LIGHTWEIGHT QUAD BAYER HYBRIDEVS DEMOSAIC ING VIA STATE SPACE AUGMENTED CROSS-ATTENTION

Anonymous authors #9614

Paper under double-blind review

ABSTRACT

Event cameras like the Hybrid Event-based Vision Sensor (HybridEVS) camera capture brightness changes as asynchronous "events" instead of frames, offering advantages over traditional cameras: high temporal resolution, wide dynamic range, and no motion blur. However, challenges arise from combining a Quad Bayer Color Filter Array (CFA) sensor with event pixels lacking color information, resulting in aliasing and artifacts on the demosaicing process before downstream application. Current methods struggle to address these issues, especially on resource-limited mobile devices. In response, we introduce **TSANet**, a lightweight Two-stage network via State space augmented cross-Attention, which can handle event pixels inpainting and Quad Bayer demosaicing separately, leveraging the benefits of dividing complex tasks into manageable subtasks and learning them through a two-step training strategy to enhance robustness. Additionally, we propose a lightweight Cross-Swin State Block (CSSB) designed to augment the model's capacity to capture global dependencies using state space models in a linear format, along with cross-modality Swin attention to integrate additional priors like CFA pattern and event map, outperforming traditional local attention mechanisms while also reducing model size. In summary, TSANet demonstrates excellent demosaicing performance on HybridEVS while maintaining a lightweight model, averaging better results than the previous state-of-the-art method DemosaicFormer across seven diverse datasets in both PSNR and SSIM, while respectively reducing parameter and computation costs by $1.86 \times$ and $3.29 \times$. Our approach presents new possibilities for efficient image demosaicing on mobile devices. Code and models are available in supplementary materials.

032 033 034

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027

028

029

031

1 INTRODUCTION

In recent years, event cameras have made sig-037 nificant progress as a new type of image sen-038 sor. Compared to conventional digital cameras, event cameras can capture event infor-040 mation by sensing the intensity change of specialized event pixels, thereby capturing infor-041 mation about objects' movement Litzenberger 042 et al. (2006); Lichtsteiner et al. (2008). This 043 additional visual data enables the precise cap-044 ture of fast-moving objects, presenting exten-045 sive potential across various domains including robotics, automotive technology, and drone 047 systems. However, basic imaging before down-048 stream applications is essential for event cameras, which has not received adequate attention. On the other hand, with the development of mobile photography technology, it has been found



Figure 1: Left: Bayer CFA vs. Quad Bayer CFA on HybridEVS. Right: PSNR vs. FLOPs running on seven datasets, each model indicates its respective parameters on the right side. Our state space (SS) based TSANet presents three different sizes, each delivering optimal performance across various complexity ranges.

that traditional Bayer CFA sensors are constrained by their design on mobile devices, making it
 difficult to attain high-quality images in low-light scenes. Therefore, non-Bayer CFA sensors have
 gradually become mainstream in mobile photography in recent years, utilizing specifically designed

CFA to enhance low-light imaging performance. Quad Bayer CFA is one of the most popular formats, which can acquire high-resolution images on common scenes and enhance low-light imaging by pixel-binning Yoo et al. (2015). However, different from Bayer CFA which has been proposed over several decades, exploration of Quad Bayer CFA is still very limited Yang et al. (2022).

058 A type of event camera design named Hybrid Event Vision Sensors (HybridEVS) Kodama et al. 059 (2023) involves utilizing a Quad Bayer CFA on a camera sensor and allocating certain pixels as 060 event pixels to capture motion data instead of RGB color information, as shown in Fig. 1. This 061 design introduces the camera with improved low-light imaging capabilities and the ability to capture 062 high-speed objects effectively. However, the non-conventional Quad Bayer arrangement and the 063 absence of color information at event pixel locations pose challenges for the demosaicing process, 064 traditional methods face difficulties in extracting patterns from such complex arrangements, resulting in reduced imaging quality and poor performance in downstream applications such as deblurring and 065 object detection. 066

067 Specifically, several challenges for demosaicing in Quad Bayer HybridEVS cameras are encountered: 068 i) How to alleviate the decrease in reconstruction quality resulting from the absence of color informa-069 tion at event pixels locations; ii) How to realize Quad Bayer demosaicing joint with denoising; iii) 070 How to reduce the model's parameter to make it valuable for practical applications in edge computing. 071 Recent studies Zhou et al. (2018); Kim & Kim (2019) have proposed some end-to-end demosaicing methods based on convolution neural networks (CNNs). However, these methods fail to produce the 072 desired results when directly applied to Non-Bayer images. Some methods have been proposed to 073 tackle the problem of joint demosaicing and denoising of Quad Bayer CFA. A Sharif et al. (2021) 074 proposed a method combining generative adversarial networks, utilizing depth and spatial attention 075 mechanisms and perceptual loss to improve the reconstruction quality of mosaic images. Zeng et 076 al. Zeng et al. (2023) proposed a dual-head network to transform noisy Quad Bayer into noise-free 077 Bayer for demosaicing tasks. However, these methods do not take into account the presence of event pixels, causing color distortion and artifacts. At the same time, these networks are mostly limited 079 by their large parameters and computational complexity, making them difficult to apply to mobile 080 devices.

081 To address the above issues, we introduce TSANet, a novel lightweight two-stage model that 082 effectively combines the position information and color information of pixel arrangements, augmented 083 by the state space model to further explore long-range relationships inside HybridEVS RAW images. 084 Specifically, to improve computational efficiency while enhancing reconstruction quality, we divide 085 the complex task of joint demosaicing and denoising for HybridEVS into two stages as shown 086 in Fig. 2. The initial stage, termed Quad-to-Quad (Q2Q), is dedicated to inpainting event pixels 087 and denoising. Subsequently, the second stage is tailored for Quad Bayer demosaicing, called 088 Quad-to-RGB (Q2R). Both stages employ U-Net like networks, with extra position branch for each sub-network. Considering the computational efficiency and task complexity, we utilize a network 089 with fewer parameters in the first stage and a network with a larger parameters in the second stage. A 090 two-step training strategy is additionally employed to effectively increase the stability and robustness. 091

092 For the design of specific models, inspired by Zheng et al. (2024); Sun et al. (2022), we use two 093 branches in the encoders of both sub-network, extend a extra position branch to utilize position information compared to common U-Net structure. This design introduces an explicit position 094 encoding, enabling the network to have prior knowledge of positional relationships. Specifically, for the fusion in Q2R stage, we introduce a Cross-Swin State Block (CSSB) (See Sec. 3.2), which 096 contains Quad Bayer Cross Swin Attention (QCSA) for local cross-modality attention within windows and the Residual Vision State Space (RVSS) Zhu et al. (2024) Model in parallel for effective long-098 range representation capture. The hybrid design fuses local window attention across position and image, while efficiently incorporating global spatial information through RVSS with linear complexity. 100 The tailored two-branch block ensures a balance between local and global information, optimizes 101 both performance and efficiency. A variant of CSSB named Conv State Block (CSB) is designed for 102 Q2R decoder, focus on inner local-global relationship of the image. We also customized a simplified 103 dot-product-based attention called Spatial Position Attention (SPA)(See Sec. 3.3) to integrate Event 104 map and Quad Bayer pattern while significantly reducing computational load.

- 105 In summary, our contributions are as follows:
- 107

- We propose TSANet, a lightweight <u>T</u>wo-stage network via <u>State</u> space augmented cross-<u>Attention</u> designed for Quad Bayer HybridEVS Demosaicing. By employing a designed subtasks allocation and dual-branch encoders, TSANet achieves state-of-the-art performance with fewer computational resources (See Fig. 1).
 - We present two unique state space augmented cross-attention blocks, termed Cross-Swin State Block (CSSB) and Conv State Block (CSB). The combined use of the Residual Vision State Space (RVSS) module with local attention and convolution demonstrates great effectiveness advantages, leading to outstanding capabilities in local-global feature extraction in a linear format.
- We design two cross-modality attention mechanisms between position and image information, named Quad Bayer Cross Swin Attention (QCSA) and Spatial Position Attention (SPA). The two modules, based on local window attention and dot product, respectively, effectively capture position features, facilitating information exchange and integration across different modalities.
- 121 122

110

111

112

113

114

115

116

117

118

119

120

123 2 RELATED WORK

125 **Baver Demosaicing** Traditional approaches for Baver demosaicing primarily rely on interpolation techniques Hirakawa & Parks (2006), utilizing methods like adaptive algorithm Hirakawa & Parks 126 (2005) and spatial-spectral correlations Li et al. (2008) to reconstruct full-color images. Recently, 127 the success of convolutional networks (CNNs) used in deep learning has led to great progress 128 in demosaicing Syu et al. (2018); Tan et al. (2017b;a; 2018); Liu et al. (2020). These methods 129 replace traditional interpolation techniques with deep neural networks, leveraging the powerful fitting 130 capabilities of neural networks to achieve excellent results. Some researchers Liu et al. (2020); 131 Guo et al. (2021); Zhang et al. (2022b) proposed a demosacing network to utilize internal image 132 information like color prior and CFA arrangement. However, most of these methods can't be extended 133 to non-Bayer CFA formats, facing limitations like artifacts and aliasing when dealing with more 134 challenging non-Bayer formats like Quad Bayer.

135 Quad Bayer Demosaicing In recent years, Quad Bayer has become a popular CFA pattern widely 136 used in mobile photography, such as smartphone cameras Yang et al. (2022). Different from traditional 137 Bayer CFA, exploration of the Quad Bayer CFA is limited. The larger gaps between pixels of the same 138 color make the task more challenging compared to Bayer CFA. Some two-stage networks Jia et al. 139 (2022); Zeng et al. (2023) are proposed for progress learning to enhance Quad Bayer demosaicing. 140 Zheng et al. (2024) proposed a dual-encoder structure to achieve better joint demosaicing and 141 denoising tasks. GAN-based networks A Sharif et al. (2021); Sharif et al. (2021) are also used to 142 strengthen the restoration of RGB images for Non-Bayer CFA sensors. The sensor's CFA offers 143 a strong positional prior, yet most methods overlook the valuable color arrangement information, 144 resulting in significant color distortion in Quad Bayer HybridEVS demosaicing.

145 **Event Camera Imaging** As a novel type of sensor that has emerged in recent years Gallego et al. 146 (2020); Son et al. (2017), the imaging research of event cameras is an active topic. Recent studies 147 have mainly focused on imaging or downstream applications like object detection Zhang et al. (2022a) 148 based on event information. Munda et al. (2018) try to reconstruct intensity images with direct event 149 integration. Scheerlinck et al. (2020) introduced a fast image restoration method with deep neural 150 networks. These methods have demonstrated excellent performance in visual tasks based on event 151 information. However, for the latest proposed hybrid event-based vision sensor Kodama et al. (2023), the field is almost blank in advanced methods for its Image Sensor Pipelines (ISPs). A crucial step in 152 ISPs is the demosaicing process, which converts Quad Bayer data into the RGB domain and directly 153 impacts the imaging quality for downstream applications. 154

State Space Models In recent years, State Space Models (SSMs) Gu et al. (2021a;b) have emerged as competitive rivals to traditional deep learning architectures like Convolution Neural Networks(CNNs) and Transformers. Pioneering works like S4 Gu et al. (2021a) and S5 Smith et al. (2022) introduced advancements on deep-state models with efficient parallel scan, modeling long-range dependency. Recently proposed Mamba Gu & Dao (2023), featuring a data-dependent SSM layer, has shown remarkable performance, surpassing Transformers in natural language tasks and demonstrating linear scalability in sequence length. Additionally, some works have applied Mamba to various vision tasks, including image classification Zhu et al. (2024), video understanding Wang et al.



179

181

Figure 2: Overview of the TSANet approach. We adopt a two-stage structure that breaks down complex tasks into manageable subtasks, while leveraging additional branches to utilize position prior. Section 3.2 introduces the Spatial Position Attention (SPA) and Cross-Swin State Block (CSSB).

(2023), image restoration Guo et al. (2024), and image segmentation Liao et al. (2024), demonstrates the potential in visual tasks for lightweight models.

199

212 213

3 PROPOSED METHOD

Initially, we present the overall pipeline of our two-stage model in Sec. 3.1. Then we introduce the proposed state space augmented-cross attention block and its variants in Sec. 3.2. Finally, we discuss details of our proposed attention mechanisms across position and image in Sec. 3.3.

191 3.1 TWO-STAGE NETWORK STRCTURE

For the task of Quad Bayer HybridEVS demosaicing, our goal is to recover a three-channel RGB image $I_{\mathbf{R}} \in \mathbb{R}^{H \times W \times 3}$ from a degraded Quad Bayer image $I_{\mathbf{Q}} \in \mathbb{R}^{H \times W \times 1}$, where the degradation includes color channel loss D_q due to the Quad Bayer CFA, pixel absence D_e due to the design of event pixels in the sensor, and noise D_n introduced during the image capture process, respectively. Most previous methods aim to directly learn the entire process through an all-in-one deep model \mathcal{M} that restores image $I_{\mathbf{R}}$ from $I_{\mathbf{Q}}$, which can be expressed as:

$$\mathbf{I}_{\mathbf{R}} = \mathcal{M}(\mathbf{I}_{\mathbf{Q}}). \tag{1}$$

200 However, these all-in-one models often struggle to extract the inner connection between position 201 and color, causing unbearable aliasing and artifacts (See Fig. 6), or require a large number of 202 parameters and computation load to achieve ideal restoration results, make it barely impossible 203 to deploy on limited-resource mobile devices. Unlike past single-model solutions, we define the composite task as two controllable sub-tasks: the former sub-task is to restore the degradation of D_e 204 and D_n from the original Quad Bayer image, inpainting absent pixels and reducing noise, producing a clean Quad Bayer image, defined as \mathcal{M}^{Q2Q} ; the later one is to restore the clean Quad Bayer image into $\mathbf{I}_{\mathbf{R}}$, defined as \mathcal{M}^{Q2R} . The overall pipeline progressively restore the RGB image from the 205 206 207 degraded Quad Bayer image. Notably, we introduce a distinctive encoder branch designed to integrate 208 position information as a dedicated prior knowledge into the network, using position information $\mathbf{P}_{\mathbf{e}} \in \mathbb{R}^{H \times W \times 1}$ from D_e and $\mathbf{P}_{\mathbf{q}} \in \mathbb{R}^{H \times W \times 3}$ from D_q to achieve better image reconstruction, the 209 210 process can be expressed as: 211

$$\mathbf{I}_{\mathbf{R}} = \mathcal{M}^{Q2R}(\mathcal{M}^{Q2Q}(\mathbf{I}_{\mathbf{Q}}, (\mathbf{P}_{\mathbf{q}}, \mathbf{P}_{\mathbf{e}})), \mathbf{P}_{\mathbf{q}}).$$
(2)

Our proposed two-stage network architecture is depicted in Fig. 2. Such a design not only assigns
 specific tasks to sub-networks but also benefits from a two-step training strategy. Prior studies have demonstrated that pretraining on sub-networks can lead to improved performance and inference

stability. Our designed two-stage architecture facilitates directional pretraining for the two sub networks by synthesizing a dummy clean Quad Bayer.

Besides, we believe fully utilizing position information is crucial for the demosaicing problem 219 of Quad Bayer, especially in the task with event pixels. Therefore, we propose an additional 220 positional encoding branch in both networks, explicitly integrating position information into the 221 network. To further enhance the model's ability to extract high-frequency texture information, 222 we designed a Fourier Feature Module (FFM) based on Fourier encoding, as shown in Fig. 2, 223 which maps the position information to a series of high-frequency features based on sine and cosine 224 functions Tancik et al. (2020), enabling the model to capture refine details and patterns in position 225 dimension, strengthening the restoration results of complex textures. At the Q2Q stage, specially, 226 before putting the Quad Bayer image into the network, we apply a coarse inpainting by averaging nearby pixels around event pixels, which aims to mitigate color loss resulting from event pixels. 227

228 3.2 STATE SPACE AUGMENTED BLOCKS

In this section, we propose two lightweight modules augmented by State Space Models for the
 encoder and decoder of the Q2R stage, respectively. We begin with a dual-branch structure named
 Cross-Swin State Block (CSSB), concurrently modeling local cross-modality attention and long-range
 dependencies with linear complexity. Then, we present its variant Conv State Block (CSB) in the
 decoder, enhancing local feature restoration while further reducing computation by convolutions.

Cross-Swin State Block Fig. 3a illustrated our proposed Cross-Swin State Block (CSSB). This 235 block is designed to capture long-range dependencies and cross-modality local attention in parallel. 236 It integrates Residual Vision State Space (RVSS) and Quad Bayer Cross Swin Attention (QCSA) 237 mentioned in Sec. 3.3. Fig. 4c shows RVSS, a simplified version of Residual State Space Block 238 of MambaIR Guo et al. (2024), preserving its core component VSSM for efficient extraction of 239 long-range dependencies, followed by a residual connection. To further reduce computational 240 complexity while capturing local positional attention intersections and global long-range dependencies 241 simultaneously, we follow SCUNet Zhang et al. (2023) and propose a parallel network module. First, 242 the image feature projected through 1x1 convolution is split and separately inputted into QCSA and 243 RVSS modules. The outputs are then concatenated and utilized for out projection through a 1×1 244 convolution, which is followed by a residual connection. For a image input $\mathbf{F}_{\mathbf{I}}$ and a position input 245 $\mathbf{F}_{\mathbf{P}}$, the process can be expressed as follows:

246
247
248
249

$$X_1, X_2 = \text{Split}(\text{Conv}_{1 \times 1}(\mathbf{F_I})),$$

 $Y_1, Y_2 = \text{QCSA}(X_1, \mathbf{F_P}), \text{RVSS}(X_2),$
 $\hat{\mathbf{F}}_{\mathbf{I}} = \text{Conv}_{1 \times 1}(\text{Concat}(Y_1, Y_2)) + \mathbf{F}_{\mathbf{I}}.$
(3)

The QCSA primarily extracts representations with position information, models spatial information
 through the two modalities of position and image while RVSS is used to effectively capture global
 information, which parallel address cross-modality local attention and long-range dependencies.
 Moreover, the dual-brunch structure is similar to group convolution, reducing the number of channels
 within the module through splitting operations, effectively reducing the computational complexity
 and parameters of the block.

Conv State Block We also pro-256 pose a variant of CSSB for the 257 decoder of the Q2R stage, as 258 shown in Fig. 3b. Instead of em-259 ploying QCSA to capture atten-260 tion features across position in-261 formation, we replace this mod-262 ule with Residual Convolution 263 (RConv) Zhang et al. (2021), 264 which focuses on restoring in-265 ternal local feature of the im-266 age, parallel with an RVSS to en-267 hance long-range feature extraction capability, thereby forming 268 a lightweight decoder block with 269 local-global dependencies.



(a) Cross-Swin State Block

(b) Conv State Block

Figure 3: Proposed state space augmented blocks. The two blocks are modified from Transformer and Convolution for the Q2R encoder and decoder, respectively. The dual-branch design parallel extracts feature with local-global dependencies while reducing computation load.

Quad Baver Cross Swin Attention (QCSA)

(a) Quad Bayer Cross Swin Attention

270

276 277 278

279

280

275

Figure 4: Illustration of designed modules. We propose two cross-attention modules and one state space model. Attentions focus on fusing position information into network and RVSS captures long-range dependencies.

(b) Spatial Position Attention

Spatial Position Attention (SPA)

1x1 Con

281 282 283

284 285

286

287 288

291

296 297

298

299

300

301 302

310 311

3.3 ATTENTION MODULES FOR POSITION FUSION

As illustrated in Fig. 4a and Fig. 4b, we propose two spatial attention modules designed for fusing position information and image feature. We utilize the position representations as extra information and explore spatial relationship of position and image.

Quad Bayer Cross Swin Attention As shown in Fig. 4a, we designed a sufficient feature fusion 289 mechanism for Q2R sub-network. It first employs an efficient Window Multi-head Cross Atten-290 tion (WMCA) inspired by Swin Transformer Liu et al. (2021). Specifically, given image feature $\mathbf{F}_{\mathbf{I}} \in \mathbb{R}^{H \times W \times d}$ and position feature $\mathbf{F}_{\mathbf{P}} \in \mathbb{R}^{H \times W \times d}$, both features are first partitioned into non-292 overlapping $M \times M$ local windows, getting $\frac{HW}{M^2} \times M^2 \times d$ features. Different from conventional attention acquiring all query(Q), key(K), value(V) from $\mathbf{F_I}$, we produce Q from $\mathbf{F_P}$. The Q, K and V for window feature $X \in \mathbb{R}^{M^2 \times d}$ from $\mathbf{F_I}$ and $Y \in \mathbb{R}^{M^2 \times d}$ from $\mathbf{F_P}$ are computed as: 293 295

$$Q = YP_Q, \quad K = XP_K, \quad V = XP_V. \tag{4}$$

Residual Vision State Space (RVSS)

(c) Residual Vision State Space

where P_Q , P_K and P_V are shared project matrices among each window. Then we have $Q, K, V \in$ $\mathbb{R}^{M^2 \times d}$ to compute in-window cross attention as:

> Attention $(Q, K, V) = \text{SoftMax}\left(\frac{QK^T}{\sqrt{d}} + B\right).$ (5)

where B represents relative positional encoding within the window, complementing the global 303 positional encoding introduced by $\mathbf{F}_{\mathbf{P}}$. This attention mechanism operates across all windows and 304 is executed in parallel h times Vaswani et al. (2017). After WMCA, the features are then fed into a 305 two-layer multi-layer perceptron (MLP) with GELU Hendrycks & Gimpel (2016) activation function 306 for further feature extraction. Both steps utilize residual connections and LayerNorm Ba et al. (2016), 307 the process can be expressed as: 308

$$\mathbf{F}_{\mathbf{I}} = WMCA(LN(\mathbf{F}_{\mathbf{I}}), LN(\mathbf{F}_{\mathbf{P}})) + \mathbf{F}_{\mathbf{I}},$$

$$\hat{\mathbf{F}}_{\mathbf{I}} = MLP(\hat{\mathbf{F}}_{\mathbf{I}}) + \hat{\mathbf{F}}_{\mathbf{I}}.$$
(6)

Then, through the shifted window mechanism Liu et al. (2021), this module achieves cross-window 312 information exchange. It incorporates global information of Quad Bayer CFA pattern into the network, 313 offering a spatial attention mechanism based on position information. Additionally, the window 314 mechanism maintains the computational complexity of the network linearly, effectively controlling 315 computational and parameter overhead compared to conventional attention methods. 316

317 **Spatial Position Attention** As shown in Fig. 4b, we designed another cross-modality attention 318 mechanism, aimed at building the relationship between position and image sufficiently in the Q2Q 319 stage. Specifically, the position branch introduces an explicit representation of the spatial dimension. 320 Firstly, the position feature $\mathbf{F}_{\mathbf{P}} \in \mathbb{R}^{H \times W \times d}$ and the image feature $\mathbf{F}_{\mathbf{I}} \in \mathbb{R}^{H \times W \times d}$ pass through their 321 respective convolutional projection layers. Subsequently, the position branch is activated by ReLU, 322 after which they demonstrate element-wise product with each other. The process can be expressed as: 323

$$\hat{\mathbf{F}}_{\mathbf{I}} = \operatorname{Conv}_{1 \times 1} \left(\operatorname{LN} \left(\mathbf{F}_{\mathbf{I}} \right) \right) \odot \operatorname{ReLU} \left(\operatorname{Conv}_{1 \times 1} \left(\operatorname{LN} \left(\mathbf{F}_{\mathbf{P}} \right) \right) + \mathbf{F}_{\mathbf{I}}.$$
(7)



Figure 5: Visualized results of all compared methods for Quad Bayer HybridEVS Demosaicing on MIPI dataset. The proposed TSANet demonstrates the best visual results among all methods, producing more vivid colors on scenes with complex coloration, better than the previous state-of-theart approach Restormer.

356

357

358

359

350

351

The position information is calculated to a weights map of image feature, determining the importance of each feature pixel. It is similar to a Gate Mechanism Zamir et al. (2022), which can control the flow of information but by position information rather than itself. Additionally, this method exhibits faster computational performance when compared to convolution or attention mechanisms.

4 EXPERIMENTS

In this section, we first introduce experimental settings including implementation details and datasets.
 Then we compared our proposed TSANet with other state-of-the-art methods on seven diverse datasets. Finally, we demonstrate an ablation study to prove the effectiveness of our methods.

4.1 EXPERIMENTAL SETTINGS

Implementation details In all experiments, we use the following hyperparameters, unless mentioned otherwise. During the training, we randomly crop the original Quad Bayer input and RGB ground truth into 128×128 patches, with batch size = 32. We use Adam as the optimizer and the learning rate starts from 2×10^{-4} and is gradually reduced to 1×10^{-7} with the cosine annealing scheme. The total iteration is set to 1×10^{6} . Specifically, to fully utilize the two-stage structure of our TSANet, we apply a pretraining step on the sub-network before end-to-end joint training. In particular, we employ published pretrained weights of DemosaicFormer for testing.

371

Datasets We train all the models with the dataset from Mobile Intelligent Photography & Imaging (MIPI) Workshop 2024 Demosaic for Hybridevs Camera challenge trackWu et al. (2024), which contains 800 pairs of Quad Bayer and RGB images with 2K resolution. The official public test set of MIPI dataset contains 26 pairs. Both the training and testing data includes real world noise. we simulate HybridEVS pattern test cases from five image datasets, including Kodak Loui et al. (2007), Urban100 Cordts et al. (2016), BSD100 Martin et al. (2001), and Wed Ma et al. (2017) (first 100 images). Additionally, we assess dynamic performance using two video datasets: REDS Nah et al.

378		Doromo	EL OD.		In	nage Datase	ets		Video I	Datasets	
379	Methods		(C)	Kodak	BSD100	Urban100	Wed	MIPI	REDS	Vid4	Average
380			(0)		PSNR/SSIM						
381	DemosaicFormer	30.28	491.1	39.32/0.982	37.65 /0.982	37.64 /0.980	34.86/0.968	39.35/0.981	42.45/0.991	36.01 /0.979	38.18/0.980
382	TSANet-l (Ours)	16.26	149.4	39.40/0.986	37.34/ 0.986	37.07/ 0.983	35.76/0.977	39.07/0.980	43.00/0.996	35.64/ 0.982	38.18/0.984
383	NAFNet	63.16	126.7	39.14/0.985	37.51/0.986	36.64/0.982	35.73/0.969	38.89/0.979	42.76/0.996	35.48/0.981	38.02/0.983
000	Restormer	26.11	282.2	39.16/0.986	37.11/0.985	36.36/0.977	35.00/0.971	38.42/0.978	41.91/0.990	35.08/0.980	37.58/0.981
384	SAGAN	22.56	341.6	36.14/0.974	30.53/0.931	29.89/0.946	28.22/0.917	34.25/0.959	38.13/0.984	32.16/0.963	32.76/0.953
385	TSANet-m (Ours)	8.74	81.94	39.24/0.986	37.25/ 0.986	36.75/0.982	35.60/ 0.976	38.93/0.979	42.87/0.996	35.53/0.982	38.02/0.984
386	PIPNet	3.46	68.8	32.20/0.960	31.97/0.950	28.92/0.942	29.19/0.929	33.73/0.950	36.19/0.981	32.20/0.964	32.06/0.954
387	CycleISP	3.23	104.9	33.09/0.970	32.18/0.969	29.78/0.942	30.22/0.944	30.04/0.934	32.96/0.975	30.46/0.964	31.25/0.957
388	TSANet-s (Ours)	4.00	37.4	38.73/0.984	36.56/0.984	36.15/0.980	35.19/0.973	38.47/0.978	41.94/0.989	35.15/0.980	37.46/0.981

Table 1: Quantitative evaluation of TSANet compared to other methods across seven diverse datasets. Methods with an orange background are Transformer-based, while those with a blue background are Convolution-based. Our TSANet provides three model sizes and consistently achieves the best average PSNR and SSIM across various computational complexity levels, while significantly reducing both parameters and FLOPs. Specifically, Our TSANet-1 achieves the best average PSNR and SSIM scores across seven datasets, surpassing DemosaicFormer with 0.004 on SSIM while reducing parameter and computation costs by $1.86 \times$ and $3.29 \times$.



Figure 6: Visualized results of all compared methods for Quad Bayer HybridEVS Demosaicing on synthesized image datasets. "Demosaicformer" is abbreviated as "DFormer". The comparison provides further validation of TSANet across various scenarios and confirms its effectiveness in restoring fine details, colors, and textures.

418

419

389

390

391

392

393

394

395

(2019) and Vid4 Liu & Sun (2011). We re-sample Quad Bayer image from RGB images and simulate event pixels to synthesize the input.

424 425 426

4.2 COMPARISON TO STATE-OF-THE-ARTS

Quantitative Comparison. We compare our proposed TSANet with several state-of-the-art methods, including two joint demosaicing and denoising methods PIPNet A Sharif et al. (2021) and SAGAN Sharif et al. (2021), four image restoration methods, Restormer Zamir et al. (2022) and DemosaicFormer Xu et al. (2024) based on Transformer, NAFNet Chen et al. (2022) and CycleISP Zamir et al. (2020) based on Convolution. We choose PSNR/SSIM scores as restoration quality metrics, parameters and FLOPs as computational complexity metrics, to illustrate models' performance. Table





CvcleISP

TSANet (Ours)

PIPNet

Restormen

Case 1 Case 2 Case 3 Case 4 Reference Case 5 Case 6 PSNR/SSIM 34.59/0.9776 34.62/0.9779 34.65/0.9780 34.90/0.9789 34.71/0.9780 34.96/0.9792

459 Figure 8: Visualization of ablation study. We displayed the visual comparison results with a difference 460 map of our ablation studies, which separately validate various components, thus proving the efficacy of our module design.

462 1 shows the quantitative comparison results on all test datasets. It is worth noting that due to its 463 large model size, Restormer and DemosaicFormer couldn't be tested on the 2k MIPI dataset, so we 464 partitioned input images into 700x700 patches for inference and recombined the results. In particular, 465 the proposed TSANet achieves state-of-the-art performance in different complexity levels across 466 seven diverse test datasets, surpassing the previous SOTA method DemosaicFormer by 0.004 in 467 SSIM, while reducing parameters and computations by 1.86x and 3.29x respectively. Additionally, 468 on seven synthetic datasets, TSANet-1 achieves the best PSNR/SSIM results in three and second-best results in four. Furthermore, smaller versions of TSANet-s and TSANet-m also achieve the best 469 results on corresponding complexity ranges. Experimental comparisons validate the effectiveness of 470 our approach, ensuring restoration performance while drastically reducing computational resources, 471 introducing a model friendly to edge devices with limited computational resources. 472

473

Visualization Comparison. Fig. 5 and Fig. 6 respectively depict the visual comparison results 474 of TSANet compared with other models on the MIPI dataset and other synthesized image datasets. 475 Benefiting from our local-global feature extraction structure, our model exhibits closer visual sim-476 ilarity to ground truth in image details, particularly in subtle color variations, resulting in more 477 vivid colors compared to other methods. While others suffer from color loss due to the Quad Bayer 478 pattern and event pixels, leading to further degradation of color information in fine textures and even 479 causing severe pixel errors (see Fig. 5). Further comparison on video datasets is demonstrated in 480 Fig.7, the finer details of our TSANet demonstrate strong competitiveness in handling high-speed 481 dynamic range scenes captured by event cameras. Overall, the perceptual visual results confirm the effectiveness of our approach. 482

- 483 4.3 ABLATION STUDY 484
- As shown in Table 2, we present ablation experiments to validate the effectiveness of the proposed 485 QCSA, SPA, RVSS modules and pretraining strategy in TSANet. Models are evaluated on the MIPI

457

458

dataset. First, we remove both the QCSA and SPA cross-attention modules, then add each module separately to assess their impact on network performance. As shown in Fig. 9, their combined inclusion improves PSNR by 0.06dB with only an 8% increase in parameters. Additionally, we evaluate the improvements introduced by RVSS and the two-step training strategy independently. The pretraining step can effectively improve model performance by 0.13dB. RVSS reduces parameters by 38% with just a 0.02dB drop in PSNR, highlighting the effectiveness of incorporating state space models. Visual comparison is demonstrated in Fig. 8. Results showed that attention modules effectively enhance model performance without notably increasing the number of parameters or computational load. Despite the introduction of RVSS leading to a slight performance decrease, the reduction in parameters and computational load is significant. Our final model achieves a balance between computational resources and effectiveness.



Figure 9: Quantitative visualization of the ablation study. The position prior branch we adopted significantly improves performance in PSNR by 0.06dB with only an 8% increase in parameters, while the state space modules reduce parameters by 38%, resulting in just a 0.02dB drop in PSNR.

Table 2: Quantitative results of ablation study. We validated the impact of the QCSA, SPA, RVSS
 modules and the two-step training strategy on TSANet-s, demonstrating the effectiveness of our
 proposed cross-attentions and training recipe while validating that employing RVSS can significantly
 reduce computational resources and maintain performance.

Case		М	odules/St	rategy	PARAMs	FLOPs	MIPI	
cust	QCSA	SPA	RVSS	Two-step Training	(M)	(G)	PSNR/SSIM	
1			\checkmark	\checkmark	3.70	32.0	38.41/0.9773	
2	\checkmark		\checkmark	\checkmark	3.64	32.5	38.42/0.9773	
3		\checkmark	\checkmark	\checkmark	4.06	36.9	38.45/0.9775	
4	\checkmark	\checkmark		\checkmark	6.42	59.1	38.49/0.9776	
5	\checkmark	\checkmark	\checkmark		4.00	37.4	38.34/0.9772	
6	\checkmark	\checkmark	\checkmark	\checkmark	4.00	37.4	38.47/0.9775	

5 CONCLUSION

We have presented a novel lightweight two-stage structure network tailored for the HybridEVS architecture, which introduces task-specific sub-networks and a corresponding two-step training strategy. Specifically, we introduce Cross-Swin State Block (CSSB) and Conv State Block (CSB) strengthened by Residual Vision State Space (RVSS) to maintain low computational complexity while simultaneously addressing local position features and long-range dependencies. We further propose the Quad Bayer Cross Swin Attention (QCSA) and Spatial Position Attention (SPA) mechanisms to effectively couple the arrangement information of Quad Bayer pattern and event points in the network, providing explicit prior encoding for global position information. To the best of our knowledge, this is the first work employing SSMs in a hybrid model for demosaicing tasks. Our approach is validated across multiple datasets, offering a lightweight and mobile-friendly model for Quad Bayer HybridEVS Demosaicing.

540	REFERENCES
541	

563

573

574

575

578

579

580

- SM A Sharif, Rizwan Ali Naqvi, and Mithun Biswas. Beyond joint demosaicking and denoising:
 An image processing pipeline for a pixel-bin image sensor. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 233–242, 2021.
- Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.
- Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration.
 In *European conference on computer vision*, pp. 17–33. Springer, 2022.
- Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2016.
- Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi,
 Stefan Leutenegger, Andrew J Davison, Jörg Conradt, Kostas Daniilidis, et al. Event-based vision:
 A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(1):154–180, 2020.
- Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv* preprint arXiv:2312.00752, 2023.
- Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured
 state spaces. *arXiv preprint arXiv:2111.00396*, 2021a.
- 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.
 Advances in neural information processing systems, 34:572–585, 2021b.
- Hang Guo, Jinmin Li, Tao Dai, Zhihao Ouyang, Xudong Ren, and Shu-Tao Xia. Mambair: A simple baseline for image restoration with state-space model. In *ECCV*, 2024.
- Shi Guo, Zhetong Liang, and Lei Zhang. Joint denoising and demosaicking with green channel prior
 for real-world burst images. *IEEE Transactions on Image Processing*, 30:6930–6942, 2021. doi:
 10.1109/TIP.2021.3100312.
 - Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). arXiv preprint arXiv:1606.08415, 2016.
- Keigo Hirakawa and Thomas W Parks. Adaptive homogeneity-directed demosaicing algorithm. *Ieee transactions on image processing*, 14(3):360–369, 2005.
 - Keigo Hirakawa and Thomas W Parks. Joint demosaicing and denoising. *IEEE Transactions on Image Processing*, 15(8):2146–2157, 2006.
- Jun Jia, Hanchi Sun, Xiaohong Liu, Longan Xiao, Qihang Xu, and Guangtao Zhai. Learning rich information for quad bayer remosaicing and denoising. In *European Conference on Computer Vision*, pp. 175–191. Springer, 2022.
- Yongnam Kim and Yunkyung Kim. High-sensitivity pixels with a quad-wrgb color filter and spatial deep-trench isolation. *Sensors*, 19(21):4653, 2019.
- Kazutoshi Kodama, Yusuke Sato, Yuhi Yorikado, Raphael Berner, Kyoji Mizoguchi, Takahiro Miyazaki, Masahiro Tsukamoto, Yoshihisa Matoba, Hirotaka Shinozaki, Atsumi Niwa, et al. 1.22 μm 35.6 mpixel rgb hybrid event-based vision sensor with 4.88 μm-pitch event pixels and up to 10k event frame rate by adaptive control on event sparsity. In 2023 IEEE International Solid-State Circuits Conference (ISSCC), pp. 92–94. IEEE, 2023.
- Xin Li, Bahadir Gunturk, and Lei Zhang. Image demosaicing: A systematic survey. In *Visual Communications and Image Processing 2008*, volume 6822, pp. 489–503. SPIE, 2008.

630

631

640

594	Weibin Liao, Yinghao Zhu, Xinyuan Wang, Cehngwei Pan, Yasha Wang, and Liantao Ma. Lightm-
595	unet: Mamba assists in lightweight unet for medical image segmentation. arXiv preprint
596	arXiv:2403.05246, 2024.
597	

- Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A 128×128 120 db 15 μ s latency 598 asynchronous temporal contrast vision sensor. IEEE Journal of Solid-State Circuits, 43(2):566-576, 2008. doi: 10.1109/JSSC.2007.914337. 600
- 601 Martin Litzenberger, Christoph Posch, D Bauer, Ahmed Nabil Belbachir, P Schon, B Kohn, and 602 H Garn. Embedded vision system for real-time object tracking using an asynchronous transient 603 vision sensor. In 2006 IEEE 12th Digital Signal Processing Workshop & 4th IEEE Signal Processing Education Workshop, pp. 173–178. IEEE, 2006. 604
- 605 Ce Liu and Deqing Sun. A bayesian approach to adaptive video super resolution. In Proceedings of 606 the IEEE Conference on Computer Vision and Pattern Recognition, pp. 209–216. IEEE, 2011. 607
- 608 Lin Liu, Xu Jia, Jianzhuang Liu, and Qi Tian. Joint demosaicing and denoising with self guidance. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 609 2240-2249, 2020. 610
- 611 Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. 612 Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the* 613 IEEE/CVF international conference on computer vision, pp. 10012–10022, 2021. 614
- Alexander Loui, Jiebo Luo, Shih-Fu Chang, Dan Ellis, Wei Jiang, Lyndon Kennedy, Keansub Lee, and 615 Akira Yanagawa. Kodak's consumer video benchmark data set: concept definition and annotation. 616 In Proceedings of the International Workshop on Workshop on Multimedia Information Retrieval, 617 MIR '07, pp. 245–254, New York, NY, USA, 2007. Association for Computing Machinery. 618 ISBN 9781595937780. doi: 10.1145/1290082.1290117. URL https://doi.org/10.1145/ 619 1290082.1290117. 620
- 621 Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei 622 Zhang. Waterloo exploration database: New challenges for image quality assessment models. *IEEE Transactions on Image Processing*, 26(2):1004–1016, 2017. doi: 10.1109/TIP.2016.2631888. 623
- 624 D. Martin, C. Fowlkes, D. Tal, and J. Malik. A database of human segmented natural images and 625 its application to evaluating segmentation algorithms and measuring ecological statistics. In 626 Proceedings Eighth IEEE International Conference on Computer Vision. ICCV 2001, volume 2, 627 pp. 416-423 vol.2, 2001. doi: 10.1109/ICCV.2001.937655.
- Gottfried Munda, Christian Reinbacher, and Thomas Pock. Real-time intensity-image reconstruction 629 for event cameras using manifold regularisation. International Journal of Computer Vision, 126 (12):1381–1393, 2018.
- 632 Seungjun Nah, Sungyong Baik, Seokil Hong, Gyeongsik Moon, Sanghyun Son, Radu Timofte, and 633 Kyoung Mu Lee. Ntire 2019 challenge on video deblurring and super-resolution: Dataset and 634 study. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops, pp. 0-0, 2019. 635
- 636 Cedric Scheerlinck, Henri Rebecq, Daniel Gehrig, Nick Barnes, Robert Mahony, and Davide Scara-637 muzza. Fast image reconstruction with an event camera. In Proceedings of the IEEE/CVF Winter 638 Conference on Applications of Computer Vision, pp. 156–163, 2020. 639
 - SMA Sharif, Rizwan Ali Naqvi, and Mithun Biswas. Sagan: adversarial spatial-asymmetric attention for noisy nona-bayer reconstruction. arXiv preprint arXiv:2110.08619, 2021.
- 642 Jimmy TH Smith, Andrew Warrington, and Scott W Linderman. Simplified state space layers for 643 sequence modeling. arXiv preprint arXiv:2208.04933, 2022. 644
- 645 B Son, Y Suh, S Kim, H Jung, JS Kim, C Shin, K Park, K Lee, J Park, J Woo, et al. A 640× 480 dynamic vision sensor with a 9 um pixel and 300 meps address-event representation. In 646 Proceedings of the IEEE International Conference on Solid-State Circuits, San Francisco, CA, 647 USA, pp. 5-9, 2017.

648 649 650	Lei Sun, Christos Sakaridis, Jingyun Liang, Qi Jiang, Kailun Yang, Peng Sun, Yaozu Ye, Kaiwei Wang, and Luc Van Gool. Event-based fusion for motion deblurring with cross-modal attention. In <i>European conference on computer vision</i> , pp. 412–428. Springer, 2022.
652 653	Nai-Sheng Syu, Yu-Sheng Chen, and Yung-Yu Chuang. Learning deep convolutional networks for demosaicing. <i>arXiv preprint arXiv:1802.03769</i> , 2018.
654 655 656 657	Daniel Stanley Tan, Wei-Yang Chen, and Kai-Lung Hua. Deepdemosaicking: Adaptive image demosaicking via multiple deep fully convolutional networks. <i>IEEE Transactions on Image Processing</i> , 27(5):2408–2419, 2018.
658 659 660	Hanlin Tan, Xiangrong Zeng, Shiming Lai, Yu Liu, and Maojun Zhang. Joint demosaicing and denoising of noisy bayer images with admm. In 2017 IEEE International Conference on Image Processing (ICIP), pp. 2951–2955. IEEE, 2017a.
661 662 663	Runjie Tan, Kai Zhang, Wangmeng Zuo, and Lei Zhang. Color image demosaicking via deep residual learning. In <i>Proc. IEEE Int. Conf. Multimedia Expo (ICME)</i> , volume 2, pp. 6, 2017b.
664 665 666 667	Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn high frequency functions in low dimensional domains. <i>Advances in neural information processing systems</i> , 33:7537–7547, 2020.
668 669 670	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. <i>Advances in neural information processing systems</i> , 30, 2017.
672 673 674	Jue Wang, Wentao Zhu, Pichao Wang, Xiang Yu, Linda Liu, Mohamed Omar, and Raffay Hamid. Se- lective structured state-spaces for long-form video understanding. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 6387–6397, 2023.
675 676 677	Yaqi Wu, Zhihao Fan, Xiaofeng Chu, Jimmy S Ren, Xiaoming Li, Zongsheng Yue, Chongyi Li, Shangcheng Zhou, Ruicheng Feng, Yuekun Dai, et al. Mipi 2024 challenge on demosaic for hybridevs camera: Methods and results. <i>arXiv preprint arXiv:2405.04867</i> , 2024.
679 680 681	Senyan Xu, Zhijing Sun, Jiaying Zhu, Yurui Zhu, Xueyang Fu, and Zheng-Jun Zha. Demosaicformer: Coarse-to-fine demosaicing network for hybridevs camera. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition Workshops</i> , pp. 1126–1135, June 2024.
682 683 684	Qingyu Yang, Guang Yang, Jun Jiang, Chongyi Li, Ruicheng Feng, Shangchen Zhou, Wenxiu Sun, Qingpeng Zhu, Chen Change Loy, Jinwei Gu, et al. Mipi 2022 challenge on quad-bayer re-mosaic: Dataset and report. In <i>European Conference on Computer Vision</i> , pp. 21–35. Springer, 2022.
685 686 687	Yoonjong Yoo, Jaehyun Im, and Joonki Paik. Low-light image enhancement using adaptive digital pixel binning. <i>Sensors</i> , 15(7):14917–14931, 2015.
688 689 690 691	Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Cycleisp: Real image restoration via improved data synthesis. In <i>Proceedings</i> of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 2696–2705, 2020.
692 693 694 695	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 <i>Proceedings of the IEEE/CVF conference on computer vision and pattern recognition</i> , pp. 5728–5739, 2022.
696 697 698 699	Haijin Zeng, Kai Feng, Jiezhang Cao, Shaoguang Huang, Yongqiang Zhao, Hiep Luong, Jan Ael- terman, and Wilfried Philips. Inheriting bayer's legacy-joint remosaicing and denoising for quad bayer image sensor. <i>arXiv preprint arXiv:2303.13571</i> , 2023.
700 701	Jiqing Zhang, Bo Dong, Haiwei Zhang, Jianchuan Ding, Felix Heide, Baocai Yin, and Xin Yang. Spiking transformers for event-based single object tracking. In <i>Proceedings of the IEEE/CVF</i> <i>conference on Computer Vision and Pattern Recognition</i> , pp. 8801–8810, 2022a.

702 703 704	Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. <i>IEEE Transactions on Pattern Analysis and Machine Intelligence</i> , 44(10):6360–6376, 2021.
705 706 707 708	Kai Zhang, Yawei Li, Jingyun Liang, Jiezhang Cao, Yulun Zhang, Hao Tang, Deng-Ping Fan, Radu Timofte, and Luc Van Gool. Practical blind image denoising via swin-conv-unet and data synthesis. <i>Machine Intelligence Research</i> , 20(6):822–836, 2023.
709 710	Tao Zhang, Ying Fu, and Cheng Li. Deep spatial adaptive network for real image demosaicing. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pp. 3326–3334, 2022b.
711 712 713 714 715	Bolun Zheng, Haoran Li, Quan Chen, Tingyu Wang, Xiaofei Zhou, Zhenghui Hu, and Chenggang Yan. Quad bayer joint demosaicing and denoising based on dual encoder network with joint residual learning. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 38, pp. 7552–7561, 2024.
716 717	Ruofan Zhou, Radhakrishna Achanta, and Sabine Süsstrunk. Deep residual network for joint demosaicing and super-resolution. <i>arXiv preprint arXiv:1802.06573</i> , 2018.
718 719 720 721	Lianghui Zhu, Bencheng Liao, Qian Zhang, Xinlong Wang, Wenyu Liu, and Xinggang Wang. Vision mamba: Efficient visual representation learning with bidirectional state space model. <i>arXiv preprint arXiv:2401.09417</i> , 2024.
722	
723	
724	
726	
727	
728	
729	
730	
731	
732	
733	
734	
735	
736	
737	
738	
739	
740	
741	
742	
743	
745	
746	
747	
748	
749	
750	
751	
752	
753	
754	
755	

A MODEL DETAILS IN PYTORCH STYLE PSEUDO-CODE

In this section, we provide a detailed explanation of the modules we designed, including Cross-Swin State Block, Conv State Block, Quad Bayer Cross Swin Attention and Spatial Position Attention using PyTorch style pseudo-code. **The specific runnable model code is provided in the supplementary material.**

762 763

809

return image

760

761

756

```
764
      Algorithm 1 Pseudo-code of Cross-Swin State Block
765
       ## Cross-Swin State Block (CSSB)
766
       class CSSB(nn.Module):
767
          def __init__ (self, rvss_dim, trans_dim, head_dim, window_size, \
768
          drop_path, type='W', input_resolution=None):
769
               super(CSSB, self). init
                                          _()
               self.rvss_dim = rvss_dim
770
               self.trans_dim = trans_dim
771
               self.head_dim = head_dim
772
               # size of local window
773
               self.window_size = window_size
774
               # drop out
               self.drop_path = drop_path
775
               # type: 'W' (Window) or 'SW' (Shifted Window)
776
               self.type = type
777
               self.input_resolution = input_resolution
778
               assert self.type in ['W', 'SW']
779
               if self.input_resolution <= self.window_size:</pre>
                   self.type = 'W'
               # input projection
781
               self.conv1_1 = nn.Conv2d(self.rvss_dim+self.trans_dim, \
782
               self.rvss_dim+self.trans_dim, 1, 1, 0, bias=True)
783
               self.conv1_2 = nn.Conv2d(self.rvss_dim+self.trans_dim, \
784
               self.rvss_dim+self.trans_dim, 1, 1, 0, bias=True)
               # employ Quad Bayer Cross Swin Attention
785
               self.trans_block = QCSA(self.trans_dim, self.trans_dim, \ \
786
               self.window_size, self.drop_path, self.type, self.input_resolution)
787
               # Visual State Space Model
788
               self.rvss = VSSBlock(
789
                       hidden_dim=self.rvss_dim,
                       drop_path=0,
                       norm_layer=nn.LayerNorm,
791
                       attn_drop_rate=0,
792
                       d_state=16,
793
                       expand=2)
794
           def forward(self, image, position):
795
               # input projection and split
796
               rvss_x, trans_x = torch.split(self.conv1_1(image), \
797
               (self.rvss_dim, self.trans_dim), dim=1)
798
               # Residual Visual State Space Module
799
               rvss_x = self.rvss(rvss_x) + rvss_x
               # Rearrange inputs into 'b h w c'
800
               trans_x = Rearrange('b c h w -> b h w c')(trans_x)
801
               position = Rearrange('b c h w -> b h w c')(position)
802
               # Quad Bayer Cross Swin Attention
803
               trans_x = self.trans_block(trans_x, position)
804
               # Rearrange ouput back
              trans_x = Rearrange('b h w c -> b c h w')(trans_x)
805
               # concatenate and output projection
806
               res = self.conv1_2(torch.cat((rvss_x, trans_x), dim=1))
807
               # residual connection
808
               image = image + res
```

```
810
      Algorithm 2 Pseudo-code of Cross-Swin State Block
811
       ## Conv State Block (CSB)
812
       class CSB(nn.Module):
813
           def __init__(self, conv_dim, rvss_dim, head_dim, window_size, \
814
           drop_path, type='W', input_resolution=None):
815
               super(CSB, self).__init__()
816
               self.conv dim = conv dim
               self.rvss_dim = rvss_dim
817
               self.head_dim = head_dim
818
819
               # input projection
820
               self.conv1_1 = nn.Conv2d(self.conv_dim+self.rvss_dim, \
               self.conv_dim+self.rvss_dim, 1, 1, 0, bias=True)
821
               self.conv1_2 = nn.Conv2d(self.conv_dim+self.rvss_dim, \
822
               self.conv_dim+self.rvss_dim, 1, 1, 0, bias=True)
823
824
               # Visual State Space Model
825
               self.rvss = VSSBlock(
                       hidden_dim=self.rvss_dim,
826
                        drop_path=0,
827
                        norm_layer=nn.LayerNorm,
828
                        attn_drop_rate=0,
829
                        d_state=16,
830
                        expand=2)
831
               # Convolution Block
832
               self.conv_block = nn.Sequential(
833
                       nn.Conv2d(self.conv_dim, self.conv_dim, 3, 1, 1, bias=False),
834
                        nn.ReLU(True),
835
                        nn.Conv2d(self.conv_dim, self.conv_dim, 3, 1, 1, bias=False)
836
                        )
837
           def forward(self, x):
838
               # input projection and split
839
               conv_x, rvss_x = torch.split(self.conv1_1(x), \
840
               (self.conv_dim, self.rvss_dim), dim=1)
841
               # Residual Convolution
               conv_x = self.conv_block(conv_x) + conv_x
842
               # Residual Visual State Space Model
843
               rvss_x = self.rvss(rvss_x) + rvss_x
844
               # concatenate and output projection
845
               res = self.conv1_2(torch.cat((conv_x, rvss_x), dim=1))
846
               # residual connection
               x = x + res
847
               return x
848
849
850
851
852
853
854
855
856
```

860 861

857

862

```
864
       Algorithm 3 Pseudo-code of Quad Bayer Cross Swin Attention
865
       ## Quad Bayer Cross Swin Attention (QCSA)
866
       class QCSA(nn.Module):
867
           def __init__(self, input_dim, output_dim, head_dim, window_size, \
868
           drop_path, type='W', input_resolution=None):
869
               super(QCSA, self).__init__()
870
               self.input dim = input dim
               self.output_dim = output_dim
871
               # Swin type
872
               assert type in ['W', 'SW']
873
               self.type = type
874
               if input_resolution <= window_size:</pre>
875
                   self.type = 'W'
               # Layer Norm
876
               self.ln1 = nn.LayerNorm(input_dim)
877
               self.ln2 = nn.LayerNorm(input_dim)
878
               self.ln3 = nn.LayerNorm(input_dim)
879
               # Window Mulit Head Cross Attention
880
               self.msa = WMCA(input_dim, input_dim, head_dim, window_size, self.type)
881
882
               # drop out
883
               self.drop_path = DropPath(drop_path) if drop_path > 0. else nn.Identity()
884
885
               # Multi Layer Perceptron
               self.mlp = nn.Sequential(
886
                   nn.Linear(input_dim, 4 * input_dim),
887
                   nn.GELU(),
888
                   nn.Linear(4 * input_dim, output_dim),
889
               )
890
           def forward(self, image, position):
891
               # Cross Swin Layer + residual connection
892
               image = image + self.drop_path(self.msa(self.ln1(image), self.ln3(position)))
893
894
               # MLP for out projection
               fused = image + self.drop_path(self.mlp(self.ln2(image)))
895
               return fused
896
897
898
```

```
17
```

```
Algorithm 4 Pseudo-code of Window Multi Head Cross Attention
919
       ## Window Multi Head Cross-attention module in Quad Bayer Cross Swin Attention
920
       class WMCA(nn.Module):
921
          def
                _init__(self, input_dim, output_dim, head_dim, window_size, type):
922
               super(WMCA, self).__init_
                                          _()
923
               self.input_dim = input_dim
924
               self.output_dim = output_dim
               self.head_dim = head_dim
925
               self.n_heads = input_dim//head_dim
926
               # scale factor
927
               self.scale = self.head_dim ** -0.5
928
               # size of local window
929
               self.window_size = window_size
               # Swin type: 'W' (Window) or 'SW' (Shifted Window)
930
               self.type=type
931
               # positional encoding inside the window
932
               self.relative_position_params = nn.Parameter(torch.zeros((2 *\
933
               window_size - 1) * (2 * window_size -1), self.n_heads))
934
               trunc_normal_(self.relative_position_params, std=.02)
               self.relative_position_params = torch.nn.Parameter(self.relative_position_params.view)
935
               window_size-1, 2*window_size-1, self.n_heads).transpose(1,2).transpose(0,1))
936
               # input projection
937
               self.embedding_layer = nn.Linear(self.input_dim, 2*self.input_dim, bias=True)
938
               self.embedding_layer_qb = nn.Linear(self.input_dim, self.input_dim, bias=True)
939
               # output projection
               self.linear = nn.Linear(self.input_dim, self.output_dim)
940
           def forward(self, x, y):
941
               # x: input image tensor with shape of [b h w c];
942
               # y: input Quad Bayer tensor with shape of [b h w c];
943
               # shift window when type = 'SW'
944
               if self.type!='W':
                   x = torch.roll(x, shifts=(-(self.window_size//2), -(self.window_size//2)), dims=(1)
945
                   y = torch.roll(y, shifts=(-(self.window_size//2), -(self.window_size//2)), dims=(1)
946
               # patrition to windows
947
               x = rearrange(x, 'b (w1 p1) (w2 p2) c \rightarrow b w1 w2 p1 p2 c', \
948
                   p1=self.window_size, p2=self.window_size)
949
               h_windows = x.size(1)
               w_windows = x.size(2)
950
               x = rearrange(x, 'b w1 w2 p1 p2 c \rightarrow b (w1 w2) (p1 p2) c', \
951
                   p1=self.window_size, p2=self.window_size)
952
               y = rearrange(y, 'b (w1 p1) (w2 p2) c -> b (w1 w2) (p1 p2) c', \
953
                   p1=self.window_size, p2=self.window_size)
954
               # input projection
               kv = self.embedding_layer(x)
955
               k, v = rearrange(kv, 'b nw np (threeh c) -> threeh b nw np c', \
956
                   c=self.head_dim).chunk(2, dim=0)
957
               q = self.embedding_layer_qb(y)
958
               q = rearrange(q, 'b nw np (h c) -> h b nw np c', c=self.head_dim)
959
               # cross window attention
               sim = torch.einsum('hbwpc,hbwqc->hbwpq', q, k) * self.scale
960
               sim = sim + rearrange(self.relative_embedding(), 'h p q -> h 1 1 p q')
961
               # masks when shifted window
962
               if self.type != 'W':
963
                   attn_mask = generate_mask(h_windows, w_windows, \
                       self.window_size, shift=self.window_size//2)
964
                   sim = sim.masked_fill_(attn_mask, float("-inf"))
965
               # attention map
966
               probs = nn.functional.softmax(sim, dim=-1)
967
                # attention to value
968
               output = torch.einsum('hbwij,hbwjc->hbwic', probs, v)
               output = rearrange(output, 'h b w p c -> b w p (h c)')
969
               # output projection
970
               output = self.linear(output)
971
               output = rearrange(output, 'b (w1 w2) (p1 p2) c -> b (w1 p1) (w2 p2) c', \
                   w1=h_windows, p1=self.window_size)
               # shift back when type = 'SW_{18}
               if self.type!='W': output = torch.roll(output, shifts=\
                   (self.window_size//2, self.window_size//2), dims=(1,2))
               return output
```

```
Algorithm 5 Pseudo-code of Spatial Position Attention
973
       ## Spatial Position Attention (SPA)
974
       class SPA(nn.Module):
975
           def __init__(self, dim, num_heads, bias=False, LayerNorm_type='WithBias'):
976
                super(SPA, self).__init__()
977
                # Layer Norm
                self.norm1_image = LayerNorm(dim, LayerNorm_type)
978
                self.norm1_position = LayerNorm(dim, LayerNorm_type)
979
                # Attention
980
                self.num_heads = num_heads
981
                self.temperature = nn.Parameter(torch.ones(dim, 1, 1))
982
                self.q = nn.Conv2d(dim, dim, kernel_size=1, bias=bias)
self.v = nn.Conv2d(dim, dim, kernel_size=1, bias=bias)
983
984
           def forward(self, image, position):
985
                # image: b, c, h, w
986
                # position: b, c, h, w
987
                # return: b, c, h, w
988
                # projection
989
                q = self.q(self.norm1_image(image)) # qb
990
                v = self.v(self.norm1_position(position)) # image
991
992
                # position attention
                x_spatial = F.relu(q) * v * self.temperature
993
994
                # residual connection
995
                fused = image + x_spatial
996
997
                return fused
998
```