# iiTransformer: A Unified Approach to Exploiting Local and Non-Local Information for Image Restoration
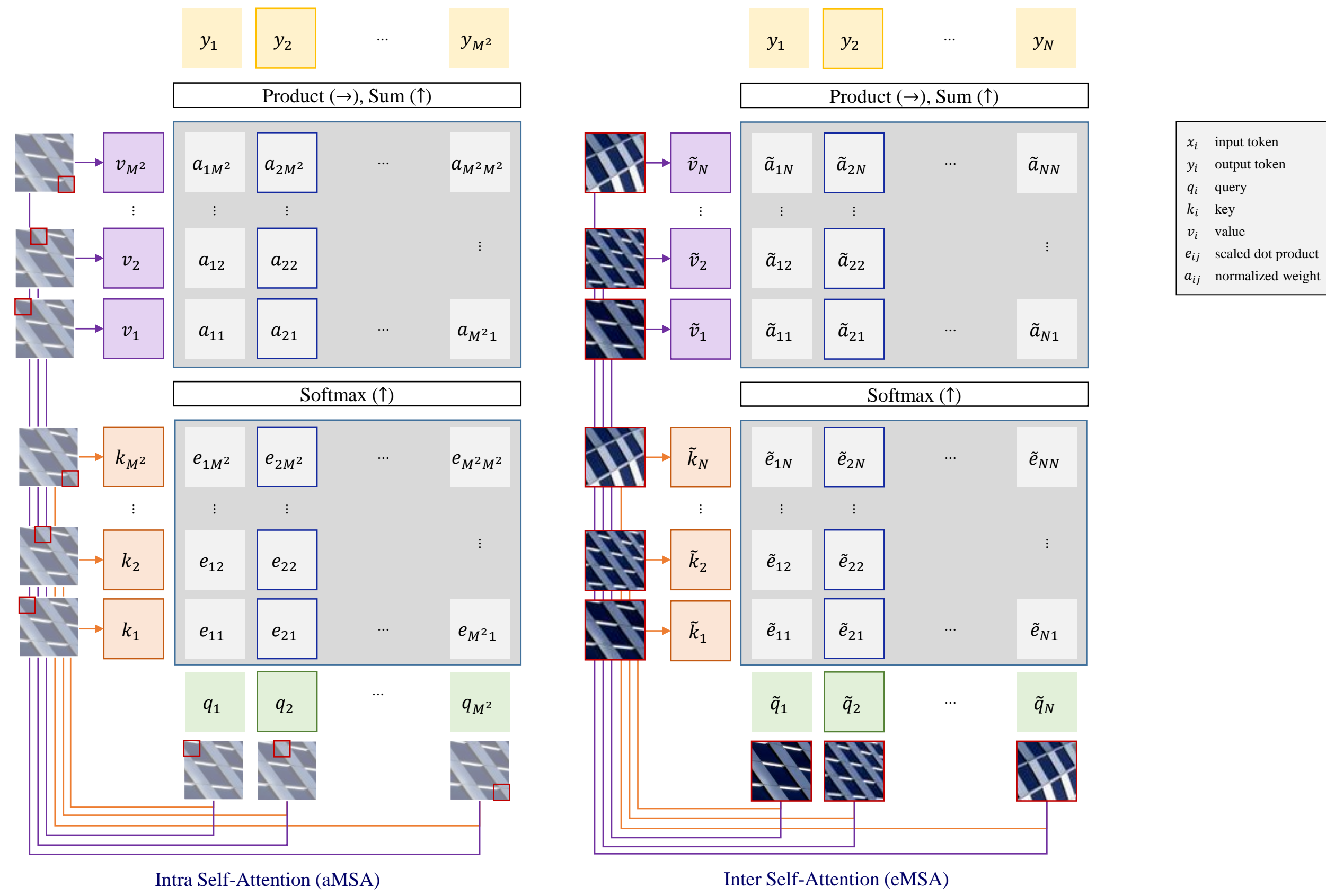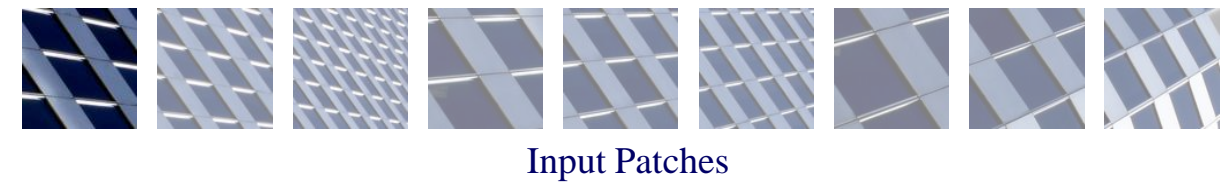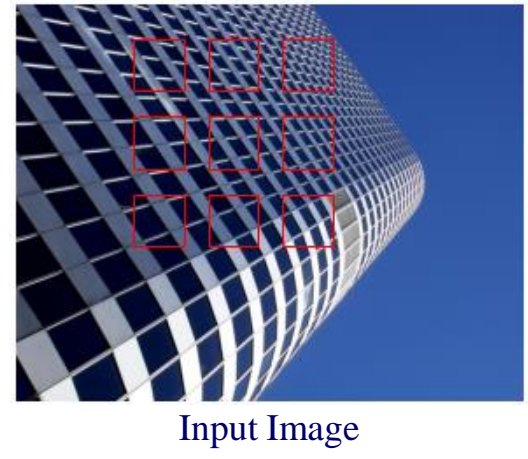
BMVC 2022

Soo Min Kang, Youngchan Song, Hanul Shin, and Tammy Lee

Samsung Research

https://github.com/SamsungLabs/iiTransformer

## Introduction

1. Pixels surrounding the degraded pixel often provide useful information, and
2. images in general contain repetitive information

- Local and non-local relationships can be captured by considering long-range dependencies at the pixel- or patch-level using the self-attention module of Transformers



- Existing inter SA module-based ViT require resolution of training and inference images to match, resulting in either
  (i) increased computational complexity due to overlapping sliding windows, or
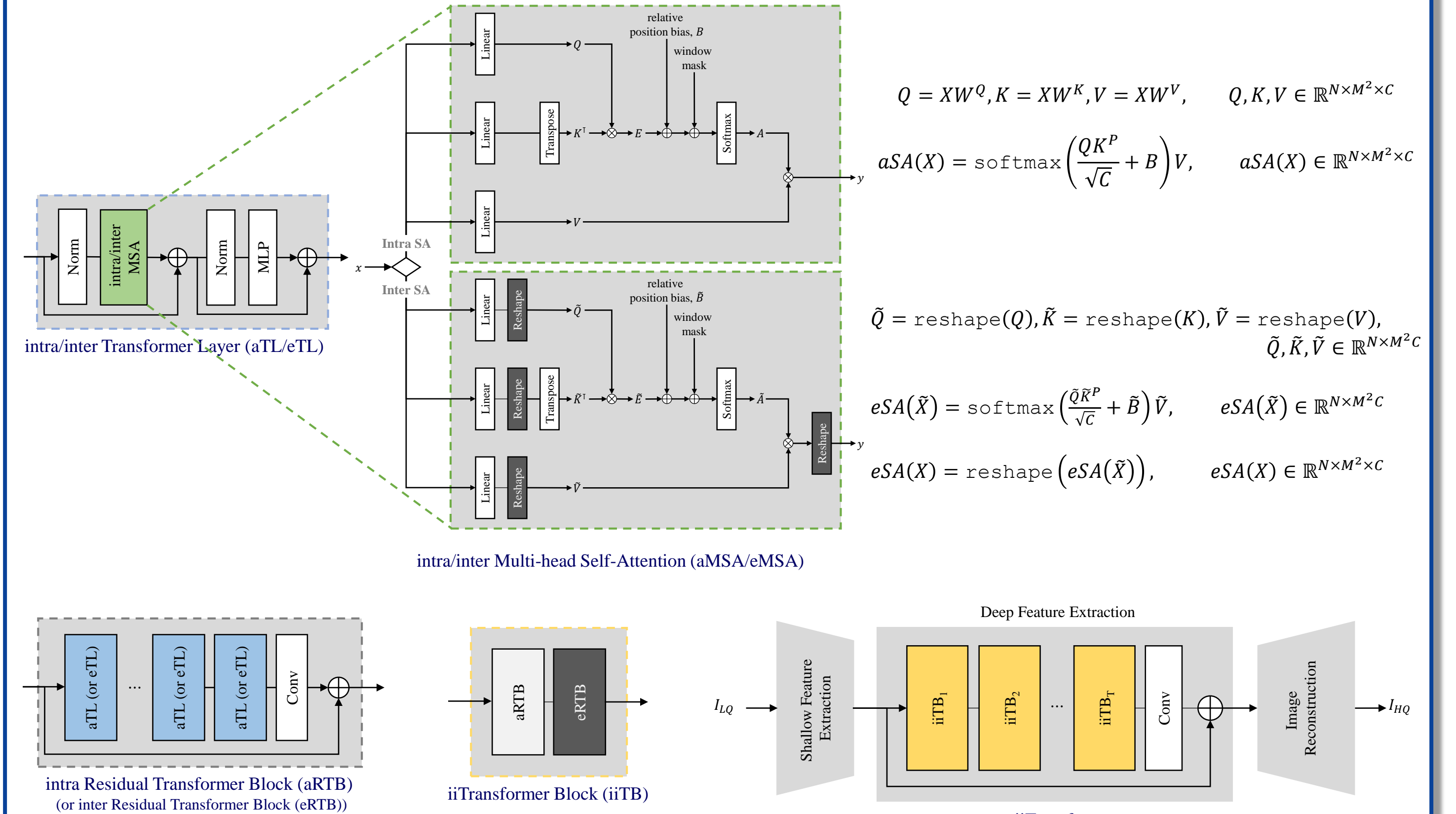  (ii) patch boundary artifacts



(a) Noisy Input Image (LQ)  (b) Presence of Boundary Artifacts  (c) Absence of Boundary Artifacts
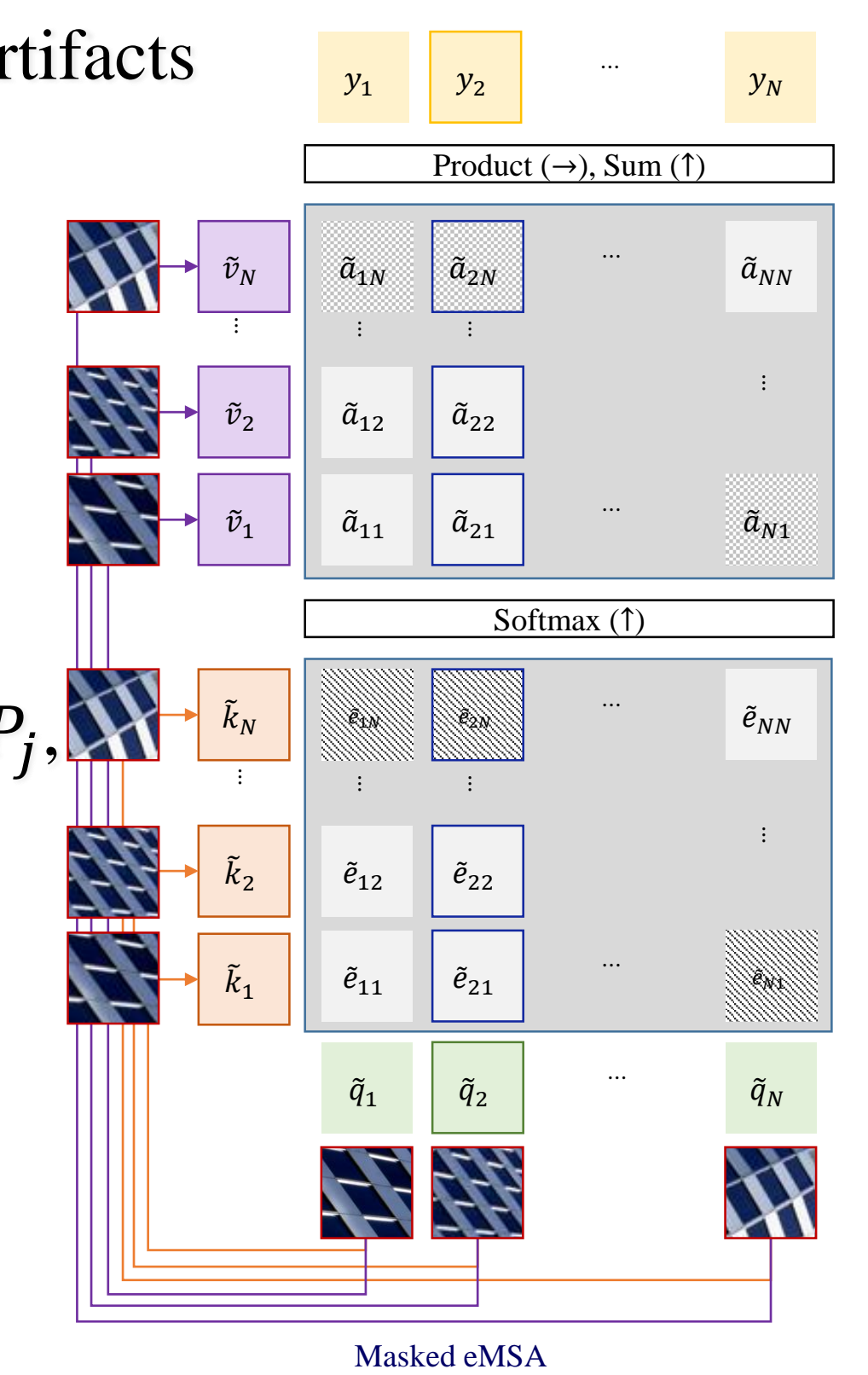
## Methodology

1. Exploiting local and non-local information for image restoration
   - intra SA (aMSA): treat <u>pixels</u> as tokens to compute <u>local pixelwise</u> correlations
   - inter SA (eMSA): treat <u>patches</u> as tokens to compute <u>non-local patchwise</u> correlations
   - <u>shape</u> of projected tokens used to compute the attention matrix differ between aMSA and eMSA



$$Q = XW^Q, K = XW^K, V = XW^V, \qquad Q,K,V \in \mathbb{R}^{N \times M^2 \times C}$$

$$aSA(X) = \text{softmax}\left(\frac{QK^P}{\sqrt{C}} + B\right)V, \qquad aSA(X) \in \mathbb{R}^{N \times M^2 \times C}$$

$$\bar{Q} = \text{reshape}(Q), \bar{K} = \text{reshape}(K), \bar{V} = \text{reshape}(V), \qquad \bar{Q}, \bar{K}, \bar{V} \in \mathbb{R}^{N \times M^2 C}$$

$$eSA(\bar{X}) = \text{softmax}\left(\frac{\bar{Q}\bar{R}^P}{\sqrt{C}} + \bar{B}\right)\bar{V}, \qquad eSA(\bar{X}) \in \mathbb{R}^{N \times M^2 C}$$

$$eSA(X) = \text{reshape}\left(eSA(\bar{X})\right), \qquad eSA(X) \in \mathbb{R}^{N \times M^2 \times C}$$

2. Support arbitrary resolutions without boundary artifacts by masking patch-to-patch distance that exceeds furthest patch-to-patch distance used during training:
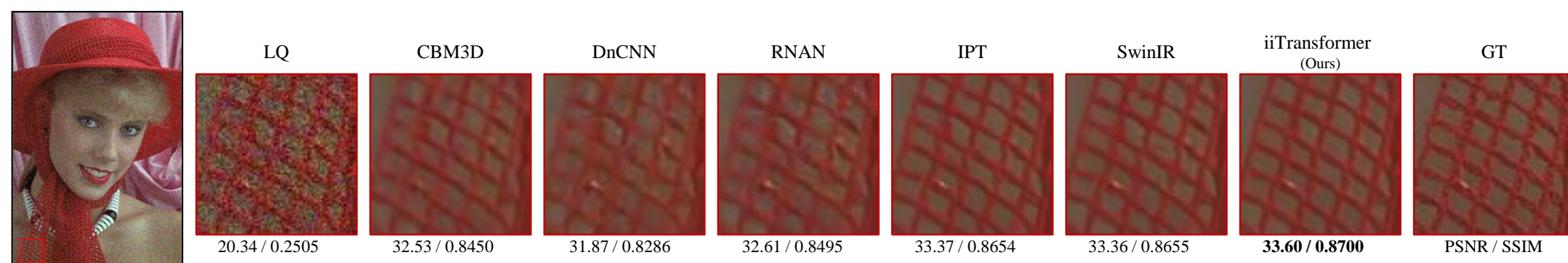
$$\bar{B}_{test}\left[d(P_i, P_j)\right] = \begin{cases} \bar{B}_{train}\left[d(P_i, P_j)\right] & \text{if } d(P_i, P_j) \leq d_{train}^{max} \\ -\infty & \text{otherwise} \end{cases}$$

$d(P_i, P_j)$ is the distance between patches $P_i$ and $P_j$,
$d_{train}^{max}$ is the distance between furthest patches in a training image, and
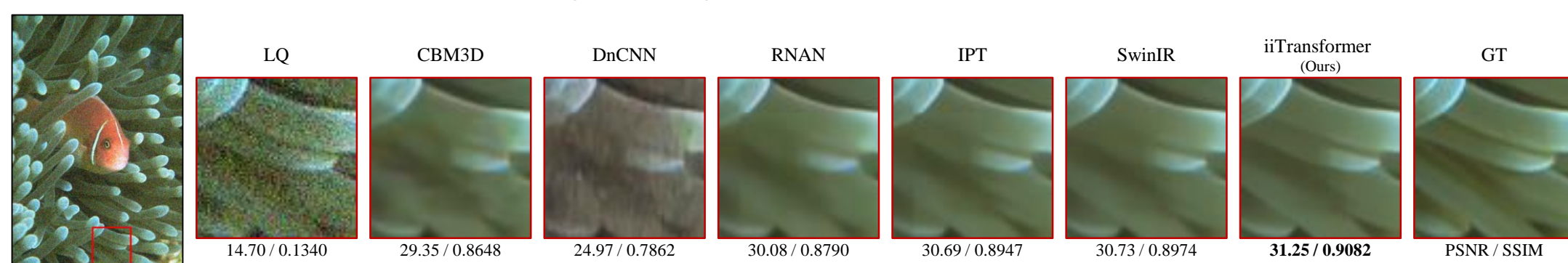$\bar{B}[k]$ is an element in $\bar{B}$ indexed at $k$



Masked eMSA

## Results and Conclusion



**AWGN Reduction**

(a) AWGN Image Denoising for $\sigma = 25$ on 'kodim04' from Kodak24

| LQ | CBM3D | DnCNN | RNAN | IPT | SwinIR | iiTransformer (Ours) | GT |
|---|---|---|---|---|---|---|---|
| 20.34 / 0.2505 | 32.53 / 0.8450 | 31.87 / 0.8286 | 32.61 / 0.8495 | 33.37 / 0.8654 | 33.36 / 0.8665 | **33.60 / 0.8700** | PSNR / SSIM |

(b) AWGN Image Denoising for $\sigma = 50$ on '210088' from BSDS68

| 14.70 / 0.1340 | 29.35 / 0.8648 | 24.97 / 0.7862 | 30.08 / 0.8790 | 30.69 / 0.8947 | 30.73 / 0.8974 | **31.25 / 0.9082** | PSNR / SSIM |

**JPEG CAR**

(a) JPEG CAR for q = 10 on 'paintedhouse' from LIVE1

| LQ | SA-DCT | DnCNN | RNAN | IPT | SwinIR | iiTransformer (Ours) | GT |
|---|---|---|---|---|---|---|---|
| 29.40 / 0.9082 | 26.12 / 0.7675 | 26.75 / 0.7846 | 26.44 / 0.7767 | 27.18 / 0.8018 | 27.41 / 0.8083 | **27.53 / 0.8103** | PSNR / SSIM |

(b) JPEG CAR for q = 20 on 'carnivaldolls' from LIVE1

| 31.07 / 0.8103 | 27.82 / 0.8453 | 31.54 / 0.9192 | 29.77 / 0.8838 | 31.74 / 0.9253 | 32.31 / 0.9327 | **32.41 / 0.9348** | PSNR / SSIM |

**SISR**

(a) SISR for s = 3 on 'YumeiroCooking' from Manga109

| LR | Bicubic | SRResNet | RCAN | SAN | IPT | SwinIR | iiTransformer (Ours) | GT |
|---|---|---|---|---|---|---|---|---|
| 26.54 / 0.8587 | 31.09 / 0.9565 | 33.14 / 0.9669 | 33.11 / 0.9666 | 32.21 / 0.9589 | 33.48 / 0.9684 | **33.99 / 0.9709** | PSNR / SSIM |

(b) SISR for s = 4 on 'barbara' from Set14

| 25.15 / 0.6851 | 25.96 / 0.7462 | 26.10 / 0.7537 | 26.06 / 0.7537 | 26.35 / 0.7518 | 26.05 / 0.7527 | **26.42 / 0.7604** | PSNR / SSIM |

### Ablation Study



LQ of 'kodim08' from Kodak24    Preferred Attention Mechanism

Regions Cropped from Each Output

| LQ | Intra | Inter | iiTransformer (Ours) | GT |
|---|---|---|---|---|
| 20.45 / 0.5606 | 30.59 / 0.9072 | 30.79 / 0.9113 | **30.89 / 0.9122** | PSNR / SSIM |

### Conclusion

- iiTransformer is a framework that combines local and non-local attention mechanisms to extract features at various sub-region levels of the image
- *Local* context is captured using the *intra* self-attention module and the *internal data repetition* is exploited using the *inter* self-attention module
- The patchwise relative position bias is *masked* to provide a boundary artifact-free solution for images of various resolutions
- State-of-the-art performance is achieved on various restoration tasks