant of the Stereo Cross-Attention Module [8], is designed to enhancing feature interaction between the left and right images provided by stereo vision systems. Its key modification lies in the introduction of the SCU and the removal of the activation function, making it spike-compatible. Due to the characteristics of stereo vision, horizontal disparities are present while there is no significant disparity in the vertical direction. As a result, performing attention computation in the horizontal direction is more efficient.

$$\mathbf{F}'_l = \mathbf{W}^3_l(\mathbf{W}^1_l\mathbf{F}_l \times (\mathbf{W}^1_r\mathbf{F}_r)^T \times \mathbf{W}^2_r\mathbf{F}_r) + \mathbf{F}_l,$$
$$\mathbf{F}'_r = \mathbf{W}^3_r((\mathbf{W}^1_l\mathbf{F}_l \times (\mathbf{W}^1_r\mathbf{F}_r)^T)^T \times \mathbf{W}^2_l\mathbf{F}_l) + \mathbf{F}_r. \quad (9)$$

In Equation (9), feature maps $\mathbf{F}_l$ and $\mathbf{F}_r$ are first reshaped into $H \times W \times C$ dimensions. The reshaped features are then processed by the learned weight matrices, $\mathbf{W}^1_l, \mathbf{W}^2_l, \mathbf{W}^1_r, \mathbf{W}^2_r, \mathbf{W}^3_l$ and $\mathbf{W}^3_r$ from $1 \times 1$ SCU. The output features $\mathbf{F}'_l$ and $\mathbf{F}'_r$ are refined by combining the spatial information from both views using the cross-attention computation and reshaped back to the $C \times H \times W$ dimensions.

**Spike Stereo Refinement Block (SSRB).** To compensate for the loss of fine details in the coarse output, we introduce a refinement stage that operates at full resolution to enhance textures and residual information without further downsampling. In Figure3, we design the SSRB, which consists of a Spike Separable Convolution (SSC) inspired by [19], adapted from the inverted separable convolution [30] for efficient spatial feature extraction. As shown in the Figure 3, after being handled by the SSC, the features from both views are first concatenated and processed by a SCU, whose output is then used to modulate each branch via element-wise multiplication and Membrane Shortcut [25], enhancing discriminative feature refinement in a lightweight manner while providing certain inter-view interaction and non-linearity.

### A. Loss Functions

Due to the distinct tasks undertaken by the two stages of the network, different loss functions are employed for each stage. The $L_1$ loss emphasizes pixel-level accuracy, making it suitable for the first stage focused on degradation removal. Perception loss [31], [32], which calculates the loss using features extracted from multiple layers to align both local high-frequency details and global low-frequency structures, is

employed in the refinement stage and denoted as $L_p$. $L_1$ and $L_p$ are computed as:

$$L_1 = \frac{1}{N} \sum_{i=1}^{N} |\mathbf{X}_{l,i} - \mathbf{O}_{l,i}| + \frac{1}{N} \sum_{i=1}^{N} |\mathbf{X}_{r,i} - \mathbf{O}_{r,i}|,$$

$$L_p = \sum_k \lambda_k (\|\mathbf{X}_l^k - \mathbf{O}_r^k\|_1 + \|\mathbf{X}_r^k - \mathbf{O}_r^k\|_1), \tag{10}$$

where $\mathbf{X}$ and $\mathbf{O}$ are input and output, $\mathbf{X}_{l,i}$ means the $i$-th pixel with a total number of N in the left view and similarly for the right view, $k$ denotes the selected layers, $\| \cdot \|_1$ denotes the $l_1$ norm, and $\lambda_k$ functions as a hyperparameter to adjust the contribution of each layer. Thus, the overall loss is the sum of the two losses, expressed as $L_{total} = L_1 + L_p$.

## IV. EXPERIMENTS

We conducted extensive analyses on three common tasks: **(a)** stereo image deraining, **(b)** stereo image low-light enhancement, and **(c)** stereo image super-resolution. Details of the datasets, training hyperparameters for each task, and more experimental results are provided in the supplementary material. In the result tables of this section, the best and second-best results are both highlighted in bold.

### A. Implementation Details

Our proposed network adopts a coarse-to-fine architecture, where the degradation removal stage uses channel dimensions of $[32, 64, 96, 128, 160]$ and in the refinement stage, four SSRBs are employed with an embedding dimension of 32. AdamW optimizer ($\beta_1 = 0.9$, $\beta_2 = 0.99$, weight decay $= 1e-4$) is used with a fixed learning rate $1e-3$. Furthermore, since low-light enhancement degradation is not high-frequency in nature, we set its spike activation threshold to 0.1, while 0.2 for other tasks. The default number of time steps $T$ for all tasks is 4. Due to the significant differences across datasets for each task, other hyperparameters vary considerably. Detailed configurations are provided in the supplementary material. Our network is implemented via SpikingJelly [26], a PyTorch-based SNN framework. All experiments are carried out on an NVIDIA GeForce RTX 3090 (24GB, 350W).

**Evaluation Metrics and Visualization.** To evaluate restoration performance, we adopt PSNR, SSIM, MS-SSIM [33], and LPIPS [34] as evaluation metrics. The number of parameters (P), computational complexity, and energy consumption (E) are used to evaluate the computational cost and resource overhead. Specifically, the computational complexity includes FLOPs and SOPs. The introduction to these two metrics can be found in Section II. In this paper, FLOPs and SOPs are measured on an image size $256 \times 256$, except for the super-resolution task, where a size of $64 \times 64$ is used. For visual comparison, we use error maps that intuitively depict the pixel-wise differences between the restored images and the ground-truth images through a red-yellow-white color gradient. In these error maps, larger errors correspond to more intense red hues, smaller errors appear closer to white, and moderate errors are represented by yellow.

**Comparison Methods.** For the stereo raindrop removal task, ANN-based methods including DRCDNet [35], Restormer

### TABLE II
COMPARISON FOR STEREO IMAGE RAIN STREAK REMOVAL ON THE RAINKT12 AND RAINKT15 DATASETS.

| Methods | P (M) | E (mJ) | RainKT12 | RainKT15 |
|---|---|---|---|---|
| DID-MDN | **0.47** | 93.23 | 29.14/0.901 | 28.97/0.899 |
| DeHRain | 50.33 | 2486.78 | **31.02/0.923** | **30.84/0.921** |
| Spikformer | 4.04 | 1255.57 | 18.73/0.600 | 22.28/0.662 |
| Spike-driven-Transformer | 4.04 | **41.70** | 23.46/0.703 | 24.32/0.728 |
| Meta-SpikeFormer | 5.54 | 201.45 | 26.78/0.846 | 27.12/0.847 |
| E-SpikeFormer | 7.30 | 281.91 | 21.65/0.709 | 21.91/0.713 |
| ESDNet | 12.81 | 217.13 | 28.52 /0.887 | 28.80/0.884 |
| SNNSIR | **3.26** | **32.69** | **30.97/0.922** | **31.43/0.920** |

### TABLE III
COMPARISON FOR STEREO IMAGE LOW-LIGHT ENHANCEMENT ON THE MIDDLEBURY AND HOLOPIX50K DATASETS.

| Methods | P (M) | E (mJ) | Middlebury | Holopix50k |
|---|---|---|---|---|
| ZeroDCE | **0.08** | 98.53 | 15.43/0.715 | 13.28/0.652 |
| RetinexNet | 0.84 | 325.73 | 22.66 0.800 | 19.20/0.779 |
| DRBN | **0.58** | 173.83 | **31.02/0.943** | **25.09/0.903** |
| Spikformer | 4.04 | 486.19 | 21.99/0.770 | 17.49/0.678 |
| Spike-driven-Transformer | 4.04 | **48.95** | 26.84/0.800 | 20.54/0.787 |
| Meta-SpikeFormer | 5.54 | 301.30 | 22.47/0.792 | 19.54/0.787 |
| ESDNet | 12.81 | 208.55 | 30.19/**0.928** | 23.75/0.883 |
| SNNSIR | 3.26 | **36.48** | **31.28**/0.923 | **24.82/0.888** |

[36], and RAttNet [37] are retrained based on official code, while results of Pix2Pix [38], Qian et al. [39], Liu et al. [40] and Quan et al. [41] are referenced from [10]. Given the scarcity of SNN-based restoration methods, Spikformer [16], Spike-driven Transformer [17], Meta-SpikeFormer [19], and E-SpikeFormer [21] are adapted from state-of-the-art SNN-based classification methods. Each method follows a 5-layer U-shaped encoder-decoder structure consistent with SNNSIR. Each layer uses the core block from its original structure, with the channel numbers adjusted to match those of SNNSIR, and the remaining parameters of each block are accordingly modified. And ESDNet [22] is a dedicated monocular restoration method. These five SNN-based methods are retrained on all tasks. Results of ANN-based methods including DID-MDN [42], DeHRain [43], ZeroDCE [44], RetinexNet [45] and DRBN [46] for stereo image rain streak removal and low-light enhancement are respectively referenced from [10] and [47]. Notably, except for RAttNet and our SNNSIR, all other methods are monocular. To be fair, monocular methods are trained separately on the left and right views. Apart from parameter sharing, all evaluation metrics are computed in the same manner as those for stereo methods. Specifically, the restoration quality is averaged across the left and right views, while the computational resource consumption is the sum of the resources required for both views.

### B. Stereo Image Deraining Results

Table I presents the results of stereo image raindrop removal on the Stereo Waterdrop dataset [37], while Table II presents
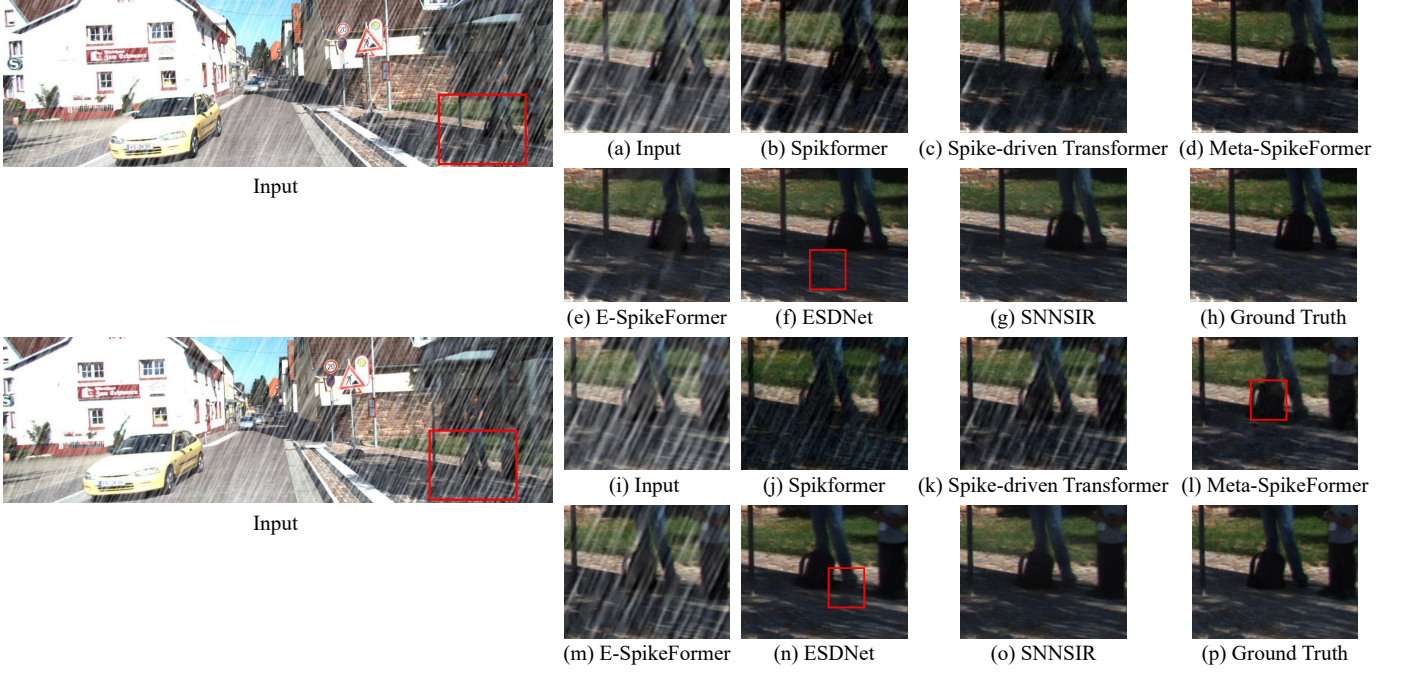
Fig. 5. Visual comparison for rain streak removal on the RainKT15 dataset. The top and bottom rows show the left and right views, respectively. Best viewed by zooming in for details.
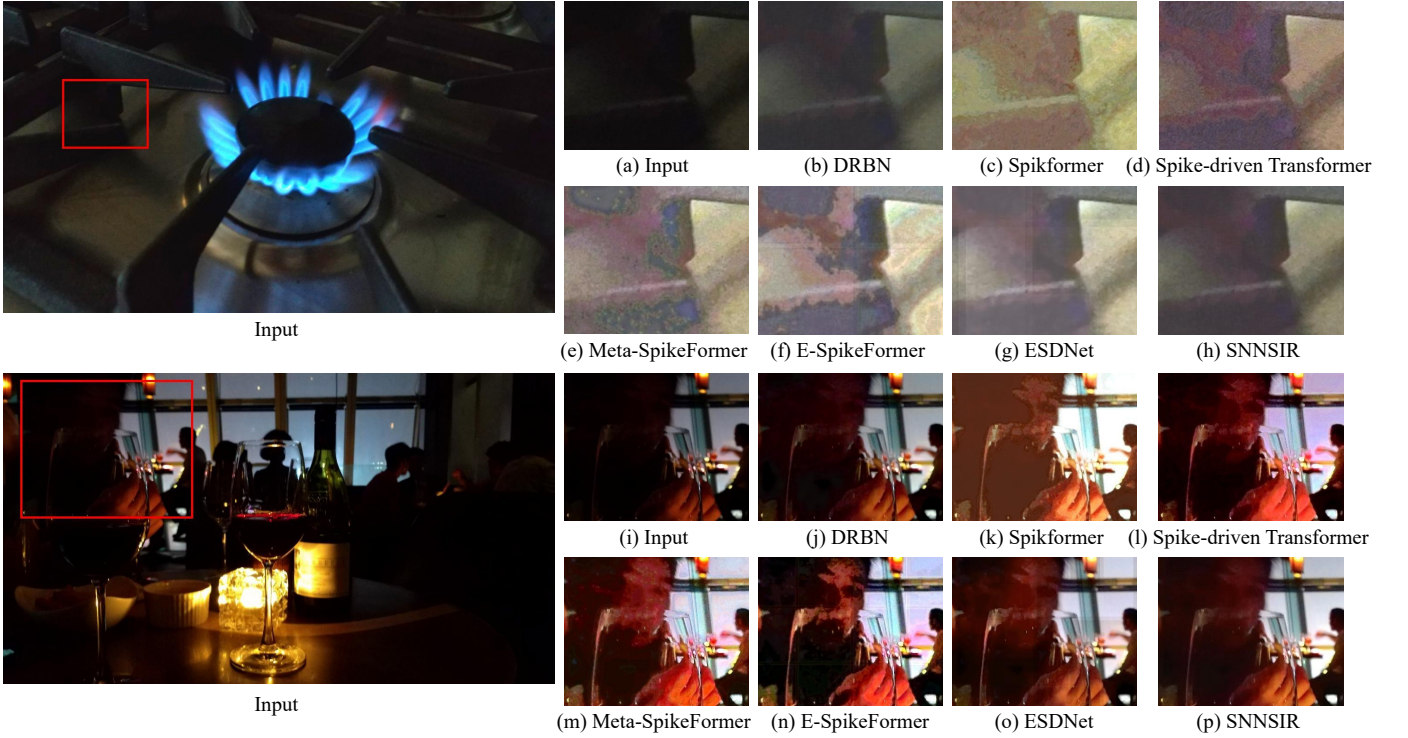


Fig. 6. Visual comparison for low-light enhancement on the test real from the Holopix50k dataset. Best viewed by zooming in for details.

the results of rain streak removal on the RainKT12 and RainKT15 datasets [9]. Compared with ANN-based methods, SNNSIR outperforms the classical methods Restormer by 0.02 dB and RAttNet by 0.51 dB. Although its restoration performance differs negligibly from the former, it reduces the P by 87.52% and E by 97.73%. This demonstrates our

proposed methods effectively leverage the spike-driven characteristics of SNNs, significantly reducing energy consumption, which shows that our methods offer a fresh perspective for image restoration. Compared with SNN-based methods, SNNSIR significantly outperforms the other five baselines in both restoration performance and resource load across the
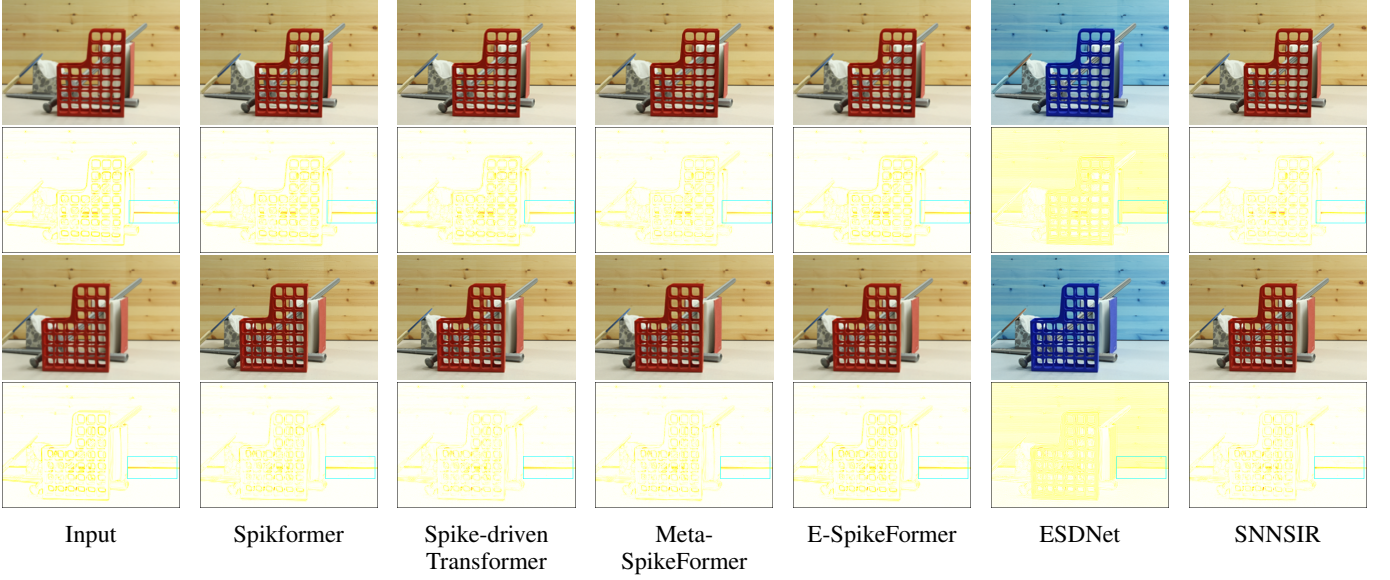
Fig. 7. Visual comparison for super-resolution on the Middlebury dataset. The top and bottom rows show the left and right views, respectively. Best viewed by zooming in for details.
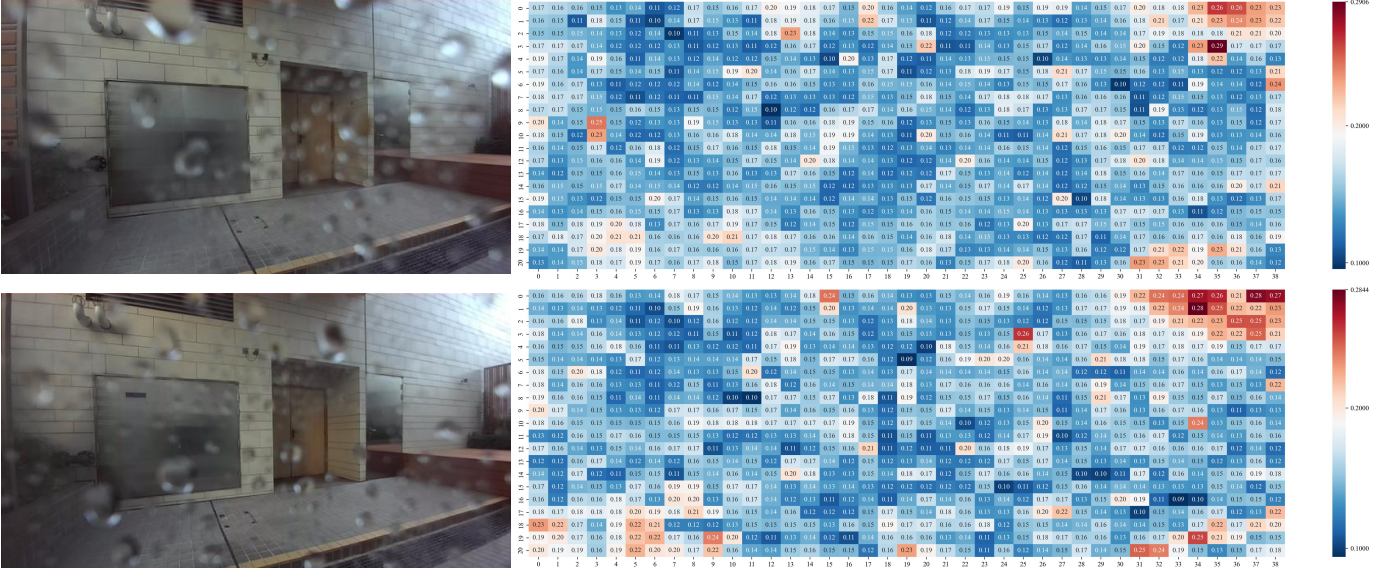


Fig. 8. Spike firing rate maps of SNNSIR on the Stereo Waterdrop dataset. The top and bottom rows show the left and right views, respectively. The blue-white-red gradient represents a gradual increase in the firing rate.

TABLE IV
COMPARISON FOR STEREO IMAGE SUPER-RESOLUTION ON THE
MIDDLEBURY AND FLICKR1024 DATASETS.

| Methods | P (M) | E (mJ) | Middlebury | Flickr1024 |
|---|---|---|---|---|
| Spikformer | **0.12** | 8.27 | 26.59/0.747 | 21.91/0.614 |
| Spike-driven-Transformer | **0.12** | **1.31** | 26.83/**0.757** | **22.06/0.627** |
| Meta-SpikeFormer | 0.39 | 17.34 | **26.90**/0.755 | 22.03/0.621 |
| E-SpikeFormer | 0.45 | 50.77 | 26.51/0.745 | 21.88/0.612 |
| ESDNet | 1.41 | 20.92 | 15.36/0.663 | 15.63/0.547 |
| SNNSIR | 0.33 | **1.59** | **27.38/0.772** | **22.29/0.640** |

two subtasks. Among these baselines, ESDNet, a monocular-

specific restoration model, lacks advantages when handling binocular tasks. Furthermore, our proposed method addresses the absence of high-quality SNN-based baselines in binocular tasks via the reasonable use of non-linearity and cross-view interaction. Figure 4 shows that SNNSIR exhibits no prominent large-area red regions, with the overall color leaning more toward white, indicating superior performance in both detail and global restoration. In Figure 5, among the results comparable to our proposed SNNSIR, we have highlighted the regions with poor restoration quality using red boxes. For ESDNet, noticeable vertical stripe artifacts appear in both the left and right views. As for Figure 5(l), Meta-SpikeFormer fails to fully restore the black backpack region. The above issues do not occur in SNNSIR.
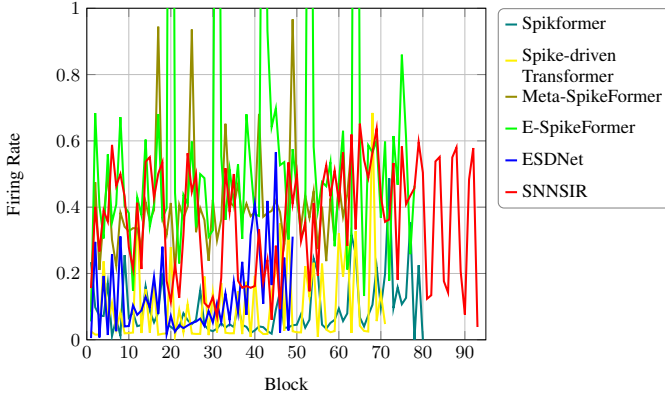
Fig. 9. Spike firing rate of SNN-based baselines.

## C. Stereo Image Low-light Enhancement Results

Experiments are conducted on the Middlebury [48] and Holopix50k [49] datasets. In Table III, SNNSIR outperforms SNN-based methods in terms of P, E and restoration quality. Specifically, it achieves superior PSNR and SSIM values with significantly lower energy consumption. Furthermore, even with a higher number of parameters, SNNSIR achieves performance comparable to DRBN while consuming only 20.99% of its energy, demonstrating exceptional energy efficiency without compromising restoration quality. To evaluate the generalization ability, we measure the results on the real-world test set from the Holopix50k dataset. The enhancement results of the five SNN-based baselines exhibit overexposure or white artifacts, which are particularly noticeable on the foreheads of individuals in Figures 6(k–o). In contrast, the results of SNNSIR appear more natural and visually consistent.

## D. Stereo Image Super-Resolution Results

We evaluate super-resolution (SR) methods on the Middlebury [48] and Flickr1024 [5] datasets. In line with the lightweight design philosophy of many SR methods, we adapt the five SNN baselines by removing up/downsampling and refinement operations, while maintaining a 32-channel width. The results in Table IV indicate ESDNet struggles with this task. While other four baselines perform comparably, their quality is limited. Our SNNSIR surpasses them, achieving the highest restoration quality with significantly lower energy usage. In terms of visual comparison, SNNSIR yields results that are generally similar to those of several existing methods. However, the regions highlighted by blue boxes demonstrate that SNNSIR achieves superior detail preservation, as shown in Figure 7. Notably, these baselines are not specifically designed for stereo SR, placing them at an inherent disadvantage. This reflects the current landscape of SNN-based stereo methods, establishing SNNSIR as the first tailored baseline for this task.

## E. Spike Firing Rate (SFR) Map

We define the SFR as the average spike firing rate along the $T$-dimension of the spike neuron outputs in the network. Since each unit block in the network contains a single spike event, we use the block as the recording point. Through SFR,

we can analyze the network's activation stability and spike distribution. Figure 9 illustrates that all SNN-based methods exhibit relatively low stability, yet they display a cyclical variation between high and low SFR, with SNNSIR showing this trend most prominently. This may be related to the activation and reset operations in SNNs, where high and low fluctuations in activation are expected. More appropriate SFR and more regular fluctuations contribute to better restoration quality and extremely low energy consumption. As shown in Figure 8, the spike distribution in SNNSIR is mostly concentrated in the blue, low-activation regions, except for the top-right corner, which is influenced by highlights. This unexpected issue may suggest that SNNs are less suitable for low-light enhancement tasks, as low-light conditions hinder activation. The scattered distribution of high-activation regions (white and red) corresponds to the isolated raindrop points, indicating that SNNSIR can effectively focus on the raindrops.

## F. Ablation Studies

For convenience, all ablation experiments are carried out on the Stereo Waterdrop dataset, which contains a suitable number of samples and diverse scenes.

**Impact of Different Designs.** Table V shows the effect of the key components. The Membrane Shortcut (MS) [25] used in SRBB brings better results to model (b) compared to model (a) [15]. Similarly, a comparison between model (d) and model (e) demonstrates that SSRB yields an improvement of 0.68 dB in PSNR while introducing only 0.03 M additional parameters. Furthermore, compared to model (b), model (c) achieves a significant performance improvement, which demonstrates the effectiveness of the non-linearity introduced by SSCM after removing the activation function from the SNNs. We observe that model (f) results in a performance loss, indicating that activation functions involving exponential operations are contrary to the characteristics of spike-driven.

**Analysis of the Effect of SNN-Based Design.** To assess whether our SNN-based design alleviates the representational limitations of binary activations, we revert SNNSIR to an ANN-based version by replacing spike neurons with ReLU and adding activation functions to SSCM and SSCA. As shown in Table VI, the SNN-based model achieves both lower energy consumption and better performance. While image restoration is inherently a static task, the temporal evolution in SNNs facilitates iterative feature refinement and selective activation, contributing to more effective noise suppression and structural recovery.

**Impact of Different Time Steps.** Table VII shows the effect of varying time steps on network performance. As the time step $T$ increases, performance initially improves, peaking at $T = 4$, after which further increases lead to degradation. This reflects the inherent characteristic of neural networks, where increasing computation and parameters does not always yield gains, and aligns with the biological nature of SNNs, where excessive load may induce neural fatigue. Clearly, by comparing with the performance at $T = 1$, it is evident that the performance improvement comes from the temporal dynamics of SNNs. However, even so, the theoretical increase

TABLE V
ABLATION STUDIES ON THE PROPOSED COMPONENTS. $\sigma$ IS THE SIGMOID FUNCTION.

| Model | P (M) | PSNR | SSIM |
|---|---|---|---|
| (a) SEW-RBB | **2.96** | 23.03 | 0.824 |
| (b) SRBB | **2.96** | 24.96 | 0.862 |
| (c) SRBB+SSCM | 3.10 | 25.91 | 0.877 |
| (d) SRBB+SSCA | 3.09 | 25.10 | 0.862 |
| (e) SRBB+SSCA+SSRB | 3.12 | 25.78 | 0.891 |
| (f) SRBB+SSCM ($\sigma$)+SSCA+SSRB | 3.26 | **26.30** | **0.899** |
| (g) SRBB+SSCM+SSCA+SSRB | 3.26 | **26.57** | **0.903** |

TABLE VI
COMPARISON OF SNN-BASED AND ANN-BASED DESIGNS OF THE PROPOSED SNNSIR.

| Methods | P (M) | PSNR | SSIM | E (mJ) |
|---|---|---|---|---|
| SNNSIR-ANN | 3.26 | 26.48 | 0.900 | 178.98 |
| SNNSIR-SNN | 3.26 | 26.57 | 0.903 | 29.32 |

TABLE VII
THE IMPACTS OF DIFFERENCE TIME STEPS.

| Time steps | OPs (G) | E (mJ) | PSNR | SSIM |
|---|---|---|---|---|
| 1 | **3.008** | **4.38** | 21.84 | 0.821 |
| 2 | **14.795** | **15.83** | 25.94 | 0.891 |
| 4 | 27.925 | 29.32 | **26.57** | **0.903** |
| 8 | 56.660 | 58.54 | **26.42** | **0.901** |

in computational cost by a factor of T does not occur due to the spike-driven nature of SNNs.

## V. CONCLUSION

In this paper, we propose SNNSIR, a novel spiking neural network for stereo image restoration. Specifically, we address the absence of nonlinearity in SNNs through a structured modulation mechanism that simultaneously highlights degraded regions. To capture cross-view dependencies, we design a spike-compatible interaction strategy, while a lightweight refinement module is employed to recover fine-grained details. Experimental results on multiple stereo restoration tasks demonstrate that our model achieves competitive performance with significantly lower energy consumption, establishing a strong baseline for future research in SNN-based stereo image restoration.

TABLE VIII
FLOPs AND SOPs OF DIFFERENT SNN-BASED METHODS FOR DIFFERENT STEREO IMAGE RESTORATION TASKS.

| Tasks | Methods | FLOPs (G) | SOPs (G) | E (mJ) |
|---|---|---|---|---|
| Raindrop removal | Spikformer | 0.566 | 416.498 | 377.45 |
| | Spike-driven-Transformer | 0.566 | 23.256 | 23.53 |
| | Meta-SpikeFormer | 0.566 | 283.480 | 257.74 |
| | E-SpikeFormer | 0.566 | 286.676 | 260.61 |
| | ESDNet | 19.284 | 95.469 | 174.63 |
| | SNNSIR | 1.132 | 26.793 | 29.32 |
| | SNNSIR-ANN | 38.908 | 0 | 178.98 |
| Rain streak | Spikformer | 0.566 | 1392.187 | 1255.57 |
| | Spike-driven-Transformer | 0.566 | 43.435 | 41.70 |
| | Meta-SpikeFormer | 0.566 | 220.940 | 201.45 |
| | E-SpikeFormer | 0.566 | 310.342 | 281.91 |
| | ESDNet | 19.284 | 142.697 | 217.13 |
| | SNNSIR | 1.132 | 30.532 | 32.69 |
| Low-light enhancement | Spikformer | 0.566 | 537.319 | 486.19 |
| | Spike-driven-Transformer | 0.566 | 51.495 | 48.95 |
| | Meta-SpikeFormer | 0.566 | 331.878 | 301.30 |
| | E-SpikeFormer | 0.566 | 498.600 | 451.34 |
| | ESDNet | 19.284 | 133.155 | 208.55 |
| | SNNSIR | 1.132 | 34.744 | 36.48 |
| Super-resolution | Spikformer | 0.046 | 8.956 | 8.27 |
| | Spike-driven-Transformer | 0.046 | 1.220 | 1.31 |
| | Meta-SpikeFormer | 0.046 | 19.031 | 17.34 |
| | E-SpikeFormer | 0.046 | 56.175 | 50.77 |
| | ESDNet | 0.385 | 21.277 | 20.92 |
| | SNNSIR | 0.046 | 1.530 | 1.59 |

## REFERENCES

[1] Y.-N. Chen, H. Dai, and Y. Ding, "Pseudo-stereo for monocular 3d object detection in autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 887–897.

[2] W. Wu, H. S. Wong, and S. Wu, "Semi-supervised stereo-based 3d object detection via cross-view consensus," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17 471–17 481.

[3] S. Brucker, S. Walz, M. Bijelic, and F. Heide, "Cross-spectral gated-rgb stereo depth estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 21 654–21 665.

[4] P. Li, J. Shi, and S. Shen, "Joint spatial-temporal optimization for stereo 3d object tracking," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 6877–6886.

[5] L. Wang, Y. Wang, Z. Liang, Z. Lin, J. Yang, W. An, and Y. Guo, "Learning parallax attention for stereo image super-resolution," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 12 250–12 259.

[6] B. Yan, C. Ma, B. Bare, W. Tan, and S. C. Hoi, "Disparity-aware domain adaptation in stereo image restoration," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 13 179–13 187.

[7] Y. Wang, X. Ying, L. Wang, J. Yang, W. An, and Y. Guo, "Symmetric parallax attention for stereo image super-resolution," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 766–775.

[8] X. Chu, L. Chen, and W. Yu, "Nafssr: Stereo image super-resolution using nafnet," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 1239–1248.

[9] K. Zhang, W. Luo, Y. Yu, W. Ren, F. Zhao, C. Li, L. Ma, W. Liu, and H. Li, "Beyond monocular deraining: Parallel stereo deraining network via semantic prior," *International Journal of Computer Vision*, vol. 130, no. 7, pp. 1754–1769, 2022.

[10] J. Nie, J. Xie, J. Cao, and Y. Pang, "Context and detail interaction network for stereo rain streak and raindrop removal," *Neural Networks*, vol. 166, pp. 215–224, 2023.

[11] C. Wang, T. Yan, W. Huang, X. Chen, K. Xu, and X. Chang, "Apanet: Asymmetrical parallax attention network for efficient stereo image deraining," *IEEE Transactions on Computational Imaging*, 2025.

[12] W. Chang, H. Chen, X. He, X. Chen, and L. Shen, "Uav-rain1k: A benchmark for raindrop removal from uav aerial imagery," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 15–22.

[13] C. Feng, Z. Chen, R. Kou, G. Gao, C. Wang, X. Li, X. Shu, Y. Dai, Q. Fu, and J. Yang, "Hazydet: Open-source benchmark for drone-view object detection with depth-cues in hazy scenes," *arXiv preprint arXiv:2409.19833*, 2024.

[14] W. Maass, "Networks of spiking neurons: the third generation of neural network models," *Neural networks*, vol. 10, no. 9, pp. 1659–1671, 1997.

[15] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, and Y. Tian, "Deep residual learning in spiking neural networks," *Advances in neural information processing systems*, vol. 34, pp. 21 056–21 069, 2021.

[16] Z. Zhou, Y. Zhu, C. He, Y. Wang, S. Yan, Y. Tian, and L. Yuan, "Spikformer: When spiking neural network meets transformer," *arXiv preprint arXiv:2209.15425*, 2022.

[17] M. Yao, J. Hu, Z. Zhou, L. Yuan, Y. Tian, B. Xu, and G. Li, "Spike-driven transformer," *Advances in neural information processing systems*, vol. 36, pp. 64 043–64 058, 2023.

[18] X. Shi, Z. Hao, and Z. Yu, "Spikingresformer: bridging resnet and vision transformer in spiking neural networks," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2024, pp. 5610–5619.

[19] M. Yao, J. Hu, T. Hu, Y. Xu, Z. Zhou, Y. Tian, B. Xu, and G. Li, "Spike-driven transformer v2: Meta spiking neural network architecture inspiring the design of next-generation neuromorphic chips," *arXiv preprint arXiv:2404.03663*, 2024.

[20] X. Luo, M. Yao, Y. Chou, B. Xu, and G. Li, "Integer-valued training and spike-driven inference spiking neural network for high-performance and energy-efficient object detection," in *European Conference on Computer Vision*. Springer, 2024, pp. 253–272.

[21] M. Yao, X. Qiu, T. Hu, J. Hu, Y. Chou, K. Tian, J. Liao, L. Leng, B. Xu, and G. Li, "Scaling spike-driven transformer with efficient spike firing approximation training," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2025.

[22] T. Song, G. Jin, P. Li, K. Jiang, X. Chen, and J. Jin, "Learning a spiking neural network for efficient image deraining," in *IJCAI*, 2024.

[23] L. Chen, X. Chu, X. Zhang, and J. Sun, "Simple baselines for image restoration," in *European conference on computer vision*. Springer, 2022, pp. 17–33.

[24] J. Kim, K. Kim, and J.-J. Kim, "Unifying activation-and timing-based learning rules for spiking neural networks," *Advances in neural information processing systems*, vol. 33, pp. 19 534–19 544, 2020.

[25] Y. Hu, L. Deng, Y. Wu, M. Yao, and G. Li, "Advancing spiking neural networks toward deep residual learning," *IEEE transactions on neural networks and learning systems*, vol. 36, no. 2, pp. 2353–2367, 2024.

[26] W. Fang, Y. Chen, J. Ding, Z. Yu, T. Masquelier, D. Chen, L. Huang, H. Zhou, G. Li, and Y. Tian, "Spikingjelly: An open-source machine learning infrastructure platform for spike-based intelligence," *Science Advances*, vol. 9, no. 40, p. eadi1480, 2023. [Online]. Available: https://www.science.org/doi/abs/10.1126/sciadv.adi1480

[27] M. Horowitz, "1.1 computing's energy problem (and what we can do about it)," in *2014 IEEE international solid-state circuits conference digest of technical papers (ISSCC)*. IEEE, 2014, pp. 10–14.

[28] H. Zheng, Y. Wu, L. Deng, Y. Hu, and G. Li, "Going deeper with directly-trained larger spiking neural networks," in *Proceedings of the AAAI conference on artificial intelligence*, vol. 35, no. 12, 2021, pp. 11 062–11 070.

[29] K. He, X. Zhang, S. Ren, and J. Sun, "Identity mappings in deep residual networks," in *European conference on computer vision*. Springer, 2016, pp. 630–645.

[30] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.

[31] J. Johnson, A. Alahi, and L. Fei-Fei, "Perceptual losses for real-time style transfer and super-resolution," in *European conference on computer vision*. Springer, 2016, pp. 694–711.

[32] Q. Chen and V. Koltun, "Photographic image synthesis with cascaded refinement networks," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 1511–1520.

[33] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *The Thirty-Seventh Asilomar Conference on Signals, Systems & Computers, 2003*, vol. 2. Ieee, 2003, pp. 1398–1402.

[34] R. Zhang, P. Isola, A. A. Efros, E. Shechtman, and O. Wang, "The unreasonable effectiveness of deep features as a perceptual metric," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 586–595.

[35] H. Wang, Q. Xie, Q. Zhao, Y. Li, Y. Liang, Y. Zheng, and D. Meng, "Rcdnet: An interpretable rain convolutional dictionary network for single image deraining," *IEEE Transactions on Neural Networks and Learning Systems*, 2023.

[36] S. W. Zamir, A. Arora, S. Khan, M. Hayat, F. S. Khan, and M.-H. Yang, "Restormer: Efficient transformer for high-resolution image restoration," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 5728–5739.

[37] Z. Shi, N. Fan, D.-Y. Yeung, and Q. Chen, "Stereo waterdrop removal with row-wise dilated attention," in *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2021, pp. 3829–3836.

[38] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1125–1134.

[39] R. Qian, R. T. Tan, W. Yang, J. Su, and J. Liu, "Attentive generative adversarial network for raindrop removal from a single image," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2482–2491.

[40] Y.-L. Liu, W.-S. Lai, M.-H. Yang, Y.-Y. Chuang, and J.-B. Huang, "Learning to see through obstructions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 215–14 224.

[41] Y. Quan, S. Deng, Y. Chen, and H. Ji, "Deep learning for seeing through window with raindrops," in *Proceedings of the IEEE/CVF international conference on computer vision*, 2019, pp. 2463–2471.

[42] H. Zhang and V. M. Patel, "Density-aware single image de-raining using a multi-stream dense network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 695–704.

[43] R. Li, L.-F. Cheong, and R. T. Tan, "Heavy rain image restoration: Integrating physics model and conditional adversarial learning," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 1633–1642.

[44] C. Guo, C. Li, J. Guo, C. C. Loy, J. Hou, S. Kwong, and R. Cong, "Zero-reference deep curve estimation for low-light image enhancement," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 1780–1789.

[45] C. Wei, W. Wang, W. Yang, and J. Liu, "Deep retinex decomposition for low-light enhancement," *arXiv preprint arXiv:1808.04560*, 2018.

[46] W. Yang, S. Wang, Y. Fang, Y. Wang, and J. Liu, "From fidelity to perceptual quality: A semi-supervised approach for low-light image enhancement," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 3063–3072.

[47] J. Huang, X. Fu, Z. Xiao, F. Zhao, and Z. Xiong, "Low-light stereo image enhancement," *IEEE Transactions on Multimedia*, vol. 25, pp. 2978–2992, 2022.

[48] D. Scharstein, H. Hirschmüller, Y. Kitajima, G. Krathwohl, N. Nešić, X. Wang, and P. Westling, "High-resolution stereo datasets with subpixel-accurate ground truth," in *Pattern Recognition: 36th German Conference, GCPR 2014, Münster, Germany, September 2-5, 2014, Proceedings 36*. Springer, 2014, pp. 31–42.

[49] Y. Hua, P. Kohli, P. Uplavikar, A. Ravi, S. Gunaseelan, J. Orozco, and E. Li, "Holopix50k: A large-scale in-the-wild stereo image dataset," *arXiv preprint arXiv:2003.11172*, 2020.