# Cross Paradigm Representation and Alignment Transformer for Image Deraining

Shun Zou
Nanjing Agricultural University &
Soochow University

Yi Zou
Xiangtan University

Juncheng Li
Shanghai University

Guangwei Gao *
Nanjing University of Posts and
Telecommunications & Soochow
University

Guojun Qi
Westlake University

## Abstract

Transformer-based networks have achieved strong performance in low-level vision tasks like image deraining by utilizing spatial or channel-wise self-attention. However, irregular rain patterns and complex geometric overlaps challenge single-paradigm architectures, necessitating a unified framework to integrate complementary global-local and spatial-channel representations. To address this, we propose a novel Cross Paradigm Representation and Alignment Transformer (CPRAformer). Its core idea is the hierarchical representation and alignment, leveraging the strengths of both paradigms (spatial-channel and global-local) to aid image reconstruction. It bridges the gap within and between paradigms, aligning and coordinating them to enable deep interaction and fusion of features. Specifically, we use two types of self-attention in the Transformer blocks: sparse prompt channel self-attention (SPC-SA) and spatial pixel refinement self-attention (SPR-SA). SPC-SA enhances global channel dependencies through dynamic sparsity, while SPR-SA focuses on spatial rain distribution and fine-grained texture recovery. To address the feature misalignment and knowledge differences between them, we introduce the Adaptive Alignment Frequency Module (AAFM), which aligns and interacts with features in a two-stage progressive manner, enabling adaptive guidance and complementarity. This reduces the information gap within and between paradigms. Through this unified cross-paradigm dynamic interaction framework, we achieve the extraction of the most valuable interactive fusion information from the two paradigms. Extensive experiments demonstrate that our model achieves state-of-the-art performance on eight benchmark datasets and further validates CPRAformer's robustness in other image restoration tasks and downstream applications.

## CCS Concepts

• **Computing methodologies → Reconstruction**.

## Keywords

Image Restoration, Single Image Deraining, Self-Attention, Transformer, Cross-Paradigm Dynamic Interaction

## 1 Introduction

Single Image Deraining (SID) is a traditional low-level vision task that aims to restore a clear, high-quality image from a given rainy

---

*Corresponding author. Email: zs@stu.njau.edu.cn, csgwgao@njupt.edu.cn.
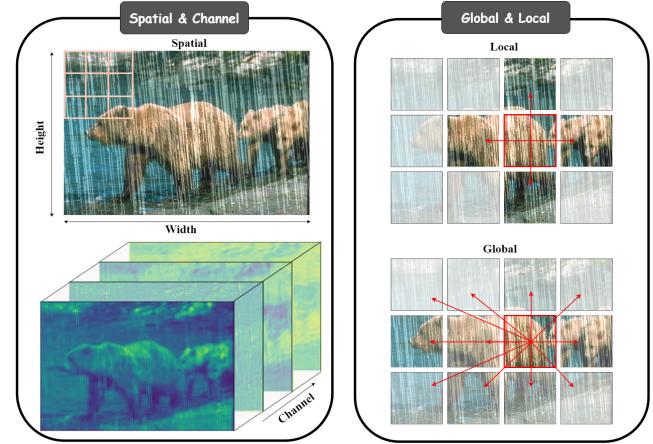


**Figure 1: Feature patterns obtained from four perspectives are distinct, two deraining paradigms offers unique advantages. Recent deraining research mainly focuses on spatial-channel or global-local paradigms, lacking a framework that effectively integrates these two paradigms.**

image. As it plays a critical role in downstream tasks across various fields, including video surveillance, autonomous driving, and medical imaging, it has garnered increasing attention from both academia and industry. Due to its ill-posed nature, early methods typically applied various priors based on the statistical characteristics of rain streaks and clean images [32, 74]. However, in complex and diverse rainy scenarios, such priors do not always hold.

Recently, many studies have proposed convolutional neural network (CNN)-based methods to address this challenge [25, 30, 31, 45, 68, 72]. However, due to the limited receptive field of convolution operators, long-range spatial modeling is hindered, which restricts model performance. Fortunately, inspired by the success of Transformers in natural language processing and advanced vision tasks [17, 51, 53, 60], researchers have developed Transformer-based architectures for the SID task [6, 52, 61]. Leveraging the self-attention mechanism, Transformer-based methods can establish global dependencies, alleviating the limitations of CNN-based approaches and demonstrating superior deraining performance. Recognizing the potential of Transformers, some researchers have explored their effective application in SID tasks from different perspectives. In terms of spatial modeling, some methods use non-overlapping spatial windows to capture more global dependencies, enhancing

spatial pixel modeling [61]. Regarding channels, "transpose" attention has been proposed [71], where self-attention is computed along the channel dimension instead of the spatial dimension. These methods, with strong feature extraction and modeling capabilities in their respective dimensions, have achieved remarkable results. Intuitively, extracting spatial features and capturing channel context information both play crucial roles in enhancing Transformer performance in image restoration [10]. Additionally, due to the intricate intertwining of rain streaks and the rain-free background, both global and local features are essential for the challenging SID task. However, the self-attention mechanism in Transformers does not fully leverage the local invariance of CNNs. To address this, some researchers have attempted to combine CNNs with Transformers [6, 8], inheriting CNNs' advantage in local modeling and Transformers' strength in capturing long-range dependencies.

However, two key questions naturally arise: (1) How can we simultaneously leverage information from all paradigms? As shown in Fig. 1, representations from different perspectives undoubtedly play a critical role in SID task performance; (2) How can we effectively align and aggregate the two feature types within each paradigm? The distinct feature models across paradigms make simple summation or concatenation prone to information loss, failing to significantly boost performance. To address these issues, we propose a novel hybrid architecture: the Cross Paradigm Representation and Alignment Transformer (CPRAformer). Its core idea is to establish a cross-paradigm representation learning framework through dimensional consistency (spatial-channel perspective) and multi-perspective integration (global-local perspective), along with alignment and hierarchical fusion of corresponding features in each paradigm. Specifically, it comprises two carefully designed components: Cross-Paradigm Interaction and Alignment Self-Attention (CPIA-SA) and Multi-Scale Flow Gating Network (MSGN).

In CPIA-SA, we introduce two types of self-attention: Sparse Prompt Channel Self-Attention (SPC-SA) and Spatial Pixel Refinement Self-Attention (SPR-SA). SPC-SA computes attention along the channel dimension and dynamically filters attention values in the dense attention matrix using prompt information. This allows the network to exploit sparsity, retaining the most valuable attention information while minimizing excessive noise interactions that could degrade image restoration quality, thus effectively extracting global-channel information. Conversely, SPR-SA utilizes an efficient CNN-based architecture to approximate self-attention, enabling the effective modeling of local fine-grained features and the relationships between neighboring spatial pixels, fully leveraging local spatial characteristics. Additionally, these two self-attention mechanisms are complementary. SPC-SA provides global information between features for SPR-SA, thereby expanding the receptive field of pixels. Meanwhile, SPR-SA enhances the spatial representation of each feature map, which aids in modeling channel context.

At the same time, to further promote alignment and interaction within each paradigm, we propose a two-stage progressive fusion strategy called the Adaptive Alignment Frequency Module (AAFM). Using adaptive weighting, it aligns corresponding branches and enhances interactions between frequency spectra to aggregate and strengthen internal feature information. Moreover, another key component of the Transformer module is the feed-forward network (FFN) [17], which typically extracts features through fully connected layers but often overlooks the critical multi-scale information needed for SID tasks [1, 6]. To address this limitation, we introduce the Multi-Scale Flow Gating Network (MSGN). Utilizing a gating mechanism, MSGN incorporates multi-scale representation learning, providing additional nonlinear information to the FFN and enhancing its capability to capture essential features.

In summary, through the design described above, our CPRAformer achieves cross-paradigm feature pattern learning and information alignment, thereby enabling robust feature representation. Our main contributions are summarized as follows:

- We propose a hybrid model for SID, CPRAformer, which integrates the advantages of spatial-channel and global-local paradigms. Through cross-paradigm dynamic interaction, it aligns and adaptively fuses feature patterns between the two paradigms.
- We utilize both SPC-SA and SPR-SA to extract features across spatial and channel dimensions, effectively modeling global dependencies while capturing local details for complementary feature integration.
- To bridge the feature gap between paradigms, we develop the AAFM, which first performs feature alignment through adaptive weighting and then achieves feature fusion via frequency-domain interaction. This facilitates better coordination and deep interaction between different types of information. In addition, the MSGN is also included to learn scale-aware spatial features.

## 2 Related works

### 2.1 Single Image Deraining

Traditional rain removal methods often rely on handcrafted priors [22, 28, 32, 40, 74], which are subjective and unable to adapt to complex rainy scenes. To address this issue, many researchers have developed CNN-based methods [25, 30, 31, 45, 55, 67, 68, 70, 72] for image deraining, achieving promising results. However, convolutions struggle to capture long-range dependencies in both spatial and channel dimensions. Inspired by the success of Transformers in advanced vision tasks [17, 51, 53, 60], they have also been applied to image deraining [6, 35, 44, 52, 54, 57, 60, 61, 63]. Transformers, as a new network backbone, show significant improvements over CNN-based methods due to their excellent global context awareness enabled by self-attention. However, a limitation of self-attention is that tokens with low attention values or irrelevant tokens can interfere with the dense attention matrix, potentially harming output features [6, 49, 65]. As shown in Fig. 3, to overcome this, we propose an adaptive sparse attention mechanism, which dynamically adjusts the sparse range using a learnable operator. This approach maximizes network sparsity and reduces excessive interference from irrelevant noise in naive self-attention.

### 2.2 Feature Aggregation

**Global & Local.** To leverage the advantage of CNNs in extracting local features and Transformers in global modeling, many hybrid models have been proposed. For example, ELF [23] was the first to unify these architectures into a lightweight deraining model based on association learning. Inspired by progressive learning, HCT-FFN [8] introduced a new staged hybrid deraining network.
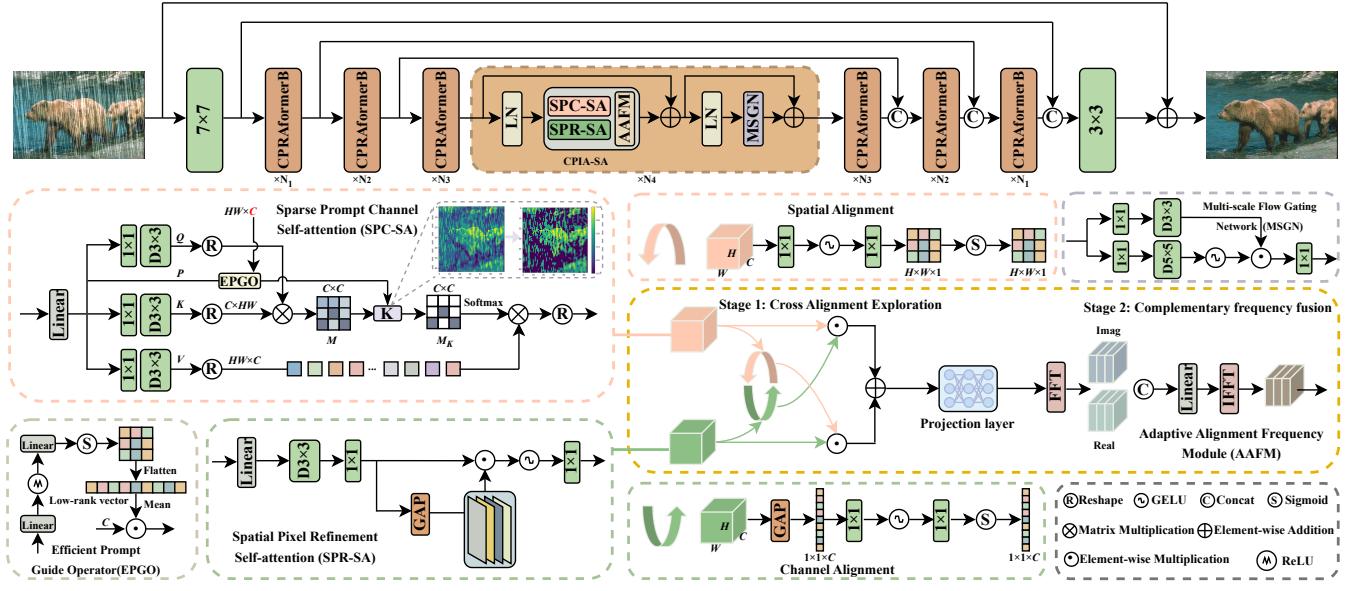
**Figure 2: The overall architecture of our proposed CPRAformer.**

SMFANet [78] uses adaptive feature aggregation to synergize local and non-local feature interactions. Dual-former [5] combines the global modeling power of self-attention with the local capability of convolutions in a unified architecture. It uses a hybrid Transformer block to model long-range spatial dependencies and handle uneven channel distributions.

**Spatial & Channel.** In CNNs, researchers apply attention mechanisms along the spatial and channel dimensions to enhance feature representation, as demonstrated by models like RESCAN [31], MSPFN [25], and MPRNet [72]. In Transformer-based approaches, spatial self-attention is primarily used to capture long-range dependencies between pixels; for instance, IDT [64] uses dual Transformers with window and spatial attention for deraining. Additionally, some studies integrate channel self-attention in Transformers to combine spatial and channel information. Notably, Restormer [71] designs an efficient Transformer model by estimating self-attention along the channel dimension, achieving significant performance gains, while DRSFormer [6] proposes a sparse Transformer along the channel dimension to fully exploit the most informative features for deraining. DAT [10] alternates between spatial and channel self-attention mechanisms across consecutive Transformer blocks, aggregating features from the spatial and channel dimensions both across and within blocks.

**Motivation.** Extensive research on both paradigms, namely global-local as well as spatial-channel, demonstrates the critical role of feature representation in complex rainy scenes. However, no study has simultaneously considered both paradigms (see Fig. 1). Therefore, we introduce two types of self-attention mechanisms into image deraining for comprehensive information exploration and representation. Additionally, we design AAFM to dynamically achieve both inter-paradigm and intra-paradigm feature aggregation, iteratively refining feature coherence across scales and dimensions.

## 3 Method

### 3.1 Overall Pipeline

The overall pipeline of our proposed CPRAformer is illustrated in Fig. 2, which adopts an encoder-decoder architecture with skip connections. Specifically, given an input rain image $I_{\text{rain}} \in \mathbb{R}^{H \times W \times 3}$, we first apply a $7 \times 7$ convolution to obtain the low-level feature embedding $F_0 \in \mathbb{R}^{H \times W \times C}$, where $H$, $W$, and $C$ represent the height, width, and channels, respectively. The low-level feature embedding is then fed into the backbone network, which consists of a 4-level encoder-decoder structure. In the encoder stage, the resolution of the high-resolution input is reduced by a factor of 2 while increasing the number of channels, whereas in the decoder stage, the process is reversed. Each encoder and decoder is composed of $N_i$ ( $i = 1, 2, 3, 4$) stacked Cross Paradigm Representation and Alignment Transformer Block (CPRAformerBs). Within each CPRAformerB, we develop the CPIA-SA to extract and aggregate features from both spatial-channel and global-local paradigms simultaneously. Additionally, we incorporate MSGN within each CPRAformerB, which leverages an elegant gating mechanism to extract multi-scale information, aiding in the image-deraining process.

As shown in Fig. 2, CPIA-SA is composed of three main components: Sparse Prompt Channel Self-Attention (SPC-SA), Spatial Pixel Refinement Self-Attention (SPR-SA), and Frequency Adaptive Interaction Module (AAFM). SPC-SA, using the Efficient Prompt Guide Operator (EPGO), effectively explores the sparsity of the neural network, adaptively retaining the most valuable attention values. This enables efficient extraction of rain-degraded features that handle spatial variations while modeling global channel context. In contrast, SPR-SA focuses on modeling the spatial background, enhancing the spatial pixel representation of each feature map, and facilitating accurate background restoration through fine-grained local features. To bridge the gap between SPC-SA and SPR-SA in terms of spatial-channel and global-local knowledge, and to fully
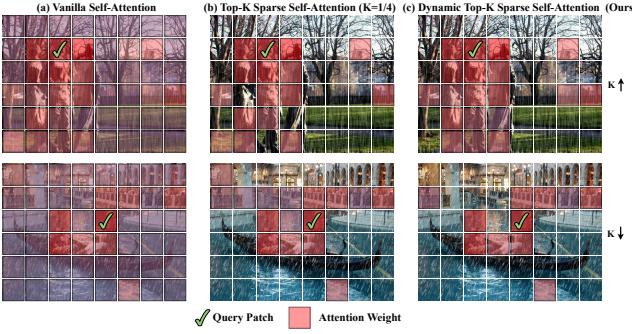
Figure 3: Comparison of different self-attention mechanisms. (a) The naive self-attention mechanism [71] computes and retains all tokens. (b) The Top-K sparse attention mechanism [6] sets a fixed K value (here, K is set to 1/4) and retains only the top K% tokens with the highest attention values while setting the remaining tokens to zero. (c) Our dynamic Top-K sparse attention mechanism adaptively modulates the K value based on input features. For instance, compared to the fixed K in (b), the K value increases in the upper image and decreases in the lower image to adapt to different images.

integrate features within both paradigms, we introduce AAFM. AAFM utilizes a dual-stage progressive strategy to alignment and fuse paradigm-specific features comprehensively.

## 3.2 Sparse Prompt Channel Self-Attention

As shown in Fig. 2, the self-attention mechanism in SPC-SA operates along the channel dimension. Specifically, given the input embedding feature $F \in \mathbb{R}^{H \times W \times C}$, point-wise convolution (PW-Conv) and 3×3 depth-wise convolution (DWConv) are applied to $F$ to aggregate cross-pixel channel information, generating the matrices for query $Q$, key $K$, and value $V$. Next, we perform a dot-product operation on the reshaped $Q$ and $K$ to generate a dense attention matrix $M \in \mathbb{R}^{C \times C}$. However, we observe that the tokens in the keys are not always relevant to those in the queries, and the self-attention values estimated using irrelevant tokens introduce noise interactions and information redundancy, affecting the quality of image recovery. To address this, we introduce a Top-k mechanism, which differs from traditional self-attention mechanisms, to filter the information in the attention matrix, retaining the most significant attention values and avoiding noise that can lead to artifacts in the deraining process. For example, with $k = 4/5$, we only retain the top 80% of attention scores, while the remaining elements are masked to zero. Notably, we also develop a novel learnable operator: the Efficient Prompt Guide Operator (EPGO), which dynamically generates prompt information based on the input, guiding the $K$ values to achieve adaptive modulation and facilitating a dynamic selection process in the attention matrix. As shown in Fig. 3, unlike previous studies [6, 71], we propose a novel dynamic sparse mechanism that fully exploits the sparsity of neural networks. Formally, the above process is expressed as follows:

$$\Omega_i^{(k)} = \arg \max_{S \subset \{1,\dots,N\}} \sum_{j \in S} M_{ij}, \tag{1}$$

$$\Omega^{(k)} = \left\{ \Omega_1^{(k)}, \Omega_2^{(k)}, \dots, \Omega_N^{(k)} \right\}, \tag{2}$$

---

**Algorithm 1** Efficient Prompt Guide Operator (EPGO)

---

**Input:** $F \in \mathbb{R}^{H \times W \times C}$                    ▷ Input feature map
**Output:** *Dynamic K*
1: $F_{\text{norm}} \leftarrow \text{LayerNorm}(F)$   ▷ Normalize the input feature map
2: $X \leftarrow \text{Linear}_1(F_{\text{norm}})$       ▷ Project $F_{\text{norm}}$ to a hidden space
3: $X \leftarrow \text{ReLU}(X)$
4: $F_{PG} \leftarrow \text{Linear}_2(X)$              ▷ Generate prompt features
5: $F_{PG} \leftarrow \text{Sigmoid}(F_{PG})$         ▷ Constrain values to $[0, 1]$
6: $F_{PG}^{\text{flat}} \leftarrow \text{Flatten}(F_{PG})$      ▷ Generate low-rank vector.
7: $p \leftarrow \frac{1}{H \cdot W \cdot C} \sum_{i=1}^{H \cdot W \cdot C} F_{PG}^{\text{flat}}[i]$       ▷ Compute global average
8: Dynamic $K \leftarrow C \times p$ ▷ Compute global scaling factor, where $C$ is the channel count
9: **return** *Dynamic K*

---

$$[M_k]_{ij} = \begin{cases} 1, & \text{if } j \in \Omega_i^{(k)} \iff M_{ij} \in \text{Top}_k(M_{i,:}), \\ 0, & \text{otherwise.} \end{cases} \tag{3}$$

Here, $\Omega_i^{(k)}$ represents the index set of the top-k elements in the i-th row, that is, the column indices of the k most important elements in the i-th row of matrix $M$. Formally, the process of SPC-SA is expressed as

$$SPC\text{-}SA = Softmax(T_k(\frac{QK^T}{\sqrt{d}}))V, \tag{4}$$

where $T_k(\cdot)$ represents the Top-k selection operation after modulation by the Efficient Prompt Guide Operator., and $\sqrt{d}$ represents an optional temperature coefficient used to control the magnitude of the dot product between $Q$ and $K$ before applying softmax. Similar to most previous works [17], we employ a multi-head strategy, concatenating all outputs of the multi-head attention and then obtaining the final result through a linear projection.

**Efficient Prompt Guide Operator.** Since a fixed $k$ leads to a rigid pattern structure that cannot dynamically update with the neural network, it struggles to adapt to the complex rain conditions in real-world scenarios. To address this, we propose the EPGO, which provides the neural network with dynamic prompt information. This guides the attention matrix towards optimal selection, allowing the model to optimize attention allocation in both sparse and dense attention scenarios, preserving the most valuable information in the current features. As shown in Fig. 2, the architecture of EPGO works as follows: given the input feature $F \in \mathbb{R}^{H \times W \times C}$, two linear layers with hidden ReLU activation and a Sigmoid function are first used to generate the prompt guide features. These features are then flattened into a low-dimensional vector. Finally, the low-dimensional vector is averaged and multiplied element-wise with the channel $C$ of the input feature $F$, adapting to the dimensional changes in the features and effectively setting a soft threshold for the $K$ values. In Algorithm 1, we outline the process of dynamically tuning $K$ using EPGO.

## 3.3 Spatial Pixel Refinement Self-Attention

Unlike SPC-SA, the motivation of SPR-SA is to efficiently model spatial information for each pixel and effectively extract local details. However, existing spatial self-attention mechanisms often come with high computational costs and have limited capability in modeling local details [61]. To address this, we propose a simple

**Table 1: Comparison of quantitative results on five datasets. Bold and underlined indicate the best and second-best results.**

| Method | Year | Test100 [77] PSNR ↑ | SSIM ↑ | Rain100H [66] PSNR ↑ | SSIM ↑ | Rain100L [66] PSNR ↑ | SSIM ↑ | Test2800 [18] PSNR ↑ | SSIM ↑ | Test1200 [75] PSNR ↑ | SSIM ↑ | Average PSNR ↑ | SSIM ↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RESCAN [31] | ECCV2018 | 21.59 | 0.726 | 18.01 | 0.467 | 24.15 | 0.791 | 24.50 | 0.765 | 24.40 | 0.759 | 22.53 | 0.702 |
| PReNet [45] | CVPR2019 | 23.17 | 0.752 | 17.63 | 0.487 | 27.76 | 0.876 | 27.20 | 0.825 | 26.05 | 0.792 | 24.36 | 0.746 |
| SPDNet [68] | ICCV2021 | 24.25 | 0.848 | 25.87 | 0.809 | 28.63 | 0.880 | 31.05 | 0.904 | 30.42 | 0.893 | 28.04 | 0.867 |
| PCNet [26] | TIP2021 | 23.29 | 0.762 | 20.83 | 0.563 | 26.64 | 0.817 | 27.10 | 0.818 | 26.53 | 0.791 | 24.88 | 0.750 |
| MPRNet [72] | CVPR2021 | 25.66 | 0.859 | 28.23 | 0.850 | 31.94 | 0.930 | 32.14 | 0.925 | 31.32 | 0.901 | 29.86 | 0.893 |
| HINet [3] | CVPRW2021 | 23.21 | 0.767 | 20.85 | 0.598 | 27.03 | 0.842 | 28.36 | 0.843 | 27.77 | 0.821 | 25.44 | 0.774 |
| DANet [24] | IJCAI2022 | 23.96 | 0.839 | 23.00 | 0.791 | 29.51 | 0.906 | 30.32 | 0.903 | 29.99 | 0.888 | 27.36 | 0.865 |
| Uformer [61] | CVPR2022 | 23.87 | 0.815 | 22.43 | 0.700 | 28.39 | 0.883 | 29.71 | 0.886 | 28.65 | 0.856 | 26.61 | 0.828 |
| ALformer [23] | ACMMM2022 | 24.41 | 0.844 | 25.10 | 0.807 | 29.39 | 0.903 | 31.36 | 0.916 | 30.40 | 0.897 | 28.13 | 0.874 |
| NAFNet [2] | ECCV2022 | 25.75 | 0.845 | 26.76 | 0.813 | 31.27 | 0.925 | 31.71 | 0.918 | 30.62 | 0.892 | 29.22 | 0.879 |
| MIRNetV2 [73] | TPAMI2022 | 25.76 | 0.867 | 28.05 | 0.846 | 32.53 | 0.935 | 32.33 | 0.925 | 32.38 | 0.915 | 30.21 | 0.897 |
| MFDNet [58] | TIP2023 | 25.90 | 0.870 | 27.06 | 0.850 | 32.76 | 0.944 | 31.92 | 0.925 | 31.15 | 0.909 | 29.76 | 0.899 |
| HCT-FFN [8] | AAAI2023 | 24.86 | 0.847 | 26.70 | 0.819 | 29.94 | 0.906 | 31.46 | 0.915 | 31.23 | 0.901 | 28.84 | 0.878 |
| DRSformer [6] | CVPR2023 | 27.86 | 0.885 | 28.16 | 0.864 | 34.79 | 0.954 | 32.80 | <u>0.931</u> | 30.99 | 0.906 | 30.92 | 0.908 |
| ChaIR [11] | KBS2023 | 28.19 | 0.879 | 28.69 | 0.862 | 34.52 | 0.953 | <u>32.85</u> | <u>0.931</u> | 31.30 | 0.903 | 31.11 | 0.906 |
| IRNeXT [14] | ICML2023 | 25.80 | 0.860 | 27.22 | 0.833 | 31.65 | 0.931 | 30.53 | 0.917 | 29.02 | 0.898 | 28.85 | 0.888 |
| OKNet [13] | AAAI2024 | 25.43 | 0.858 | 24.01 | 0.804 | 31.19 | 0.928 | 29.32 | 0.911 | 27.56 | 0.886 | 27.50 | 0.877 |
| AST [79] | CVPR2024 | 26.07 | 0.859 | 27.40 | 0.833 | 32.03 | 0.932 | 31.65 | 0.921 | 30.69 | 0.897 | 29.57 | 0.889 |
| SFHformer [27] | ECCV2024 | 25.67 | 0.856 | 27.25 | 0.832 | 32.97 | 0.944 | 32.27 | 0.925 | 31.50 | 0.904 | 29.94 | 0.892 |
| Nerd-rain [7] | CVPR2024 | 27.16 | 0.869 | 28.07 | 0.838 | 33.72 | 0.949 | 32.63 | 0.927 | 30.45 | 0.890 | 30.41 | 0.895 |
| MSDT [1] | AAAI2024 | 27.79 | 0.878 | 29.05 | 0.856 | 34.75 | 0.955 | 32.68 | 0.930 | **32.12** | **0.917** | 31.28 | 0.907 |
| FSNet [12] | TPAMI2024 | 27.95 | 0.884 | 28.70 | 0.860 | 34.10 | 0.952 | 32.68 | <u>0.931</u> | 31.26 | 0.910 | 30.94 | 0.908 |
| AdaIR [15] | ICLR2025 | <u>28.64</u> | <u>0.889</u> | <u>29.48</u> | <u>0.871</u> | <u>35.84</u> | <u>0.962</u> | 32.70 | 0.930 | 30.58 | 0.907 | <u>31.45</u> | <u>0.912</u> |
| **CPRAformer (Ours)** | – | **29.65** | **0.895** | **29.68** | **0.875** | **35.98** | **0.964** | **33.00** | **0.933** | <u>31.52</u> | <u>0.913</u> | **31.97** | **0.916** |

**Table 2: Quantitative evaluations on Raindrop dataset [41]. Bold and underlined indicate the best and second-best results.**

| Method | Year | Raindrop-A [41] PSNR ↑ | SSIM ↑ | Raindrop-B [41] PSNR ↑ | SSIM ↑ |
|---|---|---|---|---|---|
| RESCAN [31] | ECCV2018 | 25.09 | 0.837 | 22.55 | 0.727 |
| PReNet [45] | CVPR2019 | 25.61 | 0.884 | 22.99 | 0.787 |
| SPDNet [68] | ICCV2021 | 28.82 | 0.896 | 24.89 | 0.792 |
| PCNet [26] | TIP2021 | 25.68 | 0.837 | 22.89 | 0.726 |
| MPRNet [72] | CVPR2021 | 29.96 | 0.916 | 25.58 | 0.815 |
| HINet [3] | CVPRW2021 | 25.81 | 0.882 | 23.17 | 0.787 |
| DANet [24] | IJCAI2022 | 29.54 | 0.914 | 25.28 | 0.812 |
| Uformer [61] | CVPR2022 | 28.99 | 0.903 | 25.02 | 0.803 |
| ALformer [23] | ACMMM2022 | 29.11 | 0.911 | 25.11 | 0.809 |
| NAFNet [2] | ECCV2022 | 29.81 | 0.907 | 25.33 | 0.806 |
| MFDNet [58] | TIP2023 | 28.57 | 0.882 | 24.53 | 0.766 |
| HCT-FFN [8] | AAAI2023 | 28.09 | 0.891 | 24.48 | 0.791 |
| DRSformer [6] | CVPR2023 | 30.83 | 0.923 | 25.86 | 0.819 |
| ChaIR [11] | KBS2023 | 30.88 | <u>0.925</u> | 25.84 | <u>0.820</u> |
| IRNeXT [14] | ICML2023 | 30.69 | 0.924 | 25.79 | 0.819 |
| OKNet [13] | AAAI2024 | 30.39 | 0.924 | 25.65 | 0.818 |
| SFHformer [27] | ECCV2024 | 23.09 | 0.869 | 21.23 | 0.772 |
| Nerd-rain [7] | CVPR2024 | 30.96 | 0.924 | 25.96 | 0.819 |
| MSDT [1] | AAAI2024 | 30.85 | 0.922 | 25.89 | 0.818 |
| FSNet [12] | TPAMI2024 | 30.83 | <u>0.925</u> | <u>25.99</u> | 0.819 |
| AdaIR [15] | ICLR2025 | <u>30.99</u> | 0.924 | 25.97 | 0.817 |
| **CPRAformer (Ours)** | – | **31.19** | **0.926** | **26.01** | **0.821** |

yet efficient self-attention approximation module using convolutions. By capturing the positional information of each pixel, the features are reweighted, allowing each pixel to perceive different degradation signals from the same position across all channels. This approach enhances the model's ability to handle spatial variations in deraining tasks. The specific operations are illustrated in Fig. 2, and the computation is formulated as follows:

$$F_L = PW(DW^{3\times3}(Linear(F))), \quad (5)$$

$$F_{SP} = GAP(F_L), \quad (6)$$

$$F' = PW(\varphi(F_{SP} \odot F_L)), \quad (7)$$

where $PW(\cdot)$ denotes point-wise convolution, $DW(\cdot)^{x\times x}$ represents $x \times x$ depth-wise convolutions, $GAP(\cdot)$ stands for global average pooling, and $\varphi(\cdot)$ refers to the GELU function.

## 3.4 Adaptive Alignment Frequency Module

Although SPC-SA and SPR-SA capture global channel features and local spatial features respectively, effectively integrating the information from these two branches becomes a critical challenge. An intuitive observation is that there exists an uncertain knowledge gap between the convolution-based local features from CNN and the self-attention-based global features from Transformer, as well as between spatial and channel dimension features [8, 10]. Thus, simply concatenating or adding these features cannot fully maximize their potential. To address this issue, we propose the AAFM. AAFM adopts a two-stage process that progressively integrates the features from both paradigms layer by layer. First, based on the types of features from the two branches, AAFM adaptively reweights the features along either the spatial or channel dimensions to align the first paradigm (i.e., spatial-channel). Then, we introduce the features into the frequency domain, leveraging the Fourier transform to enhance the interaction across multiple frequency spaces. This allows each pixel to perceive patterns from other pixels, achieving the aggregation of global features into local features and the diffusion of local features into global ones, thus realizing the deep interaction and fusion of the second paradigm (i.e., global-local). Specifically, given the input features $F_{spc}$ and $F_{spr}$ from the two branches, AAFM first applies two interaction operations: Spatial Alignment and Channel Alignment, which generate spatial attention maps and channel attention maps, respectively. These are then reweighted onto the corresponding branch features to achieve effective alignment and interaction. The process can be expressed as follows:

$$Map_S = f(PW(\varphi(PW(F_{SPC})))), \quad (8)$$

$$Map_C = f(PW(\varphi(PW(GAP(F_{SPR}))))), \quad (9)$$

$$\hat{F} = F_{SPR} \odot Map_S + F_{SPC} \odot Map_C. \quad (10)$$

To leverage feature differences in the frequency domain and bridge information gaps, we apply the Fast Fourier Transform (FFT) to the

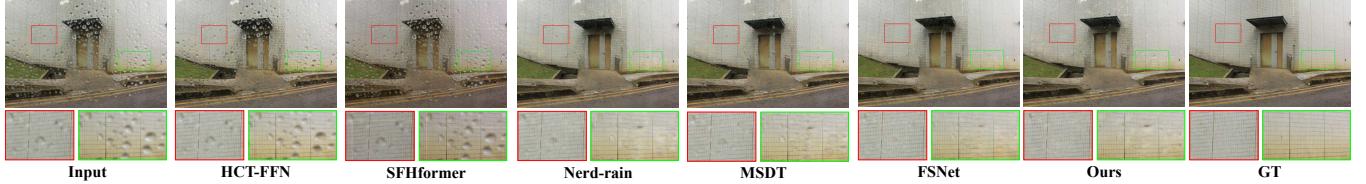**Figure 4: The qualitative comparison on Test100 [77]. See the supplements for more visualizations.**



**Figure 5: The qualitative comparison on raindrop datasets [41]. Our result has the best visual quality and details.**



**Figure 6: Comparison of visual results on a real-world dataset [59]. Our result has the best visual quality and details.**

fused feature $\hat{F}$. For simplicity, let us first consider a single-channel case, $\hat{F} \in \mathbb{R}^{H \times W}$. The 2D Fast Fourier Transform of $\hat{F}$ is defined as:

$$\mathcal{F}(\hat{F})(u, v) = \frac{1}{\sqrt{HW}} \sum_{h=0}^{H-1} \sum_{w=0}^{W-1} \hat{F}(h, w) e^{-j2\pi\left(\frac{uh}{H} + \frac{vw}{W}\right)}, \qquad (11)$$

where $u$ and $v$ denote the frequency coordinates in the transformed space, and $\mathcal{F}^{-1}$ denotes the corresponding inverse transform (IFFT).

In our implementation, before applying FFT, the feature $\hat{F} \in \mathbb{R}^{H \times W \times C}$ is first projected by a linear layer to amplify high-frequency signals, acting as a high-pass filter [27, 39, 56]. We then use FFT to decompose the feature into real and imaginary parts, $(R, I)$. Mathematically, we have:

$$R, I = \text{FFT}\big(\text{Linear}(\hat{F})\big), \qquad (12)$$

which produces a complex spectrum containing crucial global information. Next, these real and imaginary components are concatenated along the channel dimension:

$$F = \text{IFFT}\big(\text{Linear}([R, I])\big), \qquad (13)$$

where $[\cdot]$ denotes channel concatenation, and $F$ is the frequency-domain interaction-fused feature. This second linear layer further modulates and refines the frequency components, while the inverse FFT transforms them back to the spatial domain for subsequent restoration stages.

## 3.5 Multi-Scale Flow Gating Network

Traditional feed-forward networks often rely on deep convolutions to enhance the locality of latent features, but they tend to overlook the effectiveness of multi-scale feature representation in removing rain streaks [1, 6]. To address this, we combine a gating mechanism with multi-scale representation learning by introducing convolutions of different scales into both the gating and value branches (see Fig. 2). This controls the flow of expert information across the

various levels of our pipeline, facilitating the flow of cross-level, multi-scale information and effectively extracting multi-scale local details. Given the input feature $F$, the computation of MSGN can be formulated as:

$$\hat{F} = PW(F), [\hat{F}_1, \hat{F}_2] = \hat{F}, \qquad (14)$$

$$\hat{F}' = Linear(DW^{3\times3}(\hat{F}_1) \odot DW^{5\times5}(\hat{F}_2)). \qquad (15)$$

## 4 Experiments

### 4.1 Experimental Settings

**Implementation Details.** In CPRAformer, $\{N_1, N_2, N_3, N_4\}$ are set to {4, 6, 6, 8}, and the attention heads of the four levels of CPRAformerB are set to {1, 2, 4, 8}. The initial channel $C$ is set to 48. During training, we use the Adam optimizer with a patch size of 64×64, a batch size of 12, and the number of epochs set to 300. For detailed experimental settings, please refer to the supplements.

**Data and Evaluation.** We follow the majority of previous research practices to train and validate our model [3, 12, 23, 71, 72]. Specifically, we use 13,712 image pairs collected from multiple datasets [19, 33, 66, 76, 77] for training and evaluate the model on five synthetic datasets (Test100 [77], Rain100H [66], Rain100L [66], Test2800 [18], Test1200 [75]) and one real-world dataset [59]. The aforementioned datasets primarily target rain streak removal. However, raindrops represent another form of contamination in rain removal tasks. Therefore, we train and evaluate our model on the Raindrop-A and Raindrop-B datasets [41]. Consistent with existing methods [20, 25], we adopt PSNR and SSIM as evaluation metrics for the aforementioned benchmarks, both calculated on the Y channel (luminance) in the YCbCr color space.
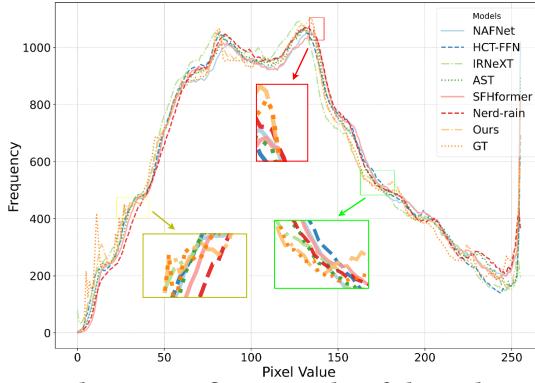
Figure 7: The average fitting results of the Y channel histogram curve in the YCbCr space on the synthetic dataset, our method produces results most similar to the ground truth.

Table 3: NIQE results under the real-world scenario [59].

| Methods | Input | DRSformer [6] | SFHformer [27] | Nerd-rain [7] | FSNet [12] | CPRAformer |
|---------|-------|---------------|----------------|---------------|------------|------------|
| NIQE ↓ | 5.923 | 5.814 | 5.745 | 5.711 | 5.667 | **5.556** |

Table 4: Ablation study of dual aggregation and dual paradigm strategy.

| SPC-SA | SPR-SA | Test100 | | Rain100H | | Average | |
|--------|--------|---------|---------|----------|---------|---------|---------|
| | | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| ✔ | | 27.62 | 0.868 | 28.90 | 0.866 | 30.66 | 0.903 |
| | ✔ | 27.40 | 0.866 | 29.01 | 0.864 | 30.91 | 0.900 |
| ✔ | ✔ | **28.80** | **0.891** | **29.22** | **0.871** | **31.56** | **0.913** |

## 4.2 Comparison with State-of-the-Art Methods

We compared CPRAformer with 22 image deraining methods: RES-CAN [31], PreNet [45], SPDNet [68], PCNet [26], MPRNet [72], HINet [3], ALformer [23], DANet [24], Uformer [61], NAFNet [2], MFDNet[58], HCT-FFN [8], DRSformer [6], FSNet [12], ChaIR [11], OKNet [13], AST [79], SFHformer [27], IRNeXT [14], Nerd-rain [7], MSDT [1] and AdaIR [15]. To ensure a fair comparison, we retrained all the above methods from scratch in our environment using their official source codes without any pretraining or fine-tuning.

**Rain Streak Synthetic Datasets.** Table 1 presents the quantitative evaluation results on five benchmark datasets, where it is evident that our CPRAformer consistently and significantly outperforms existing methods. Specifically, compared to the most recent top-performing method, AdaIR [15], CPRAformer improves the average performance across all datasets by 0.52dB. On certain datasets (such as Test100), the gain reaches up to 1.01dB. Fig. 4 shows a comparison of the visual quality of samples generated by recent methods. Thanks to the interaction and fusion of dual paradigms, CPRAformer effectively removes rain streaks while preserving details and realistic textures in the background image. Additionally, we provide a comparison of the "Y" channel histogram fitting curves in Fig. 7, confirming the consistency of the deraining results with the ground truth in statistical distribution.

**Raindrops Synthetic Datasets.** We conducted further experiments on the raindrop datasets [41], and the quantitative results are shown in Table 2, where our model achieved the highest performance. Visual comparisons in Fig 5 indicate that our method effectively removes raindrops while preserving fine texture details. This demonstrates that our model attains optimal performance
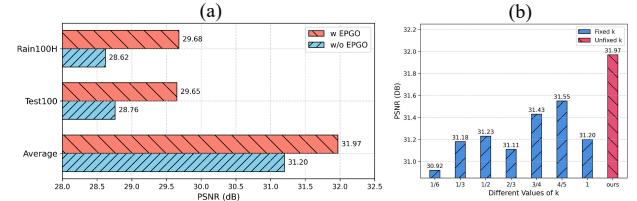


Figure 8: Ablation study of EPGO.

Table 5: Ablation study of AAFM.

| Method | Test100 | | Rain100H | | Average | |
|--------|---------|---------|----------|---------|---------|---------|
| | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| baseline | 28.80 | 0.891 | 29.22 | 0.871 | 31.56 | 0.913 |
| w/ stage1 | 28.78 | 0.892 | 29.68 | 0.873 | 31.68 | 0.913 |
| w/ stage2 | **29.65** | **0.895** | **29.68** | **0.875** | **31.97** | **0.916** |

Table 6: Ablation study of MSGN.

| Method | Test100 | | Rain100H | | Average | |
|--------|---------|---------|----------|---------|---------|---------|
| | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| SFFN [17] | 27.20 | 0.864 | 28.51 | 0.862 | 31.01 | 0.905 |
| DFN [34] | 29.45 | 0.893 | 29.66 | 0.872 | 31.58 | 0.913 |
| ConvGLU [46] | 28.33 | 0.891 | 29.26 | 0.870 | 31.46 | 0.911 |
| MSGN | **29.65** | **0.895** | **29.68** | **0.875** | **31.97** | **0.916** |

under both types of rain contamination, further confirming the generalization capability of our CPRAformer.

**Real-world Datasets.** To further demonstrate the generalization and robustness of CPRAformer, we conducted comparisons on real-world datasets [59]. As shown in Fig. 6, all other methods produced suboptimal results in either rain removal or detail restoration.

In contrast, our proposed CPRAformer outperformed the other methods, achieving visually pleasing restoration in challenging examples. This indicates that CPRAformer can generalize well to unseen real-world data types.

**Perceptual quality assessment.** We followed the method in [6, 79] to evaluate the perceptual quality of our proposed CPRAformer. The results, shown in Table 3, demonstrate that CPRAformer achieves a lower NIQE compared to other methods, indicating it delivers better perceptual quality in real rain scenes.

## 4.3 Ablation Studies

In this section, we conduct ablation experiments on five synthetic datasets to investigate the effect of each component. To ensure fair comparison, all ablation studies are performed under the same environment and training details. Due to space limits, we present results for two datasets along with the average across all five datasets. Further ablation experiments are provided in the supplements.

**Dual Aggregation and Dual Paradigm Strategy.** To investigate the effect of using both SPC-SA and SPR-SA simultaneously, we conducted multiple experiments, with results shown in Table 4. The first and second rows of the table indicate that we replaced all attention modules in CPRAformer with either SPC-SA or SPR-SA. The third row represents the scenario where both attention mechanisms are used in CPRAformer. Additionally, none of the models employed AAFM in this comparison. We observe that the best performance of 31.56 dB is achieved when both types of self-attention are utilized. This indicates that the different representations of the two paradigms are crucial for high-quality image deraining.
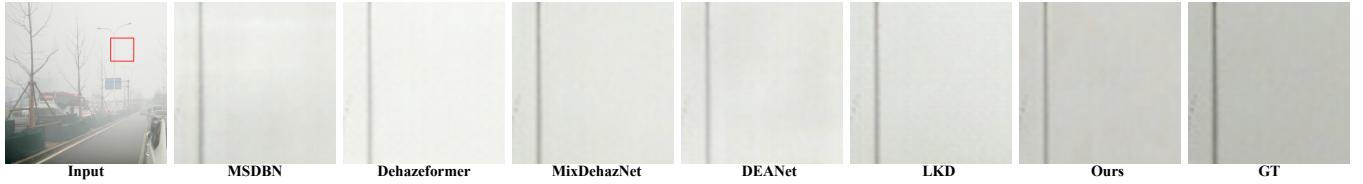
**Figure 9: The qualitative comparison on hazy images. Our result has the best visual quality and details.**

**Table 7: The averaged precision in recognizing the rain and water label from Google Vision API. Note that lower scores indicate better deraining performance [21, 50].**

| Model | Rainy Input | HCT-FFN [8] | DRSformer [6] | ChaIR [11] | IRNeXT [14] | OKNet [13] | Nerd-rain [7] | MSDT [1] | FSNet [12] | Ours |
|---|---|---|---|---|---|---|---|---|---|---|
| Percentage Score ↓ | 0.682 | 0.673 | 0.652 | 0.648 | 0.672 | 0.678 | 0.660 | 0.644 | 0.653 | **0.632** |

**Table 8: Quantitative comparisons of state-of-the-art methods on the REIDE-6K [29, 47] and Haze4K [36] datasets. Bold and underlined indicate the best and second-best results.**

| Method | Year | RESIDE-6K [29] | | Haze4K [36] | |
|---|---|---|---|---|---|
| | | PSNR ↑ | SSIM ↑ | PSNR ↑ | SSIM ↑ |
| MSBDN [16] | CVPR2020 | 25.07 | 0.896 | 25.74 | 0.918 |
| FFA-Net [42] | AAAI2020 | 24.44 | 0.925 | 27.38 | 0.942 |
| UHD [62] | CVPR2021 | 25.68 | 0.913 | 25.40 | 0.918 |
| Uformer [61] | CVPR2022 | 26.29 | 0.925 | 26.43 | 0.937 |
| gUNet [48] | Arxiv2022 | 25.67 | 0.924 | 26.50 | 0.940 |
| LKD [38] | ICME2023 | 25.42 | 0.925 | 27.39 | 0.938 |
| Dehazeformer [47] | TIP2023 | 26.25 | 0.931 | _27.45_ | _0.946_ |
| MB-TaylorFormer [43] | ICCV2023 | 26.28 | 0.923 | 26.34 | 0.933 |
| MixDehazeNet [37] | IJCNN2024 | 26.62 | 0.939 | 27.34 | 0.945 |
| DEANet [9] | TIP2024 | 26.61 | 0.932 | 26.94 | 0.942 |
| SFHformer [27] | ECCV2024 | _27.08_ | _0.940_ | 26.92 | 0.941 |
| **CPRAformer (Ours)** | – | **27.70** | **0.944** | **27.97** | **0.952** |

**Efficient Prompt Guide Operator.** To verify the effectiveness of EPGO, we first removed the Top-k mechanism and EPGO itself, with experimental results shown in Fig. 8 (a). It is observed that EPGO consistently achieves high-fidelity recovery with excellent PSNR performance. The key aspect of EPGO is to generate dynamic k values through prompt information, guiding the attention matrix to filter out important information. To this end, we tested on five datasets, comparing the dynamic k values generated by EPGO with fixed k values, with average results shown in Fig. 8 (b). The dynamic k values adaptively handle the complex and variable rain streaks in real scenarios, enhancing the model's robustness.

**Adaptive Alignment Frequency Module.** We validated the effectiveness of the AAFM through comprehensive ablation experiments, with results shown in Table 5. Specifically, we used the model from the third row of Table 4 as the baseline. First, we introduced the first alignment stage of AAFM, which achieved a gain of 0.12 dB. Subsequently, we added the second fusion stage, resulting in a gain of 0.41 dB compared to the baseline model. This indicates that AAFM maximizes the feature advantages of both paradigms, facilitating high-frequency interactions and deep fusion of information.

**Multi-Scale Flow Gating Network.** To validate the effectiveness of MSGN, we replaced it with Standard Feed-Forward Network (SFFN) [17], Depth-wise Convolution Equipped Feed-Forward Network (DFN) [34], and Convolutional Gated Linear Unit (ConvGLU) [46]. The results in Table 6 indicate that MSGN achieves optimal performance by effectively representing latent multi-scale perceptual information and enhancing information flow between classes.



**Figure 10: Semantic segmentation results on Deeplab V3 [4].**

## 4.4 Other Related Tasks

We selected image dehazing to validate the extensibility and robustness of CPRAformer in image restoration tasks. Following most previous studies [9, 37, 47, 48], we trained and validated our model. Specifically, we used two popular dehazing datasets, RESIDE-6K [29, 47] and Haze-4K [36], and compared CPRAformer with 11 state-of-the-art dehazing methods. Quantitative results, as shown in Table 8, indicate that CPRAformer achieved the best performance on both datasets. For example, on RESIDE-6K, CPRAformer improved PSNR by 0.62 dB over the previous SOTA method SFHformer, which is a significant enhancement. Furthermore, visual comparisons in Fig. 9 reveal that other methods produce images with unnatural shadows and high-frequency regions, as well as residual haze. In contrast, CPRAformer restores clear images, preserves texture and color details, and minimizes haze remnants. Additionally, we evaluated CPRAformer on downstream tasks. In object detection, 10 randomly selected images were processed with the Google Vision API [6, 21, 50, 69], and Table 7 shows that our method achieved the best average recognition. For semantic segmentation, a pre-trained DeepLab v3 [4] was used, and as shown in Fig. 10, CPRAformer's output is closest to the ground truth. See supplements for details.

## 5 Conclusion

In this paper, we propose a novel image deraining Transformer model, CPRAformer, based on dual-paradigm representation learning. Specifically, we introduce SPC-SA, which adaptively adjusts the sparsity of the neural network, enhancing global modeling capability while facilitating the flow of expert information across channels. Additionally, SPR-SA emphasizes the spatial distribution of rain variations, focusing on local feature extraction. Furthermore, we propose AAFM and MSGN to fully integrate features from both paradigms, promoting interaction between different types of representations and achieving cross-scale feature interaction. Extensive experiments on 10 benchmark datasets demonstrate that CPRAformer exhibits strong generalization and robustness.

# References

[1] Hongming Chen, Xiang Chen, Jiyang Lu, and Yufeng Li. 2024. Rethinking multi-scale representations in deep deraining transformer. In *AAAI*, Vol. 38. 1046–1053.

[2] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. 2022. Simple baselines for image restoration. In *ECCV*. Springer, 17–33.

[3] Liangyu Chen, Xin Lu, Jie Zhang, Xiaojie Chu, and Chengpeng Chen. 2021. Hinet: Half instance normalization network for image restoration. In *CVPR*. 182–192.

[4] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. 2018. Encoder-Decoder with Atrous Separable Convolution for Semantic Image Segmentation. In *ECCV*.

[5] Sixiang Chen, Tian Ye, Yun Liu, and Erkang Chen. 2024. Dual-former: Hybrid self-attention transformer for efficient image restoration. *Digital Signal Processing* 149 (2024), 104485.

[6] Xiang Chen, Hao Li, Mingqiang Li, and Jinshan Pan. 2023. Learning a sparse transformer network for effective image deraining. In *CVPR*. 5896–5905.

[7] Xiang Chen, Jinshan Pan, and Jiangxin Dong. 2024. Bidirectional multi-scale implicit neural representations for image deraining. In *CVPR*. 25627–25636.

[8] Xiang Chen, Jinshan Pan, Jiyang Lu, Zhentao Fan, and Hao Li. 2023. Hybrid cnn-transformer feature fusion for single image deraining. In *AAAI*, Vol. 37. 378–386.

[9] Zixuan Chen, Zewei He, and Zhe-Ming Lu. 2024. DEA-Net: Single image dehazing based on detail-enhanced convolution and content-guided attention. *TIP* 33 (2024), 1002–1015.

[10] Zheng Chen, Yulun Zhang, Jinjin Gu, Linghe Kong, Xiaokang Yang, and Fisher Yu. 2023. Dual Aggregation Transformer for Image Super-Resolution. In *ICCV*.

[11] Yuning Cui and Alois Knoll. 2023. Exploring the potential of channel interactions for image restoration. *Knowledge-Based Systems* 282 (2023), 111156.

[12] Yuning Cui, Wenqi Ren, Xiaochun Cao, and Alois Knoll. 2023. Image restoration via frequency selection. *TPAMI* 46, 2 (2023), 1093–1108.

[13] Yuning Cui, Wenqi Ren, and Alois Knoll. 2024. Omni-kernel network for image restoration. In *AAAI*, Vol. 38. 1426–1434.

[14] Yuning Cui, Wenqi Ren, Sining Yang, Xiaochun Cao, and Alois Knoll. 2023. Irnext: Rethinking convolutional network design for image restoration. In *ICML*.

[15] Yuning Cui, Syed Waqas Zamir, Salman Khan, Alois Knoll, Mubarak Shah, and Fahad Shahbaz Khan. 2025. AdaIR: Adaptive All-in-One Image Restoration via Frequency Mining and Modulation. In *ICLR*.

[16] Hang Dong, Jinshan Pan, Lei Xiang, Zhe Hu, Xinyi Zhang, Fei Wang, and Ming-Hsuan Yang. 2020. Multi-scale boosted dehazing network with dense feature fusion. In *CVPR*. 2157–2167.

[17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. 2020. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020).

[18] Xueyang Fu, Jiabin Huang, Delu Zeng, Yue Huang, Xinghao Ding, and John Paisley. 2017. Removing rain from single images via a deep detail network. In *CVPR*. 3855–3863.

[19] Xueyang Fu, Jiabin Huang, Delu Zeng, Yue Huang, Xinghao Ding, and John Paisley. 2017. Removing rain from single images via a deep detail network. In *CVPR*. 3855–3863.

[20] Xueyang Fu, Jie Xiao, Yurui Zhu, Aiping Liu, Feng Wu, and Zheng-Jun Zha. 2023. Continual image deraining with hypergraph convolutional networks. *TPAMI* 45, 8 (2023), 9534–9551.

[21] Ning Gao, Xingyu Jiang, Xiuhui Zhang, and Yue Deng. 2024. Efficient Frequency-Domain Image Deraining with Contrastive Regularization. In *ECCV*. Springer, 240–257.

[22] Shuhang Gu, Deyu Meng, Wangmeng Zuo, and Lei Zhang. 2017. Joint convolutional analysis and synthesis sparse representation for single image layer separation. In *ICCV*. 1708–1716.

[23] Kui Jiang, Zhongyuan Wang, Chen Chen, Zheng Wang, Laizhong Cui, and Chia-Wen Lin. 2022. Magic ELF: Image deraining meets association learning and transformer. *ACMMM* (2022).

[24] Kui Jiang, Zhongyuan Wang, Zheng Wang, Peng Yi, Junjun Jiang, Jinsheng Xiao, and Chia-Wen Lin. 2022. Danet: Image deraining via dynamic association learning.. In *IJCAI*. 980–986.

[25] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Baojin Huang, Yimin Luo, Jiayi Ma, and Junjun Jiang. 2020. Multi-scale progressive fusion network for single image deraining. In *CVPR*. 8346–8355.

[26] Kui Jiang, Zhongyuan Wang, Peng Yi, Chen Chen, Zheng Wang, Xiao Wang, Junjun Jiang, and Chia-Wen Lin. 2021. Rain-free and residue hand-in-hand: A progressive coupled network for real-time image deraining. *IEEE Transactions on Image Processing* 30 (2021), 7404–7418.

[27] Xingyu Jiang, Xiuhui Zhang, Ning Gao, and Yue Deng. 2024. When Fast Fourier Transform Meets Transformer for Image Restoration. In *ECCV*. Springer, 381–402.

[28] Li-Wei Kang, Chia-Wen Lin, and Yu-Hsiang Fu. 2011. Automatic single-image-based rain streaks removal via image decomposition. *TIP* 21, 4 (2011), 1742–1755.

[29] Boyi Li, Wenqi Ren, Dengpan Fu, Dacheng Tao, Dan Feng, Wenjun Zeng, and Zhangyang Wang. 2018. Benchmarking single-image dehazing and beyond. *TIP* 28, 1 (2018), 492–505.

[30] Pengpeng Li, Jiyu Jin, Guiyue Jin, Lei Fan, Xiao Gao, Tianyu Song, and Xiang Chen. 2022. Deep scale-space mining network for single image deraining. In *CVPR*. 4276–4285.

[31] Xia Li, Jianlong Wu, Zhouchen Lin, Hong Liu, and Hongbin Zha. 2018. Recurrent squeeze-and-excitation context aggregation net for single image deraining. In *ECCV*. 254–269.

[32] Yu Li, Robby T Tan, Xiaojie Guo, Jiangbo Lu, and Michael S Brown. 2016. Rain streak removal using layer priors. In *CVPR*. 2736–2744.

[33] Yu Li, Robby T Tan, Xiaojie Guo, Jiangbo Lu, and Michael S Brown. 2016. Rain streak removal using layer priors. In *ICCV*. 2736–2744.

[34] Yawei Li, Kai Zhang, Jiezhang Cao, Radu Timofte, and Luc Van Gool. 2021. Localvit: Bringing locality to vision transformers. *arXiv preprint arXiv:2104.05707* (2021).

[35] Yuanchu Liang, Saeed Anwar, and Yang Liu. 2022. Drt: A lightweight single image deraining recursive transformer. In *CVPR*. 589–598.

[36] Ye Liu, Lei Zhu, Shunda Pei, Huazhu Fu, Jing Qin, Qing Zhang, Liang Wan, and Wei Feng. 2021. From synthetic to real: Image dehazing collaborating with unlabeled real data. In *ACMMM*. 50–58.

[37] LiPing Lu, Qian Xiong, Bingrong Xu, and Duanfeng Chu. 2024. Mixdehazenet: Mix structure block for image dehazing network. In *IJCNN*. IEEE, 1–10.

[38] Pinjun Luo, Guoqiang Xiao, Xinbo Gao, and Song Wu. 2023. LKD-Net: Large kernel convolution network for single image dehazing. In *ICME*. IEEE, 1601–1606.

[39] Namuk Park and Songkuk Kim. 2022. How do vision transformers work? *arXiv preprint arXiv:2202.06709* (2022).

[40] Yan-Tsung Peng and Wei-Hua Li. 2023. Rain2Avoid: Self-Supervised Single Image Deraining. In *ICASSP*. IEEE, 1–5.

[41] Rui Qian, Robby T. Tan, Wenhan Yang, Jiajun Su, and Jiaying Liu. 2018. Attentive Generative Adversarial Network for Raindrop Removal From a Single Image. In *CVPR*.

[42] Xu Qin, Zhilin Wang, Yuanchao Bai, Xiaodong Xie, and Huizhu Jia. 2020. FFA-Net: Feature fusion attention network for single image dehazing. In *AAAI*, Vol. 34. 11908–11915.

[43] Yuwei Qiu, Kaihao Zhang, Chenxi Wang, Wenhan Luo, Hongdong Li, and Zhi Jin. 2023. Mb-taylorformer: Multi-branch efficient transformer expanded by taylor formula for image dehazing. In *ICCV*. 12802–12813.

[44] Chun Ren, Danfeng Yan, Yuanqiang Cai, and Yangchun Li. 2023. Semi-swinderain: Semi-supervised image deraining network using swin transformer. In *ICASSP*. IEEE, 1–5.

[45] Dongwei Ren, Wangmeng Zuo, Qinghua Hu, Pengfei Zhu, and Deyu Meng. 2019. Progressive image deraining networks: A better and simpler baseline. In *CVPR*. 3937–3946.

[46] Dai Shi. 2024. Transnext: Robust foveal visual perception for vision transformers. In *CVPR*. 17773–17783.

[47] Yuda Song, Zhuqing He, Hui Qian, and Xin Du. 2023. Vision Transformers for Single Image Dehazing. *TIP* 32 (2023), 1927–1941.

[48] Yuda Song, Yang Zhou, Hui Qian, and Xin Du. 2022. Rethinking Performance Gains in Image Dehazing Networks. *arXiv preprint arXiv:2209.11448* (2022).

[49] Jian-Nan Su, Min Gan, Guang-Yong Chen, Wenzhong Guo, and CL Philip Chen. 2024. High-similarity-pass attention for single image super-resolution. *TIP* 33 (2024), 610–624.

[50] Shangquan Sun, Wenqi Ren, Xinwei Gao, Rui Wang, and Xiaochun Cao. 2024. Restoring Images in Adverse Weather Conditions via Histogram Transformer. *ECCV* (2024).

[51] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. 2022. Maxvit: Multi-axis vision transformer. In *ECCV*. Springer, 459–479.

[52] Jeya Maria Jose Valanarasu, Rajeev Yasarla, and Vishal M Patel. 2022. Transweather: Transformer-based restoration of images degraded by adverse weather conditions. In *CVPR*. 2353–2363.

[53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *NeurIPS* 30 (2017).

[54] Cong Wang, Jinshan Pan, Wei Wang, Jiangxin Dong, Mengzhu Wang, Yakun Ju, and Junyang Chen. 2023. Promptrestorer: A prompting image restoration method with degradation perception. *NIPS* 36 (2023), 8898–8912.

[55] Cong Wang, Xiaoying Xing, Yutong Wu, Zhixun Su, and Junyang Chen. 2020. Dcsfn: Deep cross-scale fusion network for single image rain removal. In *ACMMM*. 1643–1651.

[56] Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang. 2022. Anti-Oversmoothing in Deep Vision Transformers via the Fourier Domain Analysis: From Theory to Practice. In *International Conference on Learning Representations*.

[57] Qiong Wang, Kui Jiang, Jinyi Lai, Zheng Wang, and Jianhui Zhang. 2023. Hpcnet: A hybrid progressive coupled network for image deraining. In *ICME*. IEEE, 2747–2752.

[58] Qiong Wang, Kui Jiang, Zheng Wang, Wenqi Ren, Jianhui Zhang, and Chia-Wen Lin. 2023. Multi-scale fusion and decomposition network for single image deraining. *TIP* 33 (2023), 191–204.

[59] Tianyu Wang, Xin Yang, Ke Xu, Shaozhe Chen, Qiang Zhang, and Rynson WH Lau. 2019. Spatial attentive single-image deraining with a high quality real rain dataset. In *CVPR*. 12270–12279.

[60] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. 2021. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *ICCV*. 568–578.

[61] Zhendong Wang, Xiaodong Cun, Jianmin Bao, Wengang Zhou, Jianzhuang Liu, and Houqiang Li. 2022. Uformer: A general u-shaped transformer for image restoration. In *CVPR*. 17683–17693.

[62] Boxue Xiao, Zhuoran Zheng, Xiang Chen, Chen Lv, Yunliang Zhuang, and Tao Wang. 2022. Single UHD Image Dehazing via Interpretable Pyramid Network. arXiv:2202.08589

[63] Jie Xiao, Xueyang Fu, Aiping Liu, Feng Wu, and Zheng-Jun Zha. 2022. Image de-raining transformer. *TPAMI* 45, 11 (2022), 12978–12995.

[64] Jie Xiao, Xueyang Fu, Aiping Liu, Feng Wu, and Zheng-Jun Zha. 2023. Image De-Raining Transformer. *TPAMI* 45, 11 (2023), 12978–12995. doi:10.1109/TPAMI.2022.3183612

[65] Yi Xiao, Qiangqiang Yuan, Kui Jiang, Jiang He, Chia-Wen Lin, and Liangpei Zhang. 2024. TTST: A top-k token selective transformer for remote sensing image super-resolution. *TIP* 33 (2024), 738–752.

[66] Wenhan Yang, Robby T Tan, Jiashi Feng, Jiaying Liu, Zongming Guo, and Shuicheng Yan. 2017. Deep joint rain detection and removal from a single image. In *CVPR*. 1357–1366.

[67] Wenhan Yang, Robby T Tan, Shiqi Wang, Yuming Fang, and Jiaying Liu. 2020. Single image deraining: From model-based to data-driven and beyond. *TPAMI* 43, 11 (2020), 4059–4077.

[68] Qiaosi Yi, Juncheng Li, Qinyan Dai, Faming Fang, Guixu Zhang, and Tieyong Zeng. 2021. Structure-preserving deraining with residue channel prior guidance. In *ICCV*. 4238–4247.

[69] Qiaosi Yi, Juncheng Li, Qinyan Dai, Faming Fang, Guixu Zhang, and Tieyong Zeng. 2021. Structure-Preserving Deraining with Residue Channel Prior Guidance. In *ICCV*. 4218–4227.

[70] Yi Yu, Wenhan Yang, Yap-Peng Tan, and Alex C Kot. 2022. Towards robust rain removal against adversarial attacks: A comprehensive benchmark analysis and beyond. In *CVPR*. 6013–6022.

[71] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. 2022. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*. 5728–5739.

[72] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. 2021. Multi-stage progressive image restoration. In *CVPR*. 14821–14831.

[73] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. 2022. Learning enriched features for fast image restoration and enhancement. *TPAMI* 45, 2 (2022), 1934–1948.

[74] He Zhang and Vishal M Patel. 2017. Convolutional sparse and low-rank coding-based rain streak removal. In *WACV*. IEEE, 1259–1267.

[75] He Zhang and Vishal M Patel. 2018. Density-aware single image de-raining using a multi-stream dense network. In *CVPR*. 695–704.

[76] He Zhang and Vishal M Patel. 2018. Density-aware Single Image De-raining using a Multi-stream Dense Network. In *CVPR*.

[77] He Zhang, Vishwanath Sindagi, and Vishal M Patel. 2019. Image de-raining using a conditional generative adversarial network. *TCSVT* (2019).

[78] Mingjun Zheng, Long Sun, Jiangxin Dong, and Jinshan Pan. 2024. SMFANet: A lightweight self-modulation feature aggregation network for efficient image super-resolution. In *European Conference on Computer Vision*. Springer, 359–375.

[79] Shihao Zhou, Duosheng Chen, Jinshan Pan, Jinglei Shi, and Jufeng Yang. 2024. Adapt or perish: Adaptive sparse transformer with attentive feature refinement for image restoration. In *CVPR*. 2952–2963.