# Cross-receptive Focused Inference Network for Lightweight Image Super-Resolution

Wenjie Li, Juncheng Li, Guangwei Gao, *Senior Member, IEEE,* Jiantao Zhou, *Senior Member, IEEE,* Jian Yang, *Member, IEEE* and Guo-Jun Qi, *Fellow, IEEE*

*Abstract*—With the development of deep learning, single image super-resolution (SISR) has achieved significant breakthroughs. Recently, methods to enhance the performance of SISR networks based on global feature interactions have been proposed. However, the capabilities of neurons that need to adjust their function in response to the context dynamically are neglected. To address this issue, we propose a lightweight Cross-receptive Focused Inference Network (CFIN), a hybrid network composed of a Convolutional Neural Network (CNN) and a Transformer. Specifically, a novel Cross-receptive Field Guide Transformer (CFGT) is designed to adaptively modify the network weights by using modulated convolution kernels combined with local representative semantic information. In addition, a CNN-based Cross-scale Information Aggregation Module (CIAM) is proposed to make the model better focused on potentially practical information and improve the efficiency of the Transformer stage. Extensive experiments show that our proposed CFIN is a lightweight and efficient SISR model, which can achieve a good balance between computational cost and model performance.

*Index Terms*—Cross-receptive, cross-attention, hybrid network, lightweight mode, SISR.

## I. INTRODUCTION

The task of Single Image Super-Resolution (SISR) aims to estimate a realistic High-Resolution (HR) image from the Low-Resolution (LR) one, which plays a fundamental role in various computer vision tasks, including surveillance imaging, autonomous driving, and medical imaging [1]–[3]. As an ill-posed problem, SISR is still a challenging task. To solve this task, many Convolutional Neural Networks (CNN) based methods have been proposed to directly learn the mapping between the LR and HR image pairs. For example, Dong *et al.* presented the first CNN-based model, dubbed SRCNN [4]. Although SRCNN only has three convolutional layers, its performance is significantly better than traditional solutions.

Wenjie Li and Guangwei Gao are with the Institute of Advanced Technology, Nanjing University of Posts and Telecommunications, Nanjing, China, and also with the Key Laboratory of Intelligent Information Processing, Nanjing Xiaozhuang University, Nanjing, China (e-mail: csggao@gmail.com,lewj2408@gmail.com).

Juncheng Li is with the Center for Mathematical Artificial Intelligence, Department of Mathematics, The Chinese University of Hong Kong, Hong Kong, China (e-mail: cvjunchengli@gmail.com).

Jiantao Zhou is with the State Key Laboratory of Internet of Things for Smart City, University of Macau, Taipa, Macau, and also with the Department of Computer and Information Science, Faculty of Science and Technology, University of Macau, Taipa, Macau (e-mail: jtzhou@um.edu.mo).

Jian Yang is with the School of Computer Science and Technology, Nanjing University of Science and Technology, Nanjing, China (e-mail: csjyang@njust.edu.cn).

Guo-Jun Qi is with the Department of Computer Science, University of Central Florida, Orlando, USA (e-mail: guojunq@gmail.com).
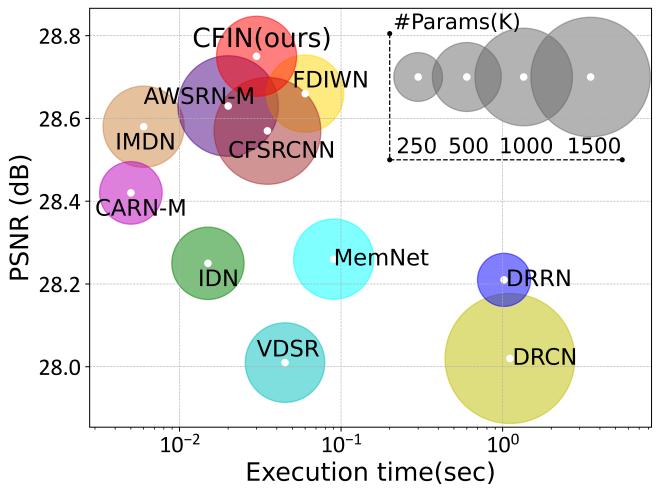


Fig. 1: Model inference time studies on Set14 ($\times 4$).

Subsequently, a series of networks with complex architectures were proposed, and deep CNN-based methods have achieved remarkable progress in SISR [5], [6]. Although these models have achieved promising results, their computational cost is often too huge (Fig. 1) to be popularized and widely used.

To solve the aforementioned problems, constructing lightweight SISR models has attracted more and more attention. Among these CNN-based lightweight models, most of them focus on efficient network architecture design, such as neural architecture search [7], multi-scale structure design [8], [9], and channel grouping strategy [10], [11]. However, convolution kernels can only extract local features, which is difficult to model the long-term dependencies of the image. As a complementary solution to CNN, Transformer has achieved excellent performance in many visual tasks with its powerful global modeling capability [12]–[15]. Recently, some Transformer-based SISR methods have been proposed [16]–[18]. For instance, SwinIR [19] introduced Transformer into SISR, relying on its advantage of using a shifted window scheme to model long-term dependencies, showing the great promise of Transformer in SISR. ESRT [20] is an efficient SISR model that combine a lightweight CNN and lightweight Transformer in an end-to-end model. However, most existing Transformer-based methods ignore the importance of dynamic modeling with context. As studied in neurology [21], the morphology of neurons should change adaptively with changes in the environment. This adjustment mechanism has been studied in many fields. For example, Jia *et al.* [22] gener-

ated convolution kernel weights from the features extracted from another network. Chen *et al.* [23] aggregated multiple convolution kernels in parallel and based on local attention to adaptively adjust the weights. Lin *et al.* [24] presented the context-gated convolution to incorporate context awareness into the convolutional layer.

Motivated by the above methods, in this paper, we introduce the power of contextual reasoning into Transformer and propose a lightweight Cross-receptive Focused Inference Network (CFIN) for SISR. CFIN is a hybrid network composed of a Convolutional Neural Network (CNN) and a Transformer. In the convolution stag, a Cross-scale Information Aggregation Module (CIAM) is designed to extract more potentially useful information with the help of the efficient Redundant Information Filter Unit (RIFU). In the Transformer stage, we propose an efficient Cross-receptive Field Guide Transformer (CFGT) to achieve cross-scale long-distance information fusion with the specially designed Context Guided Attention (CGA). In summary, the main contributions of this paper can be summarized as the following three aspects

- We propose a Redundant Information Filter Unit (RIFU), which can remove redundant information and accurately learn flexible local features. Meanwhile, an efficient Cross-scale Information Aggregation Module (CIAM) is specially designed for elaborately combining several RIFUs to ensure full use of local features.
- We propose a new Context Guided Attention (CGA) scheme, which can adaptively adjust the weights of the convolution kernel according to the representative local information in the context. In addition, a novel Cross-receptive Field Guide Transformer (CFGT) is proposed for cross-scale long-distance information fusion, which is more conducive to contextual reasoning within CGAs and suitable for long-distance modeling.
- We propose a lightweight and efficient Cross-receptive Focused Inference Network (CFIN) for SISR. CFIN elegantly integrates CNN and Transformer, achieving competitive performance and a good balance between computational cost and model performance.

The rest of this paper is organized as follows. Related works are reviewed in Section II. A detailed explanation of the proposed CFIN is given in Section III. The experimental results, ablation studies, and discussion are presented in Sections IV, V, and VI respectively. Finally, we draw a conclusion in Section VII.

## II. RELATED WORK

### A. Lightweight SISR Method

Due to the powerful learning ability of CNN, more and more CNN-based potent SISR methods have been proposed. However, most of the methods are limited in real-world applications by their huge computational cost. To handle this issue, some lightweight and efficient SISR methods have been presented [17], [25]–[31]. For example, IDN [32] used an information distillation network to selectively fuse features,

and then IMDN [33] improved it to build a lighter model. RFDN [34] combines channel splitting and residual structure to achieve better performance. IMRN [35] adopted the strategy of model pruning to alleviate the model size without significantly degrading the performance. FDIWN [36] improved the model performance by making full use of the intermediate layer features. However, most of these methods do not consider the influence of redundant features on model learning. Meanwhile, these CNN-based methods cannot observe images through global feature interaction. In this paper, we aim to explore a more efficient lightweight model that can focus on potentially practical information.

### B. Context Reasoning

With the in-depth understanding of deep learning, researchers tentatively explored how to increase the contextual information of the model. It can be roughly divided into the following categories. For example, using the attention mechanism to modify the feature representation, and the typical one is to modify the local features through the attention mechanism [37]. However, most of them can only modify the features by changing the input mapping. Recently, some works [23], [38]–[40] try to dynamically change network parameters by analyzing local or global information. However, they only consider local fragments [38], ignore weight tensors in convolutional layers [39] or have high training costs [40]. Researchers in [41] imitated the human visual system and simulates the bottom-up impact of semantic information on the model through reverse connections, but this feedback mechanism is difficult to explain in the model. In addition, none of them use local contextual information to dynamically guide global attention interactions. In this work, we aim to introduce context reasoning to further enhance the performance of the model and build a model that can adaptively modify the network weights.

## III. PROPOSED METHOD

### A. Cross-receptive Focused Inference Network

In this paper, we devise a lightweight Cross-receptive Focused Inference Network (CFIN) for SISR. As shown in Fig. 2, CFIN consists of several CNN-Transformer (CT) blocks and two PixelShuffle layers. Meanwhile, each CT block contains a convolution stage and a Transformer stage. And each two CT blocks are called once using a loop mechanism for a better trade-off between model size and performance. We define the input and output of CFIN as $I_{LR}$ and $I_{SR}$, respectively. Firstly, the dimension of the input image is rapidly increased to obtain shallow features $I_{shallow}$ for subsequent processing

$$I_{shallow} = F_{cr}(I_{LR}), \qquad (1)$$

where $F_{cr}(\cdot)$ is the channel expansion operation. Then the shallow features are sent to the CT block for feature extraction and the complete operation of each CT block can be defined as follows

$$I_{ct}^i = F_{CT}^i(I_{in}^i) = F_{cr}(F_T(F_{ce}(F_C(I_{in}^i)))) + I_{in}^i, \quad (2)$$
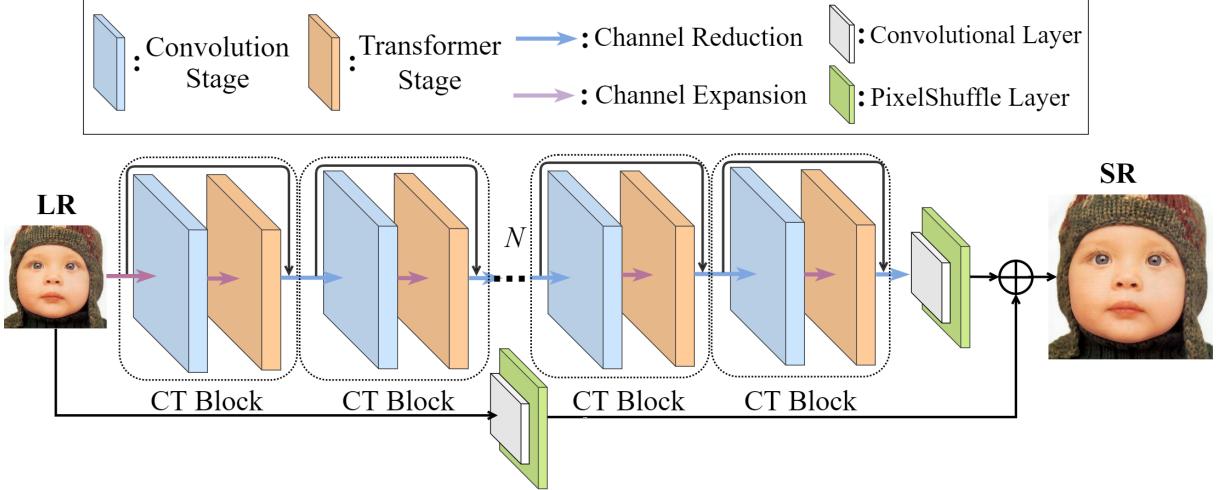
Fig. 2: The overall architecture of the proposed Cross-receptive Focused Inference Network (CFIN). It is worth noting that $1 \times 1$ convolution is used for channel reduction and channel expansion in CFIN.

where $F_{cr}(\cdot)$, $F_{ce}(\cdot)$, $F_C(\cdot)$, $F_F(\cdot)$ denote the channel reduction operation, channel expansion, the convolution stage, and the Transformer stage, respectively. $I_{in}^i$ and $I_{ct}^i$ are the input and output of the $i$-th CT block. After the operation of $N$ CT blocks, we can attain the final deep features as

$$I_{ct} = F_{CT}(I_{shallow}) = F_{CT}^N(F_{CT}^{N-1}(\cdots F_{CT}^2(F_{CT}^1(I_{shallow})))),$$ 
$$(3)$$

where $I_{ct}$ denotes the output of the $N$-th CT block. Finally, to obtain the final SR image, both $I_{ct}$ and $I_{LR}$ are simultaneously fed into the post-sampling reconstruction module

$$I_{SR} = F_{rec}(I_{ct}) + F_{rec}(I_{LR}), \qquad (4)$$

where $F_{rec}(\cdot)$ is the post-sampling reconstruction module, which composed of a $3 \times 3$ convolutional layer and a PixelShuffle layer.

### B. Convolution Stage

In the convolution stage, we propose a Cross-scale Information Aggregation Module (CIAM) to mine potential image information and make the model understand the preliminary SR information. As shown in Fig. 3, CIAM is mainly composed of three Redundant Information Filter Units (RIFU).

*1) Redundant Information Filter Unit (RIFU):* According to previous works [42], [43], we can know that it is easier to recover the smooth area that occupies most of the image area, but the complex texture information that occupies a small area of the image is difficult to recover. However, most SR methods tend to treat all areas of the image equally, which leads to the smooth area that accounts for most areas of the image being paid more attention by the network, which may miss the correct texture information, so accurate modeling cannot be achieved. In previous work [44], some methods convert the image from the time domain to the frequency domain by discrete cosine transform (DCT), and then manually set a threshold T to discard the frequency domain information greater or less than T to filter the secondary information. However, we also found that the input images used for training

vary widely. Using manual thresholding is sensitive to noise, which is not suitable for all images. Therefore, we aim to explore a method that can make the network find a mask that instructs model pays more attention to texture features. To achieve this, we propose a Redundant Information Filter Unit (RIFU). As given in Fig. 3, after the input feature $x$ enters RIFU, it first goes through a convolutional layer and an activation layer to obtain a mixed feature $X$. Then, we use a $1 \times 1$ convolutional layer to transform the number of channels of the output $R$ to $M$ ($M = 3$) and its process can be formulated as

$$X = f_{lrelu}(f_{conv}^{3x3}(x)), \qquad (5)$$

$$R = f_{n->3}(X), \qquad (6)$$

where $f_{conv}^{3x3}(\cdot)$ represents the $3 \times 3$ convolutional layer, $f_{lrelu} \cdot$ represents the LeakyRelu function, and $f_{n->3}$ represents the $1 \times 1$ convolutional layer. Next, in order to complete the end-to-end network optimization, we use the Gumbel Softmax trick [45] to generate a continuous differentiable normal distribution, which can well approximate the probability distribution represented by the network output and randomly add some sampling, so the Gumbel-Softmax enables RIFU to retain some potentially useful information in addition to salient features. It's process can be formulated as follows

$$GS_i = \frac{\exp((R_i + gs_i)/\tau)}{\sum_{m=1}^{M} \exp((R_{i_m} + gs_{i_m})/\tau)}, \qquad (7)$$

where $\tau$ defaults to 1, and $gs_i$ represents the noise obeying the $Gumbel(0,1)$ distribution. During training we need model to select only one channel from the three channels and the channel selection formula is as follows

$$y_{mask} = \arg one(GS_{i_m}), \qquad (8)$$

where $y_{mask}$ represents the single-channel output feature after masking, and $\arg one(.)$ represents the argmax branch of $GS_i$ in $m$ channels. After that, we multiply $X$ by the $y_{mask}$ feature
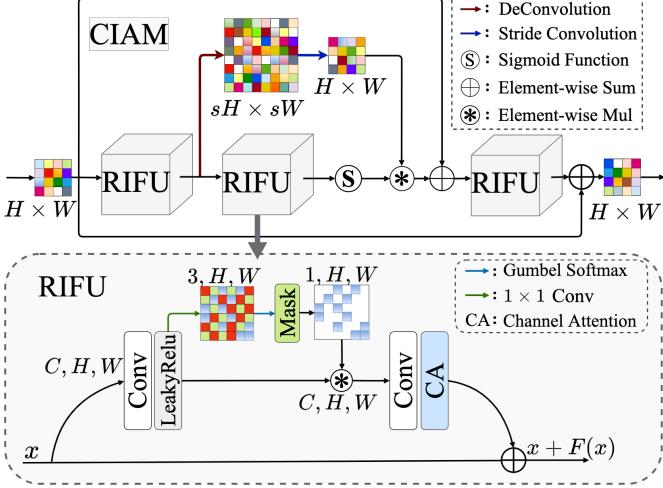
Fig. 3: The framework of the designed Cross-scale Information Aggregation Module (CIAM) and Redundant Information Filter Unit (RIFU).
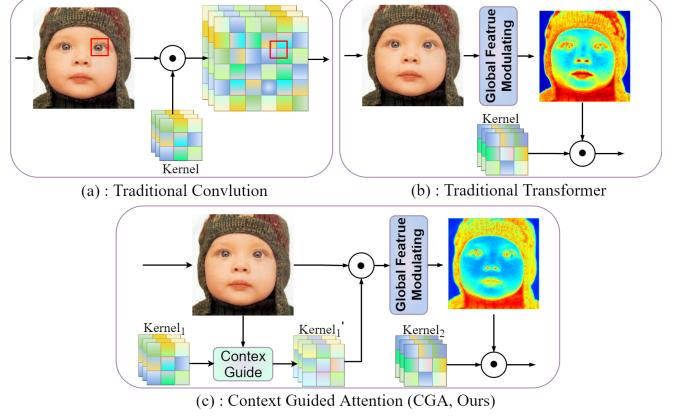


Fig. 4: (a) Traditional convolution tends to focus on local features. (b) Traditional Transformer uses global feature interactions to focus on key information. (c) Our proposed Context Guided Attention (CGA), guided by locally representative modulated convolution kernels, is able to adaptively combine global information to modify feature maps. $\odot$ denotes the convolution operation.

mask to preserve the initial details, and get the finally output $y$ of RIFU through residual connection

$$y = f_{CA}(f_{conv}^{3x3}(y_{mask} \times X)) + x, \qquad (9)$$

where $f_{CA}(\cdot)$ represents the channel attention mechanism.

*2) Cross-scale Information Aggregation Module (CIAM):*
Due to the difference in model structure, it is inevitable that the learned features are redundant. Meanwhile, with limited computing resources, we hope that the model will pay more attention to features with higher priority. So CIAM is designed to efficiently combine RIFU, its role is to extract more potentially effective information, mine deep image features, and make full use of information from different scales to observe image features.

In terms of network architecture design, simply combining RIFU in series is not conducive to gradient flow and collecting contextual information at different spatial locations. To solve this problem, we innovatively use two different scale spaces for feature transformation in the middle of the module. Among them, one is the original space, whose feature map size has the same resolution as the input, and the other is the large space after deconvolution operation. The receptive field of the transformed embedding is very large and can be used to guide the feature transformation in the original feature map. The process can be defined as

$$x_{HW}^1 = f_{RIFU}(x_{HW}), \qquad (10)$$

$$x_{HW}^2 = f_{sconv}(f_{deconv}(x_{HW})) \times f_{sig}(f_{RIFU}(x_{HW}^1)), \quad (11)$$

where $x_{HW}^i$ represents the output features with the size of $H \times W$ after the $i$-th LFFU, $f_{RIFU}(\cdot)$ represents the proposed LFFU, $f_{sconv}(\cdot)$ and $f_{deconv}(\cdot)$ represent the deconvolution and strided convolution with $s = 2$, and $f_{sig}(\cdot)$ represents the Sigmoid function.

Finally, to prevent the vanishing gradient, the dense learning mechanism is introduced into the module. Therefore,

this module has the potential for multiple RIFUs permutations and combinations, and its output $x_{HW}^3$ can be formulated as

$$x_{HW}^3 = f_{RIFU}(x_{HW}^2 + x_{HW}^1) + x_{HW}. \qquad (12)$$

*C. Transformer Stage*

In recent years, Transformer has shown great potential on SISR, which can learn global information of images through its powerful self-attention mechanism. However, existing self-attention mechanisms often neglect to incorporate context and build attention across features at different scales. To address this issue, in the Transformer stage, we propose a Cross-receptive Field Guided Transformer (CFGT). In CFGT, a specially designed Context Guided MaxConv (CGM) is introduced as its basic unit, which can provide cross-features for the model across scales. Meanwhile, it can adaptively adjust the weights of the network through semantic reasoning to meet the requirements of long-distance modeling.

*1) Context Guided Attention (CGA):* As shown in Fig. 4, traditional convolution tends to focus on local features and traditional Transformer uses self-attention to capture the global Information. However, the traditional Transformer cannot take into account the locally representative dominant features, and the self-attention mechanism also unable to incorporate contextual semantic information. To alleviate this problem, we propose a new attention, named Context Guided Attention (CGA), which is guided by the locally representative modulated convolution kernels and can adaptively combine global information to modify feature maps.

According to Fig. 5, we can clearly see that before computing the feature covariance to generate the global attention map, we introduce the Context Guided MaxConv (CGM) to emphasize the local context. Specifically, we first generate query $(Q)$, key $(K)$, and value $(V)$ projections from the input tensor $X \in \Re^{C \times H \times W}$. After that, we separately add
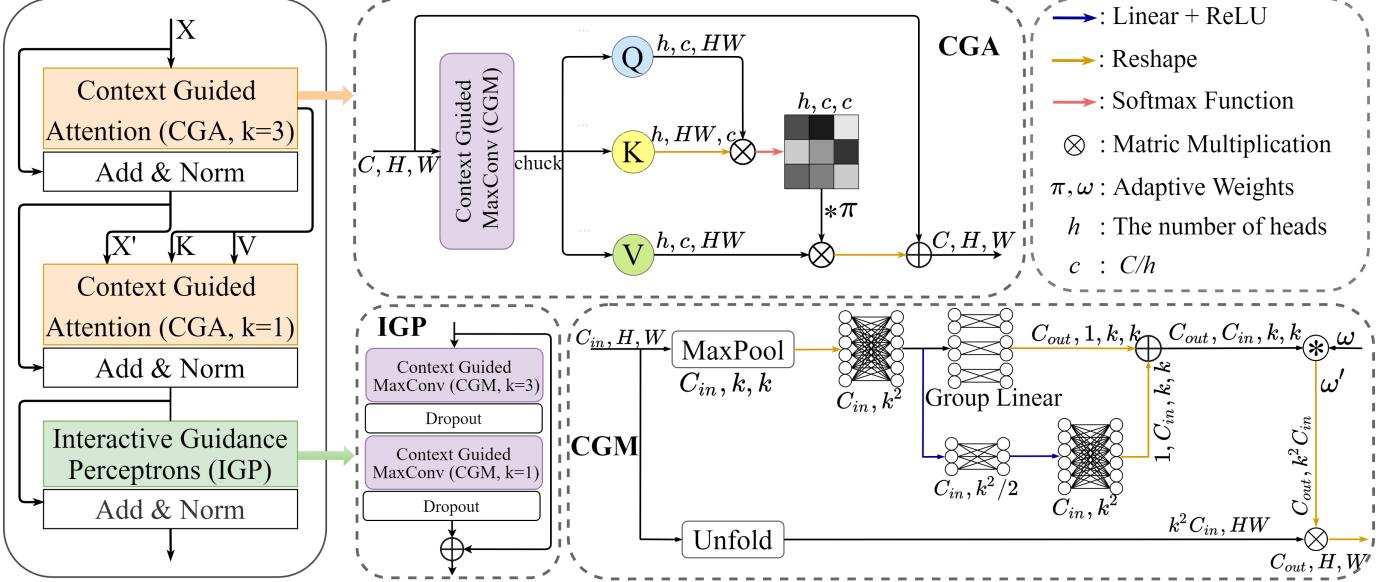
Fig. 5: The architecture of the proposed Cross-receptive Field Guide Transformer (CFGT), Context Guided Attention (CGA), Context Guided MaxConv (CGM), and Interactive Guidance Perceptrons (IGP).

representative local contexts to them, which are obtained by encoding the channel context through CGM $Q = W_{CGM}^{Q}(X)$, $K = W_{CGM}^{K}(X)$, and $V = W_{CGM}^{V}(X)$, $(Q, K, V) \in \Re^{h \times C/h \times HW}$. Among them, $h$ is the amount of attention and $W_{CGM}(\cdot)$ represents the proposed CGM, similar to traditional multi-head attention [46]. In this work, we divide the number of channels into $h$ groups for parallel learning, and the size of $Q, K, V$ is obtained after reshaping the tensor from the input image. Finally, we take the dot-product to reshape the $Q$ and $V$ projections to generate a transposed attention map of size $\Re^{h \times C/h \times C/h}$ instead of a regular feature map of size $\Re^{h \times HW \times HW}$ [47]. Overall, the attention mechanism in CGA can be formulated as

$$\underset{CGA}{Att}(Q, K, V) = V \cdot Soft((K \cdot Q)/\omega), \quad (13)$$

where $Soft(\cdot)$ represents the Softmax function, $\cdot$ represents the dot product operation, $\omega$ represents a learnable adaptive parameter. Due to the powerful function of CGM combined with local features to explicitly modify the modulation kernel, the projection vector generated in the previous space can be used to guide the generation of subsequent attention.

**Context Guided MaxConv (CGM).** Recently, methods for contextual reasoning [48], [49] have been extensively studied in computer vision. Motivated by [24], we propose a Context Guided MaxConv (CGM) to dynamically extract representative local patterns within Transformer in conjunction with contextual guidance. As shown in Fig. 5, to extract the most representative information, we scale the input image to the size of $k \times k$ by using the max-pooling operation. Then, to alleviate the time-consuming kernel modulation caused by a large number of channels, we follow the previous idea and reduce the complexity by generating two tensors through two branches. One of the branches draws on the idea of bottleneck design [37]. We project the spatial position information into

a vector of size $k$ through a linear layer of shared parameters and then generate new channel weights from this vector. In another branch, the idea of grouped convolution is applied to the linear layer, and the output $C_{out}$ is obtained by introducing a grouped linear layer with weight ($w \in \Re^{\frac{C_{in}}{g} \times \frac{C_{out}}{g}}$), where $g$ is the number of groups. Then, spatially transformer the two branches to obtain tensors with corresponding sizes of $C_{out}, 1, k, k$ and $C_{in}, 1, k, k$. Meanwhile, the two tensors are then added element-wise to get the same size as the convolution kernel $\omega(\omega \in \Re^{C_{out} \times C_{out} \times k \times k})$, which is convenient for element-wise multiplication to obtain the modulation convolution kernel. It's worth to noting that $\omega$ is an adaptive multipliers that is consistent with the size of the tensor. After multiplying(*) it with the tensor, the tensor can be converted into a set of trainable type parameters and bound to the module. During the learning process, the weights of $\omega$ can be automatically learned and modified to optimize the model. After that, the two tensors are added element by element to obtain the same size as the convolution kernel $\omega$, which is convenient for element-wise multiplication to obtain the modulation convolution kernel $\omega'$. Finally, after fully understanding the obtained semantic information, this set of optimized tensor $\omega'$ interacts with the input features to dynamically capture the local patterns of salient features.

*2) Cross-receptive Field Guide Transformer (CFGT):* In this paper, we propose a Cross-receptive Field Guide Transformer (CFGT) for long-distance modeling. As shown in Fig. 5, CFGT is mainly composed of two CGAs and one Interactive Guidance Perceptrons (IGP) in the encoder part of CFGT. Meanwhile, hierarchical normalization is performed after each block and the local residual connection is used. It is worth noting that the CGM in the two CGAs adopts different $k$ ($k$ represents the size of the receptive field that the CGM focuses on). Therefore, cross-attention can supplement

the model with cross-scale features. We assume that the input of CFGT is $T_{in}$, then the output $T_{out}$ can be formulated as

$$T_{med}^1 = Norm(CGA(T_{in})) + T_{in}, \qquad (14)$$

$$T_{med}^2 = Norm(CGA(T_{in}, K, V)) + T_{med1}, \qquad (15)$$

$$T_{out} = Norm(IGP(T_{med2})) + T_{med2}, \qquad (16)$$

where $K$ and $V$ are the key and value generated by the first CGA, **and they will serve as part of the input of the second CGA**. The operation of the second CGA can be defined as

$$\underset{CGA}{Att}(Q', K', V') = (V'+V) \cdot Soft(((K'+K) \cdot Q')/\omega). \quad (17)$$

Similarly, our IGP also adopts the idea of cross receptive fields, connecting a CGM with $k = 3$ and a CGM with $k = 1$ using a residual structure to enhance attention of the transformer.

### D. Loss Function

During training, given a training set $\left\{I_i^{LR}, I_i^{HR}\right\}_{i=1}^N$, the loss function of CFIN can be expressed by

$$Loss\,(\theta) = \arg\min_\theta \frac{1}{N} \sum_{i=1}^N \left\| F_{CFIN}(I_{LR}^i) - I_{HR}^i \right\|_1, \quad (18)$$

where $F_{CFIN}(\cdot)$ represent our proposed CFIN, $\theta$ represent the parameter set of CFIN, N represent the number of LR-HR pairs in the training set.

## IV. Experiments

In this part, we provide relevant experimental details, descriptions, and results to verify the effectiveness and excellence of the proposed CFIN.

### A. Datasets and Metrics

Following previous works, we utilize DIV2K [57] (1-800) as our training dataset. Meanwhile, we used five benchmark datasets to verify the effectiveness of the proposed model, including Set5 [58], Set14 [59], BSDS100 [60], Urban100 [61], and Manga109 [62]. Additionally, we used Peak Signal-to-Noise Ratio (PSNR) and Structural Similarity (SSIM) to evaluate the quality of our restored images on the Y channel of the YCbCr color space.

### B. Implementation Details

During training, we randomly crop patches with the size of $48 \times 48$ from the training set as input and use horizontal flipping and random rotation for data augmentation. The initial learning rate is set to $5 \times 10^{-4}$, which is finally reduced to $6.25 \times 10^{-6}$ by cosine annealing. Meanwhile, we implement our model with the PyTorch framework and update it with Adam optimizer. All our experiments are conducted on NVIDIA RTX 2080Ti GPU. In the final model, we use 8 transformer stages and 8 convolution stages, and all of which are called twice using a loop mechanism. Meanwhile, we set the initial input channel to 48 and use the weight normalization [63] after each convolutional layer in RIFUs.

### C. Comparison with Advanced Lightweight SISR Models

In this section, we compare our proposed CFIN with other advanced lightweight SISR models to verify the effectiveness of the proposed model. In TABLE I, we compare CFIN with CNN-based models. From the table, we can clearly observe that our CFIN+ and CFIN stand out from these methods and achieve the best and the second best results on almost all datasets. It is worth mentioning that our CFIN consumes less parameters and computation than most methods. This benefit from the well-designed CNN and Transformer in CFIN, which can efficiently extract the local features and integrate the global information of the image. To further demonstrate the superiority of CFIN, in TABLE II, we also provide a comprehensive comparison with some advanced Transformer-based models, including the lightweight version of SwinIR* [19], ESRT [20], and LBNet [17]. All these models are the most advanced lightweight SISR models. From the table, we can see that our CFIN achieves better results than ESRT and LBNet with fewer parameters. Compared with SwinIR*, our CFIN can still achieve close results than it with fewer parameters, Multi-adds, and execution time. It is worth mentioning that SwinIR* uses a pre-trained model for initialization, and setting the patch size as $64 \times 64$ during training. Extensive experiments have shown that the larger the patch size, the better the results. Meanwhile, some previous works [9], [64] have pointed out that the performance of models trained with multiple upsampling factors shows better results than the single one since using the inter-scale correlation between different upsampling factors can improve the model performance. Therefore, these methods can further improve the model performance. Moreover, we provide the results of the large version of CFIN, donated as CFIN-L. It can be seen that the results of CFIN-L even surpass SwinIR* on some datasets and still keep fewer parameters and less time. This is due to the rational structural design and strong modeling capabilities of CFIN. All these results fully illustrate the strong competitiveness of CFIN in balancing model size and performance.

In addition, we also provide a visual comparison of CFIN with other SISR methods in Fig. 6. Obviously, our CFIN can reconstruct high-quality images with more accurate textures details and edges. This further demonstrate the effectiveness of the proposed CFIN.

## V. Ablation Studies

In this section, we provide a series of ablation studies to further demonstrate the effectiveness of the proposed modules and model.

### A. Network Investigations

*1) The effectiveness of RIFU:* In RIFU, we use feature masks to remove redundant features. To study and verify the effectiveness of this mechanism, we remove the mask and provide the results in TABLE III. According to the table, we can see that the PSNR value increases by 0.08dB after using the masking mechanism under the slight increase in the number of parameters.

TABLE I: Average PSNR/SSIM comparison with other advance CNN-based SISR models. The Best and the second best results are highlighted and underlined, respectively. '+' indicates that the model uses a self-ensemble strategy.

| Method | Scale | Params | Multi-adds | Set5 PSNR/SSIM | Set14 PSNR/SSIM | BSDS100 PSNR/SSIM | Urban100 PSNR/SSIM | Manga109 PSNR/SSIM |
|---|---|---|---|---|---|---|---|---|
| MemNet [50] | | 677K | - | 37.78/0.9597 | 33.28/0.9142 | 32.08/0.8978 | 31.31/0.9195 | - |
| IDN [32] | | 553K | 124.6G | 37.83/0.9600 | 33.30/0.9148 | 32.08/0.8985 | 31.27/0.9196 | 38.01/0.9749 |
| CARN [10] | | 1592K | 222.8G | 37.76/0.9590 | 33.52/0.9166 | 32.09/0.8978 | 31.92/0.9256 | 38.36/0.9765 |
| IMDN [33] | | 694K | 158.8G | 38.00/0.9605 | 33.63/0.9177 | 32.19/0.8996 | 32.17/0.9283 | 38.88/0.9774 |
| AWSRN-M [51] | | 1063K | 244.1G | 38.04/0.9605 | 33.66/0.9181 | 32.21/0.9000 | 32.23/0.9294 | 38.66/0.9772 |
| MADNet [26] | | 878K | 187.1G | 37.85/0.9600 | 33.38/0.9161 | 32.04/0.8979 | 31.62/0.9233 | - |
| MAFFSRN-L [52] | | 790K | 154.4G | 38.07/0.9607 | 33.59/0.9177 | 32.23/0.9005 | 32.38/0.9308 | - |
| LAPAR-A [53] | | 548K | 171.0G | 38.01/0.9605 | 33.62/0.9183 | 32.19/0.8999 | 32.10/0.9283 | 38.67/0.9772 |
| GLADSR [27] | | 812K | 187.2G | 37.99/0.9608 | 33.63/0.9179 | 32.16/0.8996 | 32.16/0.9283 | - |
| LatticeNet+ [54] | ×2 | 756K | 165.5G | 38.15/0.9610 | 33.78/0.9193 | 32.25/0.9005 | 32.43/0.9302 | - |
| DCDN [55] | | 756K | - | 38.01/0.9606 | 33.52/0.9166 | 32.17/0.8996 | 32.16/0.9283 | 38.70/0.9773 |
| SMSR [28] | | 985K | 351.5G | 38.00/0.9601 | 33.64/0.9179 | 32.17/0.8990 | 32.19/0.9284 | 38.76/0.9771 |
| ECBSR [56] | | 596K | 137.3G | 37.90/0.9615 | 33.34/0.9178 | 32.10/0.9018 | 31.71/0.9250 | - |
| PFFN [29] | | 569K | 138.3G | 38.07/0.9607 | 33.69/0.9192 | 32.21/0.8997 | 32.33/0.9298 | 38.89/0.9775 |
| DRSAN [30] | | 690K | 159.3G | 38.11/0.9609 | 33.64/0.9185 | 32.21/0.9005 | 32.35/0.9304 | - |
| FDIWN [36] | | 629K | 112.0G | 38.07/0.9608 | 33.75/0.9201 | 32.23/0.9003 | 32.40/0.9305 | 38.85/0.9774 |
| **CFIN (Ours)** | | 675K | 116.9G | 38.14/0.9610 | 33.80/0.9199 | 32.26/0.9006 | 32.48/0.9311 | 38.91/0.9770 |
| **CFIN+ (Ours)** | | 675K | 116.9G | **38.22/0.9613** | **34.01/0.9221** | **32.35/0.9016** | **32.93/0.9347** | **39.21/0.9777** |
| MemNet [50] | | 677K | - | 34.07/0.9248 | 30.00/0.8350 | 28.96/0.8001 | 27.56/0.8376 | - |
| IDN [32] | | 553K | 56.3G | 34.11/0.9253 | 29.99/0.8354 | 28.95/0.8013 | 27.42/0.8359 | 32.71/0.9381 |
| CARN [10] | | 1592K | 118.8G | 34.29/0.9255 | 30.29/0.8407 | 29.06/0.8034 | 28.06/0.8493 | 33.43/0.9427 |
| IMDN [33] | | 703K | 71.5G | 34.36/0.9270 | 30.32/0.8417 | 29.09/0.8046 | 28.17/0.8519 | 33.61/0.9445 |
| AWSRN-M [51] | | 1143K | 116.6G | 34.42/0.9275 | 30.32/0.8419 | 29.13/0.8059 | 28.26/0.8545 | 33.64/0.9450 |
| MADNet [26] | | 930K | 88.4G | 34.16/0.9253 | 30.21/0.8398 | 28.98/0.8023 | 27.77/0.8439 | - |
| MAFFSRN-L [52] | | 807K | 68.5G | 34.45/0.9277 | 30.40/0.8432 | 29.13/0.8061 | 28.26/0.8552 | - |
| LAPAR-A [53] | | 594K | 114.0G | 34.36/0.9267 | 30.34/0.8421 | 29.11/0.8054 | 28.15/0.8523 | 33.51/0.9441 |
| GLADSR [27] | | 821K | 88.2G | 34.41/0.9272 | 30.37/0.8418 | 29.08/0.8050 | 28.24/0.8537 | - |
| LatticeNet+ [54] | ×3 | 765K | 76.3G | 34.53/0.9281 | 30.39/0.8424 | 29.15/0.8059 | 28.33/0.8538 | - |
| DCDN [55] | | 765K | - | 34.41/0.9273 | 30.31/0.8417 | 29.08/0.8045 | 28.17/0.8520 | 33.54/0.9441 |
| SMSR [28] | | 993K | 156.8G | 34.40/0.9270 | 30.33/0.8412 | 29.10/0.8050 | 28.25/0.8536 | 33.68/0.9445 |
| PFFN [29] | | 558K | 69.1G | 34.54/0.9282 | 30.42/0.8435 | 29.17/0.8062 | 28.37/0.8566 | 33.63/0.9455 |
| EMASRN+ [31] | | 427K | - | 34.48/0.9275 | 30.38/0.8422 | 29.11/0.8046 | 28.17/0.8514 | 33.71/0.9450 |
| DRSAN [30] | | 740K | 76.0G | 34.50/0.9278 | 30.39/0.8437 | 29.13/0.8065 | 28.35/0.8566 | - |
| FDIWN [36] | | 645K | 51.5G | 34.52/0.9281 | 30.42/0.8438 | 29.14/0.8065 | 28.36/0.8567 | 33.77/0.9456 |
| **CFIN (Ours)** | | 681K | 53.5G | 34.65/0.9289 | 30.45/0.8443 | 29.18/0.8071 | 28.49/0.8583 | 33.89/0.9464 |
| **CFIN+ (Ours)** | | 681K | 53.5G | **34.75/0.9298** | **30.59/0.8467** | **29.27/0.8091** | **28.85/0.8645** | **34.26/0.9484** |
| MemNet [50] | | 677K | - | 31.74/0.8893 | 28.26/0.7723 | 27.40/0.7281 | 25.50/0.7630 | - |
| IDN [32] | | 553K | 32.3G | 31.82/0.8903 | 28.25/0.7730 | 27.41/0.7297 | 25.41/0.7632 | 29.41/0.8942 |
| CARN [10] | | 1592K | 90.9G | 32.13/0.8937 | 28.60/0.7806 | 27.58/0.7349 | 26.07/0.7837 | 30.42/0.9070 |
| IMDN [33] | | 715K | 40.9G | 32.21/0.8948 | 28.58/0.7811 | 27.56/0.7353 | 26.04/0.7838 | 30.45/0.9075 |
| AWSRN-M [51] | | 1254K | 72.0G | 32.21/0.8954 | 28.65/0.7832 | 27.60/0.7368 | 26.15/0.7884 | 30.56/0.9093 |
| MADNet [26] | | 1002K | 54.1G | 31.95/0.8917 | 28.44/0.7780 | 27.47/0.7327 | 25.76/0.7746 | - |
| MAFFSRN-L [52] | | 830K | 38.6G | 32.20/0.8953 | 28.62/0.7822 | 27.59/0.7370 | 26.16/0.7887 | - |
| LAPAR-A [53] | | 659K | 94.0G | 32.15/0.8944 | 28.61/0.7818 | 27.61/0.7366 | 26.14/0.7871 | 30.42/0.9074 |
| GLADSR [27] | | 826K | 52.6G | 32.14/0.8940 | 28.62/0.7813 | 27.59/0.7361 | 26.12/0.7851 | - |
| LatticeNet+ [54] | ×4 | 777K | 43.6G | 32.30/0.8962 | 28.68/0.7830 | 27.62/0.7367 | 26.25/0.7873 | - |
| DCDN [55] | | 777K | - | 32.21/0.8949 | 28.57/0.7807 | 27.55/0.7356 | 26.09/0.7855 | 30.41/0.9072 |
| SMSR [28] | | 1006K | 89.1G | 32.12/0.8932 | 28.55/0.7808 | 27.55/0.7351 | 26.11/0.7868 | 30.54/0.9085 |
| ECBSR [56] | | 603K | 34.7G | 31.92/0.8946 | 28.34/0.7817 | 27.48/0.7393 | 25.81/0.7773 | - |
| PFFN [29] | | 569K | 45.1G | 32.36/0.8967 | 28.68/0.7827 | 27.63/0.7370 | 26.26/0.7904 | 30.50/0.9100 |
| EMASRN+ [31] | | 546K | - | 32.31/0.8964 | 28.66/0.7828 | 27.61/0.7364 | 26.15/0.7868 | 30.69/0.9105 |
| DRSAN [30] | | 730K | 49.0G | 32.30/0.8954 | 28.66/0.7838 | 27.61/0.7381 | 26.26/0.7920 | - |
| FDIWN [36] | | 664K | 28.4G | 32.23/0.8955 | 28.66/0.7829 | 27.62/0.7380 | 26.28/0.7919 | 30.63/0.9098 |
| **CFIN (Ours)** | | 699K | 31.2G | 32.49/0.8985 | 28.74/0.7849 | 27.68/0.7396 | 26.39/0.7946 | 30.73/0.9124 |
| **CFIN+ (Ours)** | | 699K | 31.2G | **32.60/0.8998** | **28.86/0.7871** | **27.76/0.7419** | **26.71/0.8028** | **31.15/0.9163** |

As we know, the max-pooling operation can filter redundant features and softmax can predict the probability distribution of different features. These functions also can make the model focus on the main features. To prove the effectiveness of the Gumbel-Softmax in RIFU, we replace the Gumbel-Softmax function with the max-pooling operation and softmax function, respectively. According to TABLE IV, we can see that the performance of the models will drop from 31.75dB to 31.71dB when only using the softmax function, and the performance will drop to 31.70dB when using the max-pooling operation.

Therefore, we choose the Gumbel-Softmax operation in the final model. All the above experiments verify the effectiveness and necessity of each mechanism in RIFU.

*2) The effectiveness of CIAM:* To verify the effectiveness of CIAM, we replace CIAM with some commonly used feature extraction blocks in lightweight SISR models. It is worth noting that we remove the Transformer stage in each CT block to speed up the training. According to TABLE V, we can found that when the model uses CIAM for feature extraction, the parameter increase does not exceed 30K, Multi-

TABLE II: Comparison with some Transformer-base methods for ×4 SR. * means this model is pre-trained based on the ×2 setup and the training patch size is set to 64 × 64 (ours is 48 × 48 and without pre-training).

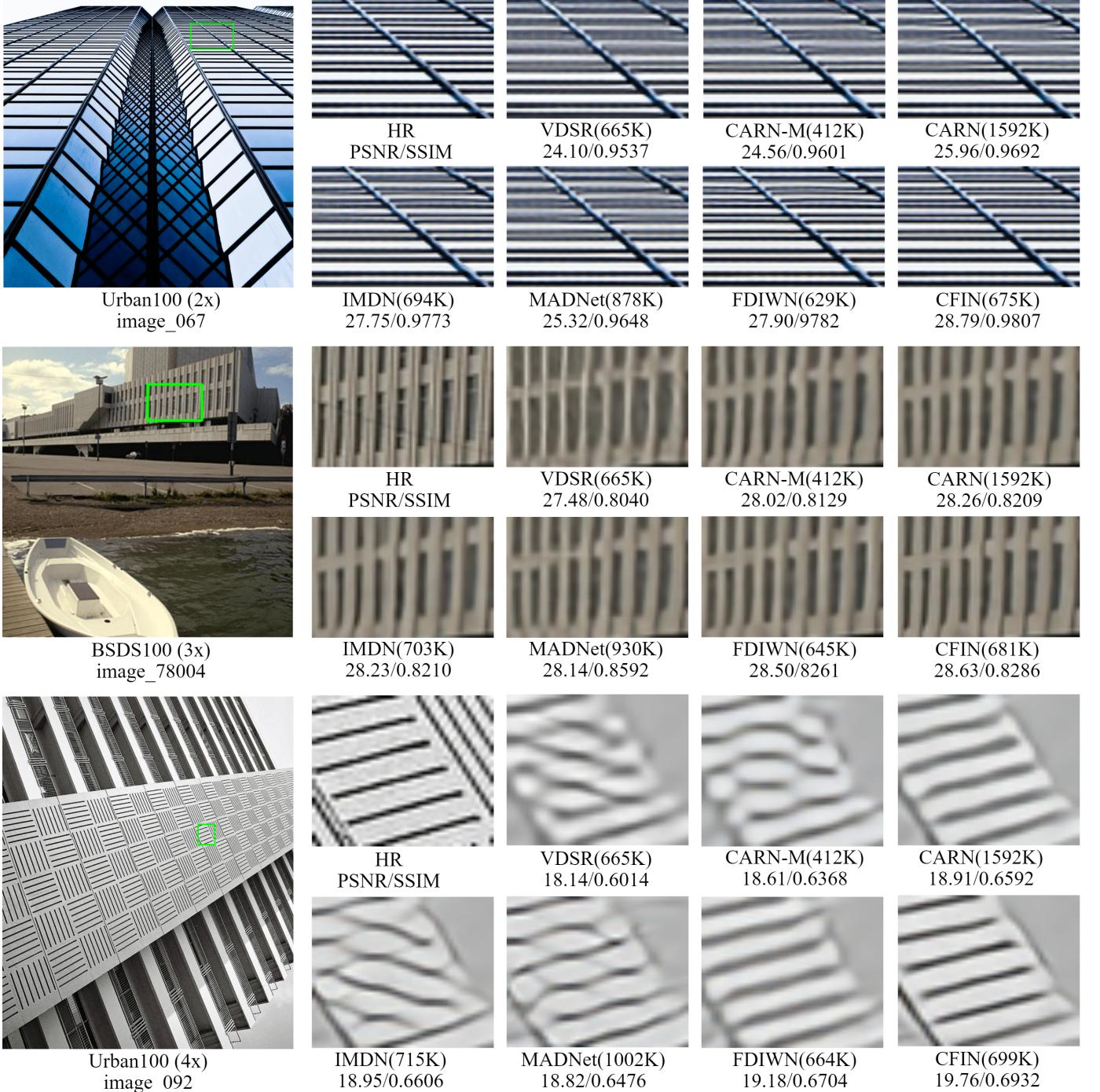| Method | Params | Multi-adds | GPU Memory | Time | Set5 | Set14 | BSD100 | Urban100 | Manga109 | Average |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | PSNR / SSIM | PSNR / SSIM | PSNR / SSIM | PSNR / SSIM | PSNR / SSIM | PSNR / SSIM |
| SwinIR* [19] | 897K | 49.6G | 10500M | 0.046s | 32.44 / 0.8976 | **28.77 / 0.7858** | **27.69 / 0.7406** | 26.47 / **0.7980** | **30.92 / 0.9151** | 29.26 / **0.8274** |
| ESRT [20] | 751K | 67.7G | 4191M | 0.032s | 32.19 / 0.8947 | 28.69 / 0.7833 | **27.69** / 0.7379 | 26.39 / 0.7962 | 30.75 / 0.9100 | 29.14 / 0.8244 |
| LBNet [17] | 742K | 38.9G | 6417M | 0.043s | 32.29 / 0.8960 | 28.68 / 0.7832 | 27.62 / 0.7382 | 26.27 / 0.7906 | 30.76 / 0.9111 | 29.12 / 0.8238 |
| CFIN | 699K | 31.2G | 11419M | 0.035s | 32.49 / 0.8985 | 28.74 / 0.7849 | 27.68 / 0.7396 | 26.39 / 0.7946 | 30.73 / 0.9124 | 29.21 / 0.8260 |
| CFIN-L | 852K | 37.8G | 14585M | 0.040s | **32.56 / 0.8988** | 28.74 / 0.7852 | **27.69 / 0.7406** | **26.49** / 0.7973 | 30.85 / 0.9134 | **29.27** / 0.8271 |



Fig. 6: Visual comparison with different lightweight SISR models. Obviously, our CFIN can reconstruct high-quality images.

TABLE III: Evaluate the effectiveness of masking mechanism.

| Scale | Mask | Params | Multi-adds | Set14 |
|---|---|---|---|---|
| | | | | PSNR / SSIM |
| ×4 | ✗ | 186.8K | 6.18220G | 28.22 / 07723 |
| | ✓ | 187.9K | 6.21538G | **28.30 / 0.7738** |

TABLE IV: Evaluate the effectiveness of Gumbel-Softmax.

| Scale | Module | Params | Multi-adds | Set5 |
|---|---|---|---|---|
| | | | | PSNR / SSIM |
| | RIFU+Maxpool | 165K | 4.89G | 31.70 / 0.8876 |
| ×4 | RIFU+Softmax | 165K | 4.89G | 31.71 / 0.8877 |
| | RIFU+Gumbel-Softmax | 165K | 4.89G | **31.75 / 0.8886** |

adds is less than most methods, but the model achieves the best results. This fully demonstrates the effectiveness of the proposed CIAM.

It is worth noting that the upsampling-downsampling branch in CIAM is also extremely important. Specifically, the performance (×4, Set5) of the ablation model (165K) with and without the upsampling-downsampling branch is 31.75dB and 31.63dB, respectively. But the position of this branch will not have much impact on model performance.

TABLE V: Evaluate the effectiveness of CIAM.

| Module | Params | Multi-adds | Set5 | Set14 | U100 |
|---|---|---|---|---|---|
| CFIN+RCAB [65] | 192K | 17.35G | 31.63 | 28.19 | 25.29 |
| CFIN+IMDB [33] | 148K | 8.47G | 31.54 | 28.19 | 25.23 |
| CFIN+RFDB [34] | 140K | 7.77G | 31.68 | 28.26 | 25.34 |
| CFIN+LB [54] | 145K | 8.16G | 31.66 | 28.22 | 25.33 |
| CFIN+HPB [20] | 150K | 8.91G | 31.59 | 28.21 | 25.32 |
| CFIN+WDIB [36] | 147K | 5.26G | 31.68 | 28.26 | 25.37 |
| CFIN+CIAM (Ours) | 188K | 6.22G | **31.87** | **28.30** | **25.44** |

*3) The effectiveness of CGM:* Context Guided MaxConv is the basic unit in our Transformer stage, which is responsible for reasoning and guiding the entire network. To verify its importance, we replace CGM with the linear layer and group convolution, respectively. Meanwhile, we set their parameters as close as possible. According to TABLE VI, we can see that our CGM achieves better performance with fewer parameters and Multi-adds. This validates the excellence and effectiveness of the proposed CGM.

TABLE VI: Evaluate the effectiveness of CGM.

| Module | Params | Multi-adds | Set5 | B100 | U100 |
|---|---|---|---|---|---|
| CFIN+GConv | 197K | 5.7981G | 31.95 | 27.41 | 25.68 |
| CFIN+Linear | 198K | 4.1104G | 29.19 | 26.37 | 23.62 |
| CFIN+CGM (Ours) | 178K | 4.1105G | **32.14** | **27.43** | **25.71** |

*4) The effectiveness of CFGT:* Compared with traditional Transformer, CFGT has two main differences in structure: the first is K, V for context communication; the second is cross-scale information exchange brought by CGM with different receptive fields. To assess the effectiveness of such a design, we provide a set of ablation experiments in TABLE VII, where **KV** stands for the contextual communication mechanism and **Cross** stands for the cross receptive fields mechanism. It can be seen that our architectural design solution can significantly improve the performance of the model with almost no addi-

TABLE VII: Evaluate the effectiveness of CFGT.

| Scale | KV | Cross | Params | Multi-adds | Set5 |
|---|---|---|---|---|---|
| | | | | | PSNR / SSIM |
| | ✗ | ✓ | 177.8K | 4.1105G | 32.05 / 0.8928 |
| ×4 | ✓ | ✗ | 177.5K | 4.1104G | 31.93 / 0.8904 |
| | ✓ | ✓ | 177.8K | 4.1105G | **32.14 / 0.8940** |

tional computational cost. This validates the effectiveness of the proposed CFGT, which can better convey context.

Meanwhile, to verify that CFGT can effectively adjust attention according to the context, we make an efficiency trade-off on the number of CFGTs. In TABLE VIII, CFGT-N indicates that N CFGTs are used in the model. It can be seen that our model achieves more than 0.4dB improvement on the Manga109 at the cost of less than 70K parameters compared to the model without CFGT. In addition, we present the visual heatmap of different numbers of CFGTs in Fig. 7. Among them, the color close to red in the image represents the part where the attention is focused. So we can see that in the absence of CFGT, the attention is only focused on a small number of texture features, and as the number of CFGTs increases, more fine-grained features in images are paid attention to. This means that more CFGT makes the model tends to recover more accurate detailed features, which is conducive to image restoration.

TABLE VIII: Efficiency trade-off of CFGT.

| Module | Params | Multi-adds | Set5 | Manga109 | Memory |
|---|---|---|---|---|---|
| CFGT-0 | 630K | 28.0G | 32.11 | 30.29 | 3551M |
| CFGT-4 | 665K | 29.6G | 32.27 | 30.49 | 7485M |
| CFGT-8 | 699K | 31.2G | **32.49** | **30.73** | 11419M |



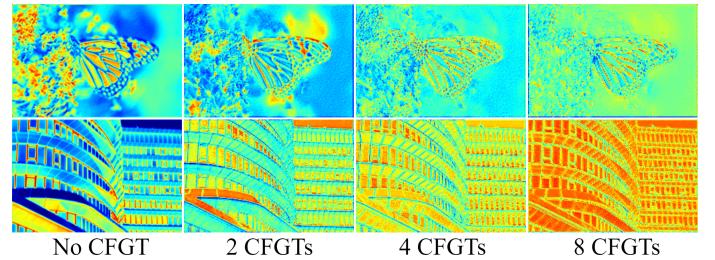| No CFGT | 2 CFGTs | 4 CFGTs | 8 CFGTs |

Fig. 7: The effect of CFGT on visual activation maps.

In TABLE IX, OnlyT indicates that only the Transformer part of the CFIN is retained. Compared with other advanced lightweight SISR models, it can be seen that our CFGT achieves promising results with only 91K parameters. This further verifies the effectiveness and advancement of CFIN.

TABLE IX: Comparison of CFGT with other methods.

| Module | Params | Multi-adds | Set5 | Set14 | B100 | U100 |
|---|---|---|---|---|---|---|
| PAN [66] | 272K | 28.2G | 32.13 | 28.61 | **27.59** | **26.11** |
| MAFFSRN [52] | 441K | 19.3G | 32.18 | 28.58 | 27.57 | 26.04 |
| RFDN [34] | 550K | 31.6G | 32.24 | 28.61 | 27.57 | 26.11 |
| FDIWN-M [36] | 454K | 19.6G | 32.17 | 28.55 | 27.58 | 26.02 |
| OnlyT (Ours) | 91K | 4.5G | **32.33** | **28.61** | 27.58 | 26.08 |

TABLE X: Evaluate the feasibility of combining CNN with Transformer.

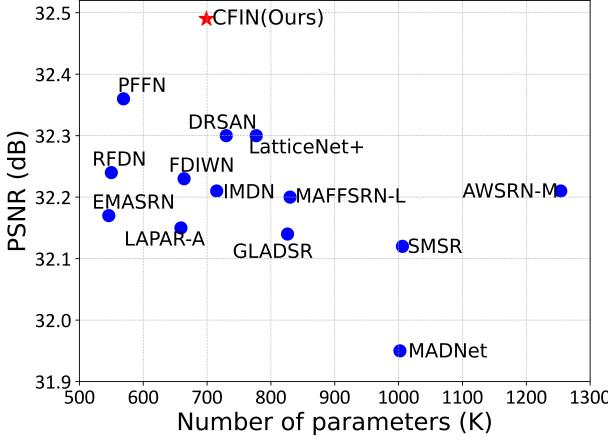| Scale | Method | Params | GPU Memory | Time | Set5 PSNR / SSIM | Set14 PSNR / SSIM | BSD100 PSNR / SSIM | Urban100 PSNR / SSIM | Manga109 PSNR / SSIM |
|---|---|---|---|---|---|---|---|---|---|
| ×4 | Pure-CNN | 720K | 3845M | 0.016s | 32.12 / 0.8939 | 28.52 / 0.7796 | 27.52 / 0.7343 | 25.96 / 0.7810 | 30.29 / 0.9057 |
| | Pure-Transformer | 94K | 11549M | 0.038s | 32.48 / 0.8980 | 28.69 / 0.7835 | 27.64 / 0.7385 | 26.24 / 0.7896 | 30.62 / 0.9102 |
| | CFIN (Ours) | 699K | 11419M | 0.035s | **32.49 / 0.8985** | **28.74 / 0.7849** | **27.68 / 0.7396** | **26.39 / 0.7946** | **30.73 / 0.9124** |



Fig. 8: Model performance and size comparison on Set5 (×4). Obviously, our CFIN achieves the best balance between model performance and size.
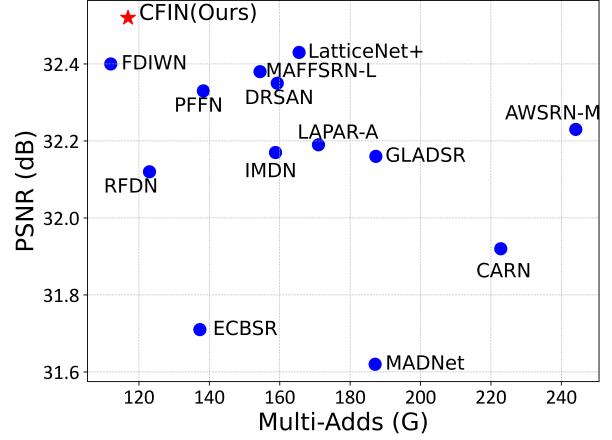


Fig. 9: Model performance and Multi-Adds comparison on Urban100 (×2). Obviously, our CFIN achieves the best balance between model performance and Multi-Adds.

### B. Complementarity of CNNs and Transformers

CNN can be used to extract local features, and Transformer has powerful global modeling capabilities, both of which are crucial for high-quality image restoration. In TABLE X, we provide a comparison of CFIN with Pure-CNN and Pure-Transformer versions. Among them, Pure-CNN and Pure-Transformer represent variant models with only the CNN part or the Transformer part, respectively. For a fair comparison, the number of parameters of Pure-CNN is set as close to CFIN as possible, and the memory consumption of Pure-Transformer is set as close to CFIN as possible. It can be seen from the table that neither Pure-CNN nor Pure-Transformer cannot achieve the performance of the original CFIN. When using CNN alone, it can reduce the consumption of GPU memory, but it is difficult to improve the performance of the model even with more parameters. When using Transformer alone, the parameters of the model can be greatly reduced, but the consumption of GPU memory will rise rapidly, and its performance still hardly exceeds the original CFIN, which is not conducive to the practical application of the model. Therefore, we chose the hybrid model of CNN and Transformer, which can achieve a good balance between the size, GPU memory consumption, and performance of the model. Therefore, we can draw a conclusion that CNN and Transformer are complementary, and the combination of these two parts is feasible.

### C. Model Complexity Analysis

We present the trade-off between our CFIN and other advanced SISR models in terms of PSNR, parameter amount, and inference time in Fig. 1. Obviously, CFIN attains the best performance among models with similar execution times and

achieves the best balance in model complexity, inference time, and performance.

In Fig. 8 and Fig. 9, we also provide the parameter, Multi-Adds, and performance comparisons of CFIN with other advanced SISR models. It can be seen that our CFIN also achieved the best PSNR results under the premise of comparable calculations. Therefore, we can draw a conclusion that our CFIN is a lightweight and efficient model, which achieves the best balance between the model size and performance.

### VI. DISCUSSION

Although CFIN is an efficient model and achieve promising results, it still has some limitations and some details that need to pay attention:

(1) In this paper, we propose a lightweight and efficient CFIN for SISR. CFIN elegantly integrates CNN and Transformer, achieving competitive performance and a good balance between computational cost and performance. Different from the previous work [67], our CGA sent $Q$ and $K$ generated by the first CGA to the second CGA, and then add them with the new $Q$ and $K$ generated in the second CGA. After that, the result of the dot product of newly generated $K$ and $Q$ are dot products with $V$ to obtain the final results. And the sizes of the kernel used in our two CGAs are different, thus these two modules have different receptive fields and can extract features with different scales. Therefore, through our method, features with different scales can be fused and guided with each other, which helps to make full use of image features to reconstruct more accurate images.

(2) We use some matrix multiplication in our model (especially matrix multiplication in CGA), so the training memory

is slightly larger. However, this does not mean the proposed method is meaningless. In summary, we proposed a new Transformer that can use modulated convolution kernel combined with locally representative semantic information to adaptively modify network weights, which can achieve excellent SR performance. We also note that some researchers [20] claimed that their methods can greatly reduce the memory required for Transformer training. In future works, We will further explore the effectiveness of these strategies and introduce them into CFIN to further improve the model.

## VII. Conclusions

In this paper, we proposed a lightweight and efficient Cross-receptive Focused Inference Network (CFIN) for SISR. Specifically, we proposed a Cross-scale Information Aggregation Module (CIAM) composed of Redundant Information Filter Units (RIFUs) for filtering redundant information. Meanwhile, a Cross-receptive Field Guide Transformer (CFGT) is proposed for fine-grained information learning. In CFGT, Context Guided Max (CGM) is introduced to dynamically adjust the attention combined with neurons, and the cross-connection strategy is used to fuse context information across scales to meet the requirements of long-distance modeling. Extensive experiments have shown that CFIN can strike a good balance between the performance and complexity of the model.

## References

[1] Y. Hu, J. Li, Y. Huang, and X. Gao, "Channel-wise and spatial feature modulation network for single image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 30, no. 11, pp. 3911–3927, 2020.

[2] Z.-S. Liu, W.-C. Siu, and Y.-L. Chan, "Photo-realistic image super-resolution via variational autoencoders," *IEEE Transactions on Circuits and Systems for video Technology*, vol. 31, no. 4, pp. 1351–1365, 2021.

[3] J. Zhang, C. Long, Y. Wang, H. Piao, H. Mei, X. Yang, and B. Yin, "A two-stage attentive network for single image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 3, pp. 1020–1033, 2022.

[4] C. Dong, C. C. Loy, K. He, and X. Tang, "Image super-resolution using deep convolutional networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 38, no. 2, pp. 295–307, 2015.

[5] M. Haris, G. Shakhnarovich, and N. Ukita, "Deep back-projection networks for super-resolution," in *CVPR*, 2018, pp. 1664–1673.

[6] Y. Zhang, Y. Tian, Y. Kong, B. Zhong, and Y. Fu, "Residual dense network for image super-resolution," in *CVPR*, 2018, pp. 2472–2481.

[7] X. Chu, B. Zhang, and R. Xu, "Multi-objective reinforced evolution in mobile neural architecture search," in *ECCV*, 2020, pp. 99–113.

[8] J. Li, F. Fang, K. Mei, and G. Zhang, "Multi-scale residual network for image super-resolution," in *ECCV*, 2018, pp. 517–532.

[9] J. Li, F. Fang, J. Li, K. Mei, and G. Zhang, "Mdcn: Multi-scale dense cross network for image super-resolution," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 7, pp. 2547–2561, 2020.

[10] N. Ahn, B. Kang, and K.-A. Sohn, "Fast, accurate, and lightweight super-resolution with cascading residual network," in *ECCV*, 2018, pp. 252–268.

[11] Z. Wang, G. Gao, J. Li, Y. Yu, and H. Lu, "Lightweight image super-resolution with multi-scale feature interaction network," in *ICME*, 2021, pp. 1–6.

[12] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *ECCV*, 2020, pp. 213–229.

[13] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *ICCV*, 2021, pp. 10 012–10 022.

[14] H. Chen, Y. Wang, T. Guo, C. Xu, Y. Deng, Z. Liu, S. Ma, C. Xu, C. Xu, and W. Gao, "Pre-trained image processing transformer," in *CVPR*, 2021, pp. 12 299–12 310.

[15] R. Yu, D. Du, R. LaLonde, D. Davila, C. Funk, A. Hoogs, and B. Clipp, "Cascade transformers for end-to-end person search," in *CVPR*, June 2022, pp. 7267–7276.

[16] Z. Wang, X. Cun, J. Bao, W. Zhou, J. Liu, and H. Li, "Uformer: A general u-shaped transformer for image restoration," in *CVPR*, June 2022, pp. 17 683–17 693.

[17] G. Gao, Z. Wang, J. Li, W. Li, Y. Yu, and T. Zeng, "Lightweight bimodal network for single-image super-resolution via symmetric cnn and recursive transformer," *IJCAI*, 2022.

[18] X. Chen, X. Wang, J. Zhou, and C. Dong, "Activating more pixels in image super-resolution transformer," *arXiv preprint arXiv:2205.04437*, 2022.

[19] J. Liang, J. Cao, G. Sun, K. Zhang, L. Van Gool, and R. Timofte, "Swinir: Image restoration using swin transformer," in *ICCVW*, 2021, pp. 1833–1844.

[20] Z. Lu, J. Li, H. Liu, C. Huang, L. Zhang, and T. Zeng, "Transformer for single image super-resolution," in *CVPR Workshop*, 2022, pp. 457–466.

[21] C. D. Gilbert and W. Li, "Top-down influences on visual processing," *Nature Reviews Neuroscience*, vol. 14, no. 5, pp. 350–363, 2013.

[22] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool, "Dynamic filter networks," in *NeurIPS*, 2016, pp. 20.1–29.9.

[23] Y. Chen, X. Dai, M. Liu, D. Chen, L. Yuan, and Z. Liu, "Dynamic convolution: Attention over convolution kernels," in *CVPR*, 2020, pp. 11 030–11 039.

[24] X. Lin, L. Ma, W. Liu, and S.-F. Chang, "Context-gated convolution," in *ECCV*, 2020, pp. 701–718.

[25] A. Lahiri, S. Bairagya, S. Bera, S. Haldar, and P. K. Biswas, "Lightweight modules for efficient deep learning based image restoration," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 31, no. 4, pp. 1395–1410, 2021.

[26] R. Lan, L. Sun, Z. Liu, H. Lu, C. Pang, and X. Luo, "Madnet: a fast and lightweight network for single-image super resolution," *IEEE Transactions on Cybernetics*, vol. 51, no. 3, pp. 1443–1453, 2020.

[27] X. Zhang, P. Gao, S. Liu, K. Zhao, G. Li, L. Yin, and C. W. Chen, "Accurate and efficient image super-resolution via global-local adjusting dense network," *IEEE Transactions on Multimedia*, vol. 23, pp. 1924–1937, 2021.

[28] L. Wang, X. Dong, Y. Wang, X. Ying, Z. Lin, W. An, and Y. Guo, "Exploring sparsity in image super-resolution for efficient inference," in *CVPR*, 2021, pp. 4917–4926.

[29] D. Zhang, C. Li, N. Xie, G. Wang, and J. Shao, "Pffn: Progressive feature fusion network for lightweight image super-resolution," in *ACMMM*, 2021, pp. 3682–3690.

[30] K. Park, J. W. Soh, and N. I. Cho, "Dynamic residual self-attention network for lightweight single image super-resolution," *IEEE Transactions on Multimedia*, 2021.

[31] X. Zhu, K. Guo, S. Ren, B. Hu, M. Hu, and H. Fang, "Lightweight image super-resolution with expectation-maximization attention mechanism," *IEEE Transactions on Circuits and Systems for Video Technology*, vol. 32, no. 3, pp. 1273–1284, 2022.

[32] Z. Hui, X. Wang, and X. Gao, "Fast and accurate single image super-resolution via information distillation network," in *CVPR*, 2018, pp. 723–731.

[33] Z. Hui, X. Gao, Y. Yang, and X. Wang, "Lightweight image super-resolution with information multi-distillation network," in *ACMMM*, 2019, pp. 2024–2032.

[34] J. Liu, J. Tang, and G. Wu, "Residual feature distillation network for lightweight image super-resolution," in *ECCV*, 2020, pp. 41–55.

[35] X. Jiang, N. Wang, J. Xin, X. Xia, X. Yang, and X. Gao, "Learning lightweight super-resolution networks with weight pruning," *Neural Networks*, vol. 144, pp. 21–32, 2021.

[36] G. Gao, W. Li, J. Li, F. Wu, H. Lu, and Y. Yu, "Feature distillation interaction weighting network for lightweight image super-resolution," *arXiv preprint arXiv:2112.08655*, 2021.

[37] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *CVPR*, 2018, pp. 7132–7141.

[38] F. Wu, A. Fan, A. Baevski, Y. N. Dauphin, and M. Auli, "Pay less attention with lightweight and dynamic convolutions," *arXiv preprint arXiv:1901.10430*, 2019.

[39] X. Zhu, H. Hu, S. Lin, and J. Dai, "Deformable convnets v2: More deformable, better results," in *CVPR*, 2019, pp. 9308–9316.

[40] Y. Jo, S. W. Oh, J. Kang, and S. J. Kim, "Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation," in *CVPR*, 2018, pp. 3224–3232.

[41] Y. Yang, Z. Zhong, T. Shen, and Z. Lin, "Convolutional neural networks with alternately updated clique," in *CVPR*, 2018, pp. 2413–2422.

[42] C. Xie, Y. Wu, L. v. d. Maaten, A. L. Yuille, and K. He, "Feature denoising for improving adversarial robustness," in *CVPR*, 2019, pp. 501–509.

[43] R. Zhang, "Making convolutional networks shift-invariant again," in *ICML*, 2019, pp. 7324–7334.

[44] S. A. Magid, Y. Zhang, D. Wei, W.-D. Jang, Z. Lin, Y. Fu, and H. Pfister, "Dynamic high-pass filtering and multi-spectral attention for image super-resolution," in *ICCV*, 2021, pp. 4288–4297.

[45] E. Jang, S. Gu, and B. Poole, "Categorical reparameterization with gumbel-softmax," *arXiv preprint arXiv:1611.01144*, 2016.

[46] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *NeurIPS*, 2017, pp. 30.1–30.11.

[47] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[48] X. Chen, L.-J. Li, L. Fei-Fei, and A. Gupta, "Iterative visual reasoning beyond convolutions," in *CVPR*, 2018, pp. 7239–7248.

[49] K. Li, Y. Zhang, K. Li, Y. Li, and Y. Fu, "Visual semantic reasoning for image-text matching," in *ICCV*, 2019, pp. 4654–4662.

[50] Y. Tai, J. Yang, and X. Liu, "Image super-resolution via deep recursive residual network," in *CVPR*, 2017, pp. 3147–3155.

[51] C. Wang, Z. Li, and J. Shi, "Lightweight image super-resolution with adaptive weighted learning network," *arXiv preprint arXiv:1904.02358*, 2019.

[52] A. Muqeet, J. Hwang, S. Yang, J. Kang, Y. Kim, and S.-H. Bae, "Multi-attention based ultra lightweight image super-resolution," in *ECCV*, 2020, pp. 103–118.

[53] W. Li, K. Zhou, L. Qi, N. Jiang, J. Lu, and J. Jia, "Lapar: Linearly-assembled pixel-adaptive regression network for single image super-resolution and beyond," in *NeurIPS*, 2020, pp. 20 343–20 355.

[54] X. Luo, Y. Xie, Y. Zhang, Y. Qu, C. Li, and Y. Fu, "Latticenet: Towards lightweight image super-resolution with lattice block," in *ECCV*, 2020, pp. 272–289.

[55] Y. Li, J. Cao, Z. Li, S. Oh, and N. Komuro, "Lightweight single image super-resolution with dense connection distillation network," *ACM TOMM*, vol. 17, no. 1s, pp. 1–17, 2021.

[56] X. Zhang, H. Zeng, and L. Zhang, "Edge-oriented convolution block for real-time super resolution on mobile devices," in *ACMMM*, 2021, pp. 4034–4043.

[57] R. Timofte, E. Agustsson, L. Van Gool, M.-H. Yang, and L. Zhang, "Ntire 2017 challenge on single image super-resolution: Methods and results," in *CVPRW*, 2017, pp. 114–125.

[58] M. Bevilacqua, A. Roumy, C. Guillemot, and M. L. Alberi-Morel, "Low-complexity single-image super-resolution based on nonnegative neighbor embedding," in *BMVC*, 2012, pp. 135.1–135.10.

[59] R. Zeyde, M. Elad, and M. Protter, "On single image scale-up using sparse-representations," in *ICCS*, 2010, pp. 711–730.

[60] D. Martin, C. Fowlkes, D. Tal, and J. Malik, "A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics," in *ICCV*, 2001, pp. 416–423.

[61] J.-B. Huang, A. Singh, and N. Ahuja, "Single image super-resolution from transformed self-exemplars," in *CVPR*, 2015, pp. 5197–5206.

[62] Y. Matsui, K. Ito, Y. Aramaki, A. Fujimoto, T. Ogawa, T. Yamasaki, and K. Aizawa, "Sketch-based manga retrieval using manga109 dataset," *Multimedia Tools and Applications*, vol. 76, no. 20, pp. 21 811–21 838, 2017.

[63] T. Salimans and D. P. Kingma, "Weight normalization: A simple reparameterization to accelerate training of deep neural networks," in *NeurIPS*, 2016, pp. 29.1–29.9.

[64] B. Lim, S. Son, H. Kim, S. Nah, and K. Mu Lee, "Enhanced deep residual networks for single image super-resolution," in *CVPRW*, 2017, pp. 136–144.

[65] Y. Zhang, K. Li, K. Li, L. Wang, B. Zhong, and Y. Fu, "Image super-resolution using very deep residual channel attention networks," in *ECCV*, 2018, pp. 286–301.

[66] H. Zhao, X. Kong, J. He, Y. Qiao, and C. Dong, "Efficient image super-resolution using pixel attention," in *ECCVW*, 2020, pp. 56–72.

[67] R. Song, W. Ni, W. Cheng, and X. Wang, "Csanet: Cross-temporal interaction symmetric attention network for hyperspectral image change detection," *IEEE Geoscience and Remote Sensing Letters*, 2022.