# Sharing Key Semantics in Transformer Makes Efficient Image Restoration

**Bin Ren**[1,2,3*]   **Yawei Li**[4†]   **Jingyun Liang**[4]   **Rakesh Ranjan**[5]   **Mengyuan Liu**[6]
**Rita Cucchiara**[7]   **Luc Van Gool**[3]   **Ming-Hsuan Yang**[8]   **Nicu Sebe**[2]
[1]University of Pisa   [2]University of Trento   [3]INSAIT, Sofia University,
[4]ETH Zürich   [5]Meta Reality Labs   [6]Peking University
[7]University of Modena and Reggio Emilia   [8]University of California, Merced

## Abstract

Image Restoration (IR), a classic low-level vision task, has witnessed significant advancements through deep models that effectively model global information. Notably, the emergence of Vision Transformers (ViTs) has further propelled these advancements. When computing, the self-attention mechanism, a cornerstone of ViTs, tends to encompass all global cues, even those from semantically unrelated objects or regions. This inclusivity introduces computational inefficiencies, particularly noticeable with high input resolution, as it requires processing irrelevant information, thereby impeding efficiency. Additionally, for IR, it is commonly noted that small segments of a degraded image, particularly those closely aligned semantically, provide particularly relevant information to aid in the restoration process, as they contribute essential contextual cues crucial for accurate reconstruction. To address these challenges, we propose boosting IR's performance by sharing the key semantics via Transformer for IR (*i.e.*, SemanIR) in this paper. Specifically, SemanIR initially constructs a sparse yet comprehensive key-semantic dictionary within each transformer stage by establishing essential semantic connections for every degraded patch. Subsequently, this dictionary is shared across all subsequent transformer blocks within the same stage. This strategy optimizes attention calculation within each block by focusing exclusively on semantically related components stored in the key-semantic dictionary. As a result, attention calculation achieves linear computational complexity within each window. Extensive experiments across 6 IR tasks confirm the proposed SemanIR's state-of-the-art performance, quantitatively and qualitatively showcasing advancements. The visual results, code, and trained models are available at https://github.com/Amazingren/SemanIR.

## 1 Introduction

Image restoration (IR) stands as a fundamental task within low-level computer vision, aiming to enhance the quality of images affected by numerous factors, including noise, blur, low resolution, compression artifacts, mosaic patterns, adverse weather conditions, and other forms of distortion. This capability holds broad utility across various domains, facilitating information recovery in medical imaging, surveillance, and satellite imagery. Furthermore, it bolsters downstream vision tasks like object detection, recognition, and tracking [60, 70, 75]. Despite notable progress in recent years, prevalent IR methods encounter challenges in effectively addressing complex distortions or preserving/recovering crucial image details [48]. Achieving high-quality image recovery necessitates meticulous exploration of the rich information present in degraded counterparts.

---

*Work done during visiting at ETH Zürich and INSAIT Sofia University.
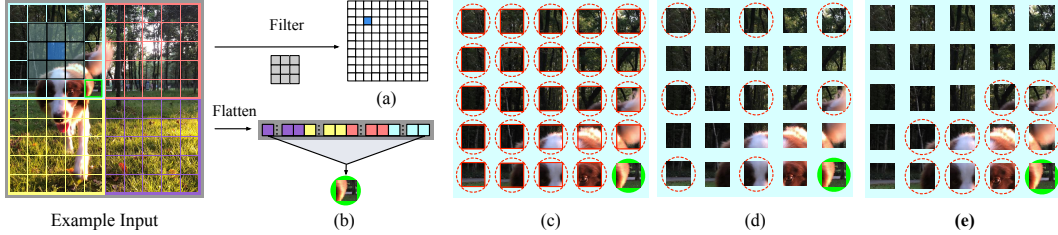†Project Leader & Corresponding Author. Email: li.yawei.ai@gmail.com

Figure 1: (a) The CNN filter captures information only within a local region. (b) The standard MLP/Transformer architectures take full input in a long-sequence manner. (c) The window-size multi-head self-attention (MSA) mechanism builds a full connection within each window. (d) Position-fixed sparse connection. (e) The proposed *Key-Semantic* connection.

In modern IR systems, representative networks for learning rich image information are typically constructed using three fundamental architectural paradigms. *i.e.*, convolutional neural networks (CNNs) [42, 95], Multilayer perceptrons (MLPs) [5, 83], and Vision Transformers (ViTs) [20, 85]. The input image is treated as a regular grid of pixels in the Euclidean space for CNNs (Fig. 1(a)) or a sequence of patches for MLPs and ViTs (Fig. 1(b)). However, degraded inputs usually contain irregular and complex objects. While these choices perform admirably in scenarios with regular or well-organized object boundaries, they have limitations when applied to degraded images with more flexible and complex geometrical contexts.

Additionally, CNNs struggle to model long-range dependencies due to their limited receptive field (Fig. 1(a)). In contrast, MLPs and ViTs can capture long-range relations effectively, although at the cost of losing inductive bias and incurring a heavy computational burden, *i.e.*, quadratic complexity increases with higher input resolution[20, 71, 83, 85]. To address these limitations, recent IR methods have explored strategies for complexity reduction. A common approach is to implement MSA within local image regions [50, 53]. For example, SwinIR [50] and GRL [48] employ full MSA or region-fixed anchored stripe MSA, but they still struggle with irregular object connections. Furthermore, prior research [112] highlights that smooth image content is more prevalent than complex image details, underscoring the need for distinct treatment based on semantic content.

In this paper, we introduce a novel approach, SemanIR, to address the limitations above. Specifically, within each transformer stage, we first construct a key-semantic dictionary, which stores only the top-k semantically related relations for each given degraded patch with the k-nearest neighbors (*i.e.*, KNN) algorithm. Then, the attention operation within each transformer layer occurs only among the top-k patches. This design brings two main advantages, *i.e.*, 1) Each degraded patch benefits from its semantically similar patches, typically containing comparable contextual or textural information, while excluding the side effects from other patches that contain entirely unrelated information. 2) Compared to the conventional window-wise attention, which built a dense connection between all the patches (Fig. 1(c)) that leads to highly computationally demanding, or a sparse but position-fixed manner (Fig. 1(d)) which introduces irrelevant semantics. Our key-semantic connection (Fig. 1(e)) leads to a sparse yet more meaningful attention operation, which allows our method to achieve the same receptive field as previous ViTs-based methods while maintaining lower computational costs. This is not like previous token merging or pruning methods [6, 69, 92] that may merge unrelated information or prune some semantically related information. In addition, to make the proposed method more efficient, instead of creating a key-semantic dictionary for each transformer layer, we create it just once at the beginning of each transformer stage and then share it with all the following transformer layers within the same stage. This not only largely reduced the computation burden but also made our methods different from other token pruning and merging methods [90, 107, 108], which include dynamic patch skipping/selection within each attention layer or involve an additional offset generation network. Meanwhile, merging or pruning tokens will lead to a loss of information in corresponding patches, which is not preferred in image restoration [86]. In addition, such a sharing strategy allows each degraded patch to be continuously optimized by its semantically related patches within each stage.

It is also worth noting that the implementation of the attention layer of our SemanIR is achieved in three interesting manners (*i.e.*, Triton [19][3], torch-mask, and torch-gather), which are discussed in our ablation studies. Overall, our method's suitability for image restoration comes from the

---

[3]Open-source GPU programming tool https://openai.com/research/triton.

2

utilization of semantic information, the preservation of the details, and the effective KNN strategy. The contributions of this work are:

1. For each degraded input patch, we propose to construct a key-semantic dictionary that stores its most semantically relevant $k$ patches in a sparse yet representative manner. This strategy excludes the side effects of a given degraded patch from semantically unrelated parts.
2. Based on the constructed key-semantic dictionary, we propose to share the key semantic information across all the attention layers within each transformer stage, which not only makes each degraded patch well-optimized but also largely reduces the computational complexity compared to conventional attention operations.
3. Extensive experimental results show that the proposed SemanIR achieves state-of-the-art performance on 6 IR tasks, *i.e.*, deblurring, JPEG compression artifact removal (JPEG CAR), denoising, IR in adverse weather conditions (AWC), demosaicking, and classic image super-resolution (SR).

## 2   Related Work

**Image Restoration (IR),** as a long-standing ill-posed inverse problem, is designed to reconstruct the high-quality image from the corresponding degraded counterpart with numerous applications [3, 49, 73]. Initially, IR was addressed through model-based solutions, involving the search for solutions to specific formulations. However, learning-based approaches have gained much attention with the significant advancements in deep neural networks. Numerous approaches have been developed, including regression-based [41, 48, 50, 51, 105] and generative model-based pipelines [25, 56, 88, 93]. In this paper, we propose a regression-based method for image restoration.

**Non-Local Priors Modeling in IR.** Tradition model-based IR methods reconstruct the image by regularizing the results (*e.g.*, Tikhonov regularization [27]) with formulaic prior knowledge of natural image distribution. However, it is challenging for these methods to recover realistic detailed results with hand-designed priors. Besides, some other classic method finds that self-similarity is an effective prior, which leads to an impressive performance [7, 17]. Apart from traditional methods, the non-local prior has also been utilized in modern deep learning networks [48, 87, 110], typically captured by the self-attention mechanism. More recently, the overwhelming success of transformers [85] in the natural language processing domain [35] and the classic vision community [8, 10, 20, 82, 91] has led to the development of numerous ViT-based IR methods. These methods aim to enhance the learning ability for modeling non-local priors [14, 15, 48, 50, 94, 99] and consequently archives better performance. Meanwhile, this raises a question: are all non-local priors essential for IR?

**Key-Semantic Non-local Prior Exploration for IR.** To answer the question, we found many methods demonstrating the effectiveness of modeling key semantics within ViTs. For example, KiT [43] proposed increasing the non-local connectivity between patches at different positions through KNN matching. This approach aims to better capture the non-local relations between the base patch and other patches in each attention calculation. However, it results in significant additional computational costs due to the KNN matching. DRSformer [13] proposed a top-k selection strategy that chooses the most relevant tokens to model the non-local priors for draining after each self-attention calculation without reducing the computation complexity, since after each attention calculation, the DRSFormer utilized (mask, top-k, scatter) operations at each transformer layer. Consequently, this inevitably increases the computation cost. Similar conclusions can be also drawn from the graph perspective solutions [28, 32, 61, 74] for various IR tasks, like facial expression restoration [54], image denoising [78], and artifact reduction [61]. [32] construct the graph with transformer-based architecture where each patch is connected to all other patches. All these methods suggest that if the semantically related information can be addressed, the degraded image can be restored with better performance. However, the efficiency issue, which is extremely unignorable, remains untouched within the aforementioned methods. It is particularly crucial for ViTs-based image restoration methods, which often need to address high-resolution degraded input images. LaViT [108] reduces computational costs by storing attention scores from a few initial layers and reusing them in subsequent layers. However, this approach does not change the computation cost of attention itself; it merely reuses previously computed scores. In this paper, we propose sharing key semantics within each transformer stage, demonstrating its efficiency and effectiveness through experimental and theoretical analysis. Our method, SemanIR, reduces computation in both training and inference

**(a)** The architecture of the proposed **SemanIR** for image super-resolution

**(b)** The proposed Key-Semantic Transformer Layer     **(c)** Key-Semantic Dictionary     **(d)** Key-Semantic Attention

Feature Extractor    Image Reconstructor    **D** Key-Seman Dictionary Construction    ⊕ Element-wise addition    ⊗ Matrix Multiplication
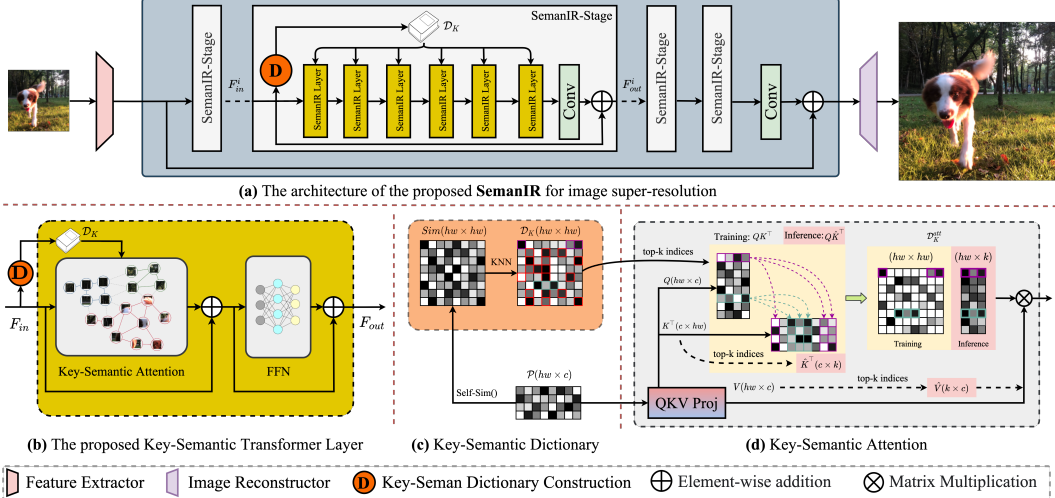
Figure 2: The proposed SemanIR mainly consists of a convolutional feature extractor, the main body of SemanIR for representation learning, and an image reconstructor. The main body in columnar shape shown here is for image SR, while the U-shaped structure (shown in *Appx.* A.2) is used for other IR tasks. (b) The transformer layer of our SemanIR. The toy example of $k$=3 for (c) the Key-semantic dictionary construction and (d) the attention of each Layer.

by using a semantic dictionary to filter out irrelevant patches during training and optimizing attention operations with Triton kernels during inference.

## 3 Methodology

To comprehensively study the effectiveness of the proposed method that is architecture-agnostic for various IR tasks, we adopted two of the most commonly used architectures *i.e.* the columnar architecture (shown in Fig. 2 (a)) for image SR and the U-shaped architecture (shown in the *Appendix*, *i.e.*, *Appx.* A.2) for other IR tasks. In the following, we first show how to construct the key-semantic dictionary in Sec. 3.1. Based on the key-semantic dictionary, then we explain why sharing it works for IR, and we introduce the basic unit, the key-semantic transformer layer in Sec. 3.2. Finally, two interesting discussions (Sec.3.3) are introduced regarding the implementation style of the Key-Graph attention and two top-k settings during the training. The efficiency analysis is provided in *Appx.* A.3.

### 3.1 Key-Semantic Dictionary Construction

Consider the input feature $F_{in} \in \mathbb{R}^{H \times W \times C}$, where $H$, $W$, and $C$ denote the height, the width, and the channel. ViTs are good at modeling global dependencies for $F_{in}$. This is achieved by the MSA, the core of ViTs, by connecting all other patches to a certain patch. Specifically, $F_{in}$ is first split into $N$ non-overlapping patches, forming its patch representation $\mathcal{P} = \{p_i | p_i \in \mathbb{R}^{hw \times c}, i = 1, 2, 3, ..., N\}$, where $h$, $w$, and $c$ are the height, the width, and the channel of each patch. To achieve such global connectivity $\mathcal{D}$, $\mathcal{P}$ is linearly projected into Query ($Q$), Key ($K$), and Value ($V$) matrices, which are denoted as $Q = \mathcal{P}\mathbf{W}_{qry}$, $K = \mathcal{P}\mathbf{W}_{key}$, and $V = \mathcal{P}\mathbf{W}_{val}$. $\mathbf{W}_{qry/key/val}$ represents the learnable projection weights. Then $\mathcal{D}$ is performed by a softmax function as follows:

$$\mathcal{D}_{ij} = \frac{\exp(Q_i K_j^\top)}{\sum_{k=1...j} \exp(Q_i K_k^\top / \sqrt{d})}, i = 1, 2, 3, ..., N, \tag{1}$$

where $d$ is the dimension of $Q$ and $K$. Then each patch is aggregated via $\sum_i \mathcal{D}_{ij} V_i$. However, $\mathcal{D}_{ij}$ functions as a full semantic dictionary, where each patch is connected to all other patches regardless of their semantic relatedness. *e.g.*, given a sub-graph with a green dog patch shown in Fig. 1(c), the tree-related patches are also considered. Since such an operation occurs at each attention calculation step in ViTs, it inevitably increases the computational cost, especially for large-scale inputs. In addition, for IR, a degraded patch usually benefits from its most semantically related patches, as they share similar texture and geometric information. This naturally raises the question: *Can we build a key-semantic dictionary, $\mathcal{D}_K$, where each patch is connected only to its most related patches?*

4

---

**Algorithm 1** Key-Semantic Transformer Stage (*i.e.*, SemanIR Stage)

---

**Input:** input feature $F_{in}$, numbers of SemanIR layer $N_{layer}$, KNN value $k$, the patched feature $\mathcal{P}$
**Output:** aggregated feature $F_{out}$
1: $\mathcal{D}_K \leftarrow$ KeySemanDictionary_Constructor$(\mathcal{P}, k)$ // Sec. 3.1
2: **for** $i = 1$ to $N_{layer}$ **do**
3:     $Q, K, V \leftarrow$ Linear_Proj$(\mathcal{P})$
4:     $\hat{\mathcal{P}} \leftarrow$ SemanIR_Att$(Q, K, V, \mathcal{D}_K)$ // Sec. 3.2
5:     $\mathcal{P} \leftarrow \hat{\mathcal{P}} + $ FFN$(\hat{\mathcal{P}})$ // Sec. 3.2
6: **end for**
7: $F_{out} \leftarrow F_{in} + $ Conv(UnPatch$(\mathcal{P})$)
8: **return** $F_{out}$

---

To mitigate this problem, given $\mathcal{P}$, we first construct a fully connected dictionary $\mathcal{D}$ by calculating its self-similarity Sim() via a naive dot product operation as $\mathcal{D}(i, j) = \text{Sim}(i, j) = p_i \cdot p_j^\top$, which describes the correlation among all the patches, with higher values indicating stronger correlations. To reduce the side influence of patches with low correlation (*e.g.*, the tree-related patches at the upper left part in Fig. 1 (c)) for the green background dog destination patch, we keep only $k$ highly related patches and exclude the remaining. This is achieved by a KNN algorithm from $\mathcal{D}$ as follows:

$$\mathcal{D}_K(i, j) = \begin{cases} \mathcal{D}(i, j), & \mathcal{D}(i, j) \geq \text{Sim}(i, )_k \text{ and } i \neq j \\ 0, & \text{otherwise,} \end{cases} \tag{2}$$

where $\text{Sim}(i, )_k$ denotes the $k_{th}$ largest connectivity value of patch $p_i$. As a result, $\mathcal{D}_K$ contains only the patches with high correlation (*e.g.*, dog-related patches in Fig. 1(e)) for the destination patch (*e.g.*, the green dog patch). We formalize the key-semantic dictionary construction as KeySemanDictionary_Constructor() in Alg. 1. Although such a dictionary allows the subsequent attention operation to focus on the most semantically related patches, constructing it before each attention operation significantly increases computational costs. Meanwhile, we observed that IR architectures typically use transformers stage-by-stage (See Fig. 2(a)). This means that in each stage, several transformer layers are directly connected sequentially and operate at the same semantic level. Inspired by this, we propose to *share the same key-semantic dictionary* for all the transformer layers.

### 3.2 Sharing Key Semantics Cross Transformer Layers

The structure of each transformer layer is shown in Fig. 2(b), which consists of a key-semantic attention block followed by a feed-forward network (FFN). Specifically, given an input $F_{in}$, we form each transformer layer as $z = \text{FFN}(f_\theta(F_{in}))$, where $f$ means the transformer layer, $z$ denotes the output, and $\theta$ is the trainable parameters. Previous methods tried to reduce the computation cost mainly by applying some techniques (*i.e.*, $\mathcal{T}$) like token merging or pruning after each attention calculation or the entire transformer layer, and it can be formalized as $z = \mathcal{T}(\text{FFN}(f_\theta(F_{in})))$ or $z = (\text{FFN}(\mathcal{T}f_\theta(F_{in})))$. However, the main computation cost from MSA is still untouched.

Owing to the permutation-invariant property (*i.e.*, $f_\theta(\mathcal{T}x) = \mathcal{T}f_\theta(x)$, here $\mathcal{T} \in \mathbb{R}^{N \times N}$ means any token level permutation matrix) inherent in both the MSA and the FFN [44, 85], the transformer layer consistently produces identical representations for patches that share the same attributes, regardless of their positions or the surrounding structures [9]. In other words, patches at the same location are consistently connected to other patches possessing the same attributes as they traverse through the various layers within the same stage. It enables $\mathcal{D}_K$ to serve as a reference permutation for each attention in the subsequent transformer layers, facilitating efficient yet highly semantics-related attention operations. This distinguishes our method from previous token merging/pruning [90, 107] or sparse attention-based methods (Fig.1(d)) that only activate patches in a grid-fixed manner[99].

The workflow is intuitively illustrated in Fig. 2 (c) and (d). Initially, the patch $\mathcal{P}$ is linear projected via Linear_Proj() (The 3rd step in Alg. 1) into $Q$, $K$, and $V$. For each patch $p_i$ in $Q$, instead of calculating the self-attention with all $hw$ patches in $K\&V$, only $k$ essential patches are selected via the semantic lookup via the indices from $\mathcal{D}_K$ in them, forming the $\hat{K}\&\hat{V}$. Then the attention matrix $\mathcal{D}_K^{att}$ is obtained by: $\mathcal{D}_K^{att} = \text{Softmax}_K(Q\hat{K}^\top/\sqrt{d})$, which captures the pair-wise relation between each destination patch $p_i$ in $Q$ with only $k$ patches in $K\&V$ that are semantically highly related to $p_i$. For other unselected patches in $K\&V$, we aim to maintain their position in their corresponding

Table 1: GPU memory footprint of different implementations (*i.e.*, Triton, Torch-Gather, and Torch-Mask) of our key-graph attention block. $N$ is the number of tokens and $k$ is the number of nearest neighbors. OOM denotes "out of memory".

| $N$ | Triton | Torch-Gather | Torch-Mask |
|---|---|---|---|
| 512 | 0.27 GB | 0.66 GB | 0.36 GB |
| 1024 | 0.33 GB | 1.10 GB | 0.67 GB |
| 2048 | 0.68 GB | 2.08 GB | 1.91 GB |
| 4096 | 2.61 GB | 4.41 GB | 6.83 GB |
| 8192 | 10.21 GB | 10.57 GB | 26.42 GB |

| $k$ | Triton | Torch-Gather | Torch-Mask |
|---|---|---|---|
| 32 | 5.51 GB | 15.00 GB | 13.68 GB |
| 64 | 5.82 GB | 27.56 GB | 13.93 GB |
| 128 | 6.45 GB | OOM | 14.43 GB |
| 256 | 7.70 GB | OOM | 15.43 GB |
| 512 | 10.20 GB | OOM | 17.43 GB |

Table 2: ***Single-image motion deblurring*** results. GoPro [62] dataset is used for training.

| Method | GoPro | | HIDE | | Average | |
|---|---|---|---|---|---|---|
| | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| DeblurGAN [39] | 28.70 | 0.858 | 24.51 | 0.871 | 26.61 | 0.865 |
| Nah *et al.* [62] | 29.08 | 0.914 | 25.73 | 0.874 | 27.41 | 0.894 |
| DeblurGAN-v2 [40] | 29.55 | 0.934 | 26.61 | 0.875 | 28.08 | 0.905 |
| SRN [81] | 30.26 | 0.934 | 28.36 | 0.915 | 29.31 | 0.925 |
| Gao *et al.* [24] | 30.90 | 0.935 | 29.11 | 0.913 | 30.01 | 0.924 |
| DBGAN [104] | 31.10 | 0.942 | 28.94 | 0.915 | 30.02 | 0.929 |
| MT-RNN [64] | 31.15 | 0.945 | 29.15 | 0.918 | 30.15 | 0.932 |
| DMPHN [97] | 31.20 | 0.940 | 29.09 | 0.924 | 30.15 | 0.932 |
| Suin *et al.* [80] | 31.85 | 0.948 | 29.98 | 0.930 | 30.92 | 0.939 |
| CODE [111] | 31.94 | - | 29.67 | - | 30.81 | - |
| SPAIR [65] | 32.06 | 0.953 | 30.29 | 0.931 | 31.18 | 0.942 |
| MIMO-UNet+ [16] | 32.45 | 0.957 | 29.99 | 0.930 | 31.22 | 0.944 |
| IPT [10] | 32.52 | - | - | - | - | - |
| MPRNet [95] | 32.66 | 0.959 | 30.96 | 0.939 | 31.81 | 0.949 |
| KiT [43] | 32.70 | 0.959 | 30.98 | 0.942 | 31.84 | 0.951 |
| NAFNet [11] | 32.85 | 0.960 | - | - | - | - |
| Restormer [94] | 32.92 | 0.961 | 31.22 | 0.942 | 32.07 | 0.952 |
| Ren *et al.* [72] | 33.20 | 0.963 | 30.96 | 0.938 | 32.08 | 0.951 |
| SemanIR (ours) | 33.44 | 0.964 | 31.05 | 0.941 | 32.25 | 0.953 |

places without any computation. Based on $\mathcal{D}_K^{att}$, the attention outputs the updated feature $\hat{\mathcal{P}}$ via: $\hat{\mathcal{P}} = \mathcal{D}_K^{att}\hat{V}$. We formulate these two procedures as $\text{SemanIR\_Att}()$ in the 4th step of Alg. 1. This differs from the conventional MSA, which calculates the relation of each patch in $Q$ and all patches in $K\&V$. Finally, with FFN, the output of each transformer layer is achieved via the 5th step of Alg. 1.

Conversely, our design offers two advantages. Firstly, the computational cost can be significantly reduced within each attention window (detailed analysis can be found in the *Appx.* A.3), enhancing efficiency. Additionally, sharing the key semantics across transformer layers within each stage acts as a loop that continuously optimizes a degraded patch with its most semantically related patches, ensuring the performance of the proposed method (supported by our experimental results in Sec. 4).

### 3.3 Discussion

**Fixed top-k *vs.* Random top-k Training Strategies.** In the fixed top-k approach, $k$ remains constant at 512 during training. In contrast, in the random top-k method, $k$ is randomly selected from the set $[64, 128, 192, 256, 384, 512]$. It is important to note that even in the random top-k setting, a fixed k value is maintained for all patches/pixels in each iteration. During inference, the random top-k strategy offers more flexibility and requires training only a single model, making it more user-friendly and less resource-intensive.

**Implementation of the Attention of SemanIR.** To achieve the attention operation of the proposed SemanIR, we explored three different manners for the implementation, *i.e.*, (i) *Triton*, (ii) *Torch-Gather*, and (iii) *Torch-Mask*. Specifically, (i) is based on FlashAttention [19], and a customized GPU kernel is written for the operators proposed in this paper. Parallel GPU kernels are called for the nodes during run time. (ii) means that we use the 'torch.gather()' function in PyTorch to choose the corresponding $Q_{gather}$ and $K_{gather}$ based on $\mathcal{D}_K$, then the attention operation is conducted between $Q_{gather}$ and $K_{gather}$. (iii) denotes that we keep only the value of selected patches of $\mathcal{D}_K$ and omitting other patches with low correlation via assigning those values to $-\infty$ guided by $\mathcal{D}_K$. Discussions of the pros and cons regarding these manners are provided in Sec. 4.1.

## 4 Experiments

In this section, we first analyze three important ablation studies of our SemanIR, followed by extensive experiments on **6** IR tasks, *i.e.*, deblurring, JPEG CAR, image denoising, IR in AWC, image demosaicking, and image SR. More details about the architecture design, training protocols, the training/testing dataset, and full quantitative/additional qualitative results are shown in *Appx.* A to E. The best and the 2nd-best results are reported in red and blue, respectively. Note that † denotes a single model that is trained to handle multiple degradation levels *i.e.*, noise levels, and quality factors.
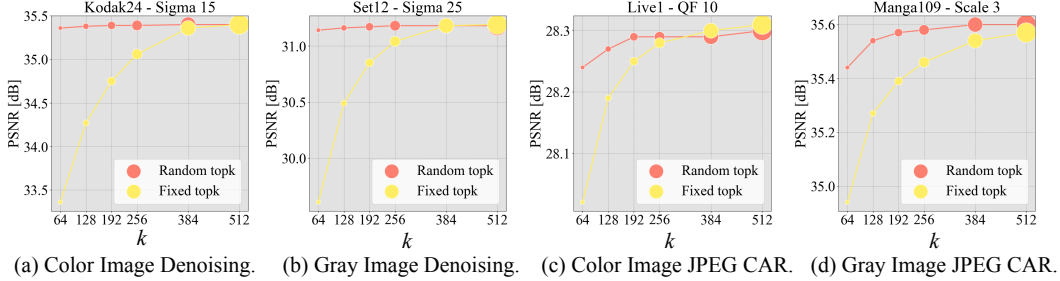
(a) Color Image Denoising.  (b) Gray Image Denoising.  (c) Color Image JPEG CAR.  (d) Gray Image JPEG CAR.

Figure 3: The impact of $k$ with different inference $k$ value. Circle size represents FLOPs.



(a) query region    (b) without top-k    (c) with top-k=8    (d) with top-k = 16    (e) with top-k = 64    (f) with top-k = 256
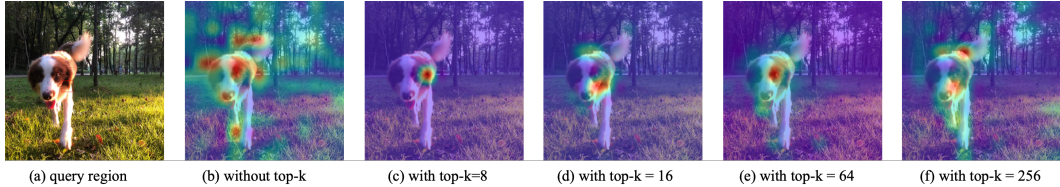
Figure 4: The impact of $k$ with different inference $k$ value.

Table 3: The **efficiency comparisons** results on Urban100 dataset.

| Task | Method | Architecture | Params [M]↓ | FLOPs [G]↓ | Runtime [ms]↓ | PSNR↑ |
|---|---|---|---|---|---|---|
| ×4 SR | SwinIR [50] | Columnar | 11.90 | 215.32 | 152.24 | 27.45 |
| | CAT [15] | Columnar | 16.60 | 387.86 | 357.97 | 27.89 |
| | HAT [14] | Columnar | 20.77 | 416.90 | 368.61 | 28.37 |
| | SemanIR-S (Ours) | Columnar | 12.02 | 290.20 | 211.94 | 28.34 |
| Denoising ($\sigma = 50$) (The same architecture for other IR task) | SwinIR [50] | Columnar | 11.75 | 752.06 | 1772.84 | 27.98 |
| | Restormer [94] | U-shape | 26.10 | 154.88 | 210.44 | 28.29 |
| | GRL [48] | Columnar | 19.81 | 1361.77 | 3944.17 | 28.59 |
| | SemanIR (Ours) | U-shape | 25.85 | 135.26 | 240.05 | 28.63 |

## 4.1 Ablation Study

**The impact of the implementation of SemanIR Attention** is assessed in terms of (i) *Triton*, (ii) *Torch-Gather*, and (iii) *Torch-Mask* under different numbers of N (various from 512 to 8192) and K (various from 32 to 512). The results of the GPU memory footprint are shown in Tab. 1, which indicate that *Torch-Gather* brings no redundant computation while requiring a large memory footprint. Though *Torch-Mask* brings the GPU memory increase, the increment is affordable compared to *Torch-Gather* and also easy to implement. *Triton* largely saves the GPU memory while at the cost of slow inference and difficult implementation for the back-propagation process. To optimize the efficiency of our SemanIR, we recommend employing *Torch-Mask* during training and *Triton* during inference, striking a balance between the efficiency and the GPU memory requirement.

**The Impact of the $k$ in Key-Semantic Dictionary Construction.** Three interesting phenomena are observed from the results shown in Fig. 3 regarding the two top-k training strategies (Sec. 3.3). (1) The PSNR can largely increase with the increase of $k$ in a fixed manner. (2) When $k$ reaches a certain number (*i.e.*, 384), the performance improvements become marginal, supporting our statement that only the most semantically related patches contribute significantly to the restoration. (3) The randomly sampled strategy has a very stable and better performance compared to the fixed top-k manner especially when the inference $k$ is fixed to a small number (*i.e.*, 64, 128, 256). We conclude that a random sampled strategy is more general and stable. It can also make the inference process more flexible regarding different computation resources. Meanwhile, we set a query region in the input and provided a detailed comparison from the attention-based activation map together with the input query region in Fig. 4. Fig. 4(a) shows the query region input. Fig. 4(b) displays the activation map generated using standard attention mechanisms. Fig. 4(c-f) illustrate activation maps using our key-semantic dictionary with different $k$ values ([8, 16, 64, 256]) during inference. The comparisons indicate that increasing $k$ allows for connections to more semantically related regions. However, when $k$ is set too high (*e.g.*, $k = 256$ as shown in Fig. 4(f)), the activation map may include some semantically unrelated regions. This aligns with the findings and the results depicted in Fig 3, where increasing the $k$ beyond a certain point (e.g., from 396 to 512) does not further improve PSNR. More ablation results can be found in our *Appx.* D about the effect of the noise level and quality factor for denoising and JPEG CAR.

Table 4: ***Color image JPEG compression artifact removal*** results.

| Set | QF | JPEG | | QGAC [21] † | | FBCNN [33] † | | DRUNet [100] † | | SemanIR (Ours) † | | SwinIR [50] | | GRL-S [48] | | SemanIR (Ours) | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ | PSNR↑ | SSIM↑ |
| LIVE1 | 10 | 25.69 | 0.7430 | 27.62 | 0.8040 | 27.77 | 0.8030 | 27.47 | 0.8045 | 28.19 | 0.8146 | 28.06 | 0.8129 | 28.13 | 0.8139 | 28.31 | 0.8176 |
| | 20 | 28.06 | 0.8260 | 29.88 | 0.8680 | 30.11 | 0.8680 | 30.29 | 0.8743 | 30.53 | 0.8781 | 30.44 | 0.8768 | 30.49 | 0.8776 | 30.61 | 0.8792 |
| | 30 | 29.37 | 0.8610 | 31.17 | 0.8960 | 31.43 | 0.8970 | 31.64 | 0.9020 | 31.89 | 0.9051 | 31.81 | 0.9040 | 31.85 | 0.9045 | 31.94 | 0.9058 |
| | 40 | 30.28 | 0.8820 | 32.05 | 0.9120 | 32.34 | 0.9130 | 32.56 | 0.9174 | 32.81 | 0.9201 | 32.75 | 0.9193 | 32.79 | 0.9195 | 32.85 | 0.9204 |
| BSDS500 | 10 | 25.84 | 0.7410 | 27.74 | 0.8020 | 27.85 | 0.7990 | 27.62 | 0.8001 | 28.25 | 0.8076 | 28.22 | 0.8075 | 28.26 | 0.8083 | 28.37 | 0.8102 |
| | 20 | 28.21 | 0.8270 | 30.01 | 0.8690 | 30.14 | 0.8670 | 30.39 | 0.8711 | 30.55 | 0.8738 | 30.54 | 0.8739 | 30.57 | 0.8746 | 30.63 | 0.8750 |
| | 30 | 29.57 | 0.8650 | 31.33 | 0.8980 | 31.45 | 0.8970 | 31.73 | 0.9003 | 31.90 | 0.9026 | 31.90 | 0.9025 | 31.92 | 0.9030 | 31.96 | 0.9035 |
| | 40 | 30.52 | 0.8870 | 32.25 | 0.9150 | 32.36 | 0.9130 | 32.66 | 0.9168 | 32.84 | 0.9190 | 32.84 | 0.9189 | 32.86 | 0.9192 | 32.88 | 0.9193 |



(a) Color Image Denoising, One-for-multi Sigma
(b) Gray Image Denoising, One-for-multi Sigma
(c) Color Image JPEG CAR, One-for-multi QF
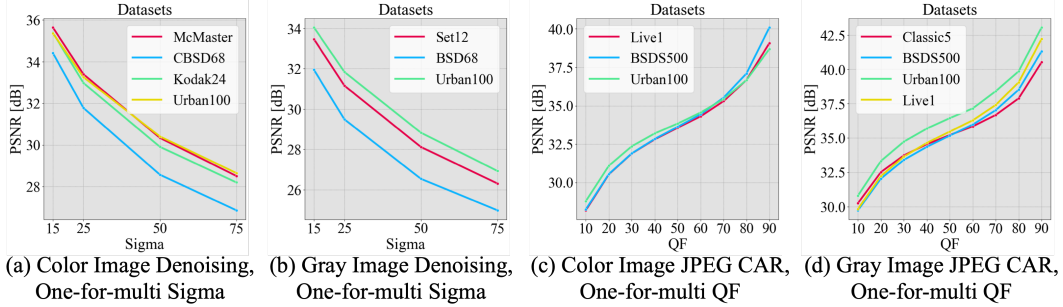(d) Gray Image JPEG CAR, One-for-multi QF

Figure 5: One model is trained to handle multi-degradations for denoising (a-b) and JPEG CAR (c-d).

Table 5: ***Gray image denoising*** PSNR.

| Window Size | Set12 | BSD68 | Urban100 |
|---|---|---|---|
| 8 | 31.01 | 29.49 | 31.33 |
| 16 | 31.06 | 29.51 | 31.45 |
| 32 | 31.17 | 29.50 | 31.88 |

Table 6: ***Color image denoising*** PSNR.

| Window Size | Mcmaster | CBSD68 | Kodak24 | Urban100 |
|---|---|---|---|---|
| 8 | 33.20 | 31.72 | 32.84 | 32.89 |
| 16 | 33.24 | 31.74 | 32.86 | 32.95 |
| 32 | 33.38 | 31.75 | 32.97 | 33.27 |

**Efficiency Analysis.** We compare our SemanIR method with four recent promising approaches on both 4x super-resolution (SR) and grayscale denoising tasks using the Urban100 dataset. Key metrics, including trainable parameters, FLOPs, runtime, and PSNR, are reported in Tab. 3. The results show that HAT and SemanIR achieve top-tier PSNR performance, reaching 28.37 dB and 28.34 dB, respectively, while SemanIR-S is much faster and uses 41.7% fewer parameters than HAT, making it more computationally efficient. SwinIR runs slightly faster than SemanIR-S but at the cost of a significant 0.89 dB loss in PSNR. In comparison, SemanIR-S offers better accuracy and speed than both CAT and HAT, solidifying its effectiveness. Furthermore, the comparison between SwinIR and SemanIR for denoising reveals that, although SwinIR has fewer trainable parameters, its FLOPs are substantially higher, indicating that SemanIR is more efficient in terms of computation. These results highlight that SemanIR-S strikes an optimal balance between performance and efficiency, making it highly competitive for both SR and denoising tasks.

**The Impact of One Model is Trained to Handle Multiple Degradation Levels.** The experiments were conducted for both denoising and JPEG CAR with both color and grayscale images. For denoising, $\sigma$ is set to $[15, 25, 50, 75]$. For JEPG CAR, QF is set to $[10, 20, 30, 40, 50, 60, 70, 80, 90]$. The results in Fig. 5 indicate that the PSNR for both tasks across all the datasets, under both color and grayscale settings, decreases when the degraded level increases. It means that maintaining a decent generalization ability with one model to handle multiple degradation levels is not easy. Meanwhile, the proposed SemanIR can still outperform other methods on various tasks (See Tab. 7 and Tab. 4), which means that capturing the key semantics is also essential for model's generalization ability.

**The Impact of the window size.** The windows indeed contain mixed information from different semantic parts. Yet, it is precisely this semantic distinction that motivates us to develop a selection mechanism for semantic information using KNN. We have conducted ablation studies of the window size (*i.e.*, on both gray and color image denoising with $\sigma = 25$). The results are summarized in Tab. 5 and Tab. 6. With the increase of the window size, the semantic relevant information for each token is increased, thus leading to a PSNR gain for different IR tasks. For gray image denoising, larger window sizes lead to improved PSNR performance, consistently yielding higher PSNR values across different datasets. Similarly, in color image denoising, larger window sizes also result in better outcomes. For example, in the McMaster dataset, the PSNR increases from 33.20 dB with a

Table 7: ***Color and grayscale image denoising*** PSNR results.

| Method | # P | Color | | | | | | | | | Grayscale | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | CBSD68 | | | McMaster | | | Urban100 | | | Set12 | | | BSD68 | | | Urban100 | | |
| | | $\sigma$=15 | $\sigma$=25 | $\sigma$=50 | $\sigma$=15 | $\sigma$=25 | $\sigma$=50 | $\sigma$=15 | $\sigma$=25 | $\sigma$=50 | $\sigma$=15 | $\sigma$=25 | $\sigma$=50 | $\sigma$=15 | $\sigma$=25 | $\sigma$=50 | $\sigma$=15 | $\sigma$=25 | $\sigma$=50 |
| †DnCNN [36] | 0.56 | 33.90 | 31.24 | 27.95 | 33.45 | 31.52 | 28.62 | 32.98 | 30.81 | 27.59 | 32.67 | 30.35 | 27.18 | 31.62 | 29.16 | 26.23 | 32.28 | 29.80 | 26.35 |
| †FFDNet [102] | 0.49 | 33.87 | 31.21 | 27.96 | 34.66 | 32.35 | 29.18 | 33.83 | 31.40 | 28.05 | 32.75 | 30.43 | 27.32 | 31.63 | 29.19 | 26.29 | 32.40 | 29.90 | 26.50 |
| †DRUNet [100] | 32.64 | 34.30 | 31.69 | 28.51 | 35.40 | 33.14 | 30.08 | 34.81 | 32.60 | 29.61 | 33.25 | 30.94 | 27.90 | 31.91 | 29.48 | 26.59 | 33.44 | 31.11 | 27.96 |
| †Restormer [94] | 26.13 | 34.39 | 31.78 | 28.59 | 35.55 | 33.31 | 30.29 | 35.06 | 32.91 | 30.02 | 33.35 | 31.04 | 28.01 | 31.95 | 29.51 | 26.62 | 33.67 | 31.39 | 28.33 |
| †SemanIR (Ours) | 25.82 | 34.42 | 31.78 | 28.57 | 35.65 | 33.40 | 30.34 | 35.37 | 33.26 | 30.41 | 33.47 | 31.16 | 28.12 | 31.95 | 29.49 | 26.54 | 34.05 | 31.84 | 28.83 |
| DnCNN [36] | 0.56 | 33.90 | 31.24 | 27.95 | 33.45 | 31.52 | 28.62 | 32.98 | 30.81 | 27.59 | 32.86 | 30.44 | 27.18 | 31.73 | 29.23 | 26.23 | 32.64 | 29.95 | 26.26 |
| EDT-B [47] | 11.48 | 34.39 | 31.76 | 28.56 | 35.61 | 33.34 | 30.25 | 35.22 | 33.07 | 30.16 | - | - | - | - | - | - | - | - | - |
| DRUNet [100] | 32.64 | 34.30 | 31.69 | 28.51 | 35.40 | 33.14 | 30.08 | 34.81 | 32.60 | 29.61 | 33.25 | 30.94 | 27.90 | 31.91 | 29.48 | 26.59 | 33.44 | 31.11 | 27.96 |
| SwinIR [50] | 11.75 | 34.42 | 31.78 | 28.56 | 35.61 | 33.20 | 30.22 | 35.13 | 32.90 | 29.82 | 33.36 | 31.01 | 27.91 | 31.97 | 29.50 | 26.58 | 33.70 | 31.30 | 27.98 |
| Restormer [94] | 26.13 | 34.40 | 31.79 | 28.60 | 35.61 | 33.34 | 30.30 | 35.13 | 32.96 | 30.02 | 33.42 | 31.08 | 28.00 | 31.96 | 29.52 | 26.62 | 33.79 | 31.46 | 28.29 |
| Xformer [98] | 25.23 | 34.43 | 31.82 | 28.63 | 35.68 | 33.44 | 30.38 | 35.29 | 33.21 | 30.36 | 33.46 | 31.16 | 28.10 | 31.98 | 29.55 | 26.65 | 33.98 | 31.78 | 28.71 |
| SemanIR (Ours) | 25.82 | 34.43 | 31.79 | 28.60 | 35.65 | 33.43 | 30.38 | 35.38 | 33.29 | 30.51 | 33.48 | 31.18 | 28.14 | 31.97 | 29.52 | 26.53 | 34.09 | 31.87 | 28.86 |

Table 8: ***IR in adverse weather conditions***.

| Type | Test1 (rain+fog) | | SnowTest100k-L | | RainDrop | |
|---|---|---|---|---|---|---|
| | Method | PSNR↑ | Method | PSNR↑ | Method | PSNR↑ |
| Task Specific | pix2pix [31] | 19.09 | DesnowNet [52] | 27.17 | AttGAN [66] | 30.55 |
| | HRGAN [45] | 21.56 | JSTASR [12] | 25.32 | Quan [68] | 31.44 |
| | SwinIR [50] | 23.23 | SwinIR | 28.18 | SwinIR | 30.82 |
| | MPRNet [95] | 21.90 | DDMSNET [103] | 28.85 | CCN [67] | 31.34 |
| Multi Task | All-in-One [46] | 24.71 | All-in-One | 28.33 | All-in-One | 31.12 |
| | TransWea. [84] | 27.96 | TransWea. | 28.48 | TransWea. | 28.84 |
| | SemanIR (Ours) | 29.57 | SemanIR (Ours) | 30.76 | SemanIR (Ours) | 30.82 |

Table 9: ***Image demosaicking*** PSNR results.

| Datasets | Kodak | McMaster |
|---|---|---|
| Matlab | 35.78 | 34.43 |
| MMNet [38] | 40.19 | 37.09 |
| DDR [89] | 41.11 | 37.12 |
| DeepJoint [26] | 42.00 | 39.14 |
| RLDD [29] | 42.49 | 39.25 |
| DRUNet [100] | 42.68 | 39.39 |
| RNAN [110] | 43.16 | 39.70 |
| GRL [48] | 43.57 | 40.22 |
| SemanIR (Ours) | 43.62 | 40.68 |

window size of 8 to 33.38 dB with a window size of 32. These results suggest that larger window sizes enhance performance by capturing more contextual information.

## 4.2 Evaluation of SemanIR on Various IR Tasks

**Evaluation on Image deblurring.** Tab. 2 shows the quantitative results for single image motion deblurring on synthetic datasets (GoPro [62], HIDE [77]). Compared to the previous state-of-the-art Restormer [94], our SemanIR achieves significant PSNR improvement (*i.e.*, 0.52 dB) on the GoPro dataset and the second-best on the HIDE dataset. The visual results are shown in the *Appx.* E.

**Evaluation on JPEG CAR.** The experiments for color images are conducted with 4 image quality factors ranging from 10 to 40 under two settings (*i.e.*, † a single model is trained to handle multiple quality factors, and each model for each quality). The quantitative results shown in Tab. 4 indicate that our SemanIR achieves the best results on all the test sets across various quality factors among all the comparison methods for the color images. The visual comparisons in the *Appx.* E further supports the effectiveness of our method.

**Evaluation on Image Denoising.** We show color and grayscale image denoising results in Tab. 7 under two settings (*i.e.*, † one model for all noise levels $\sigma = \{15, 25, 50\}$ and each model for each noise level). For a fair comparison, both parameters and accuracy are reported for all the methods. For †, our SemanIR performs better on all test sets for color and grayscale image denoising than others. It is worth noting that we outperform DRUNet and Restormer with lower trainable parameters. For another setting, the proposed SemanIR also archives better results on CBSD68 and Urban100 for color image denoising, and on Set12 and Urban100 for grayscale denoising. These interesting comparisons validate the effectiveness of the proposed SemanIR and also indicate that our method has a higher generalization ability. The visual results in *Appx.* E also support that the proposed SemanIR can remove heavy noise corruption and preserve high-frequency image details, resulting in sharper edges and more natural textures without over-smoothness or over-sharpness problems.

**Evaluation in AWC.** We validate SemanIR in adverse weather conditions, including rain+fog (Test1), snow (SnowTest100K), and raindrops (RainDrop). PSNR is reported in Tab. 8. Our method achieves the best performance on Test1 (*i.e.*, 5.76% improvement) and SnowTest100k-L (*i.e.* 8.01% improvement), while the second-best PSNR on RainDrop compared to all other methods. See *Appx.* E for Visual comparisons.

**Evaluation on Image Demosaicking.** The quantitative results shown in 9 indicate that the proposed SemanIR performs best on both the Kodak and MaMaster test sets, especially 0.05dB and 0.45dB absolute improvement compared to the current state-of-the-art.

Table 10: ***Classical image SR*** results. Both lightweight and accurate models are summarized.

| Method | Scale | Params [M] | Set5 PSNR↑ | Set5 SSIM↑ | Set14 PSNR↑ | Set14 SSIM↑ | BSD100 PSNR↑ | BSD100 SSIM↑ | Urban100 PSNR↑ | Urban100 SSIM↑ | Manga109 PSNR↑ | Manga109 SSIM↑ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RCAN [109] | 2× | 15.44 | 38.27 | 0.9614 | 34.12 | 0.9216 | 32.41 | 0.9027 | 33.34 | 0.9384 | 39.44 | 0.9786 |
| SAN [18] | 2× | 15.71 | 38.31 | 0.9620 | 34.07 | 0.9213 | 32.42 | 0.9028 | 33.10 | 0.9370 | 39.32 | 0.9792 |
| HAN [63] | 2× | 63.61 | 38.27 | 0.9614 | 34.16 | 0.9217 | 32.41 | 0.9027 | 33.35 | 0.9385 | 39.46 | 0.9785 |
| IPT [10] | 2× | 115.48 | 38.37 | - | 34.43 | - | 32.48 | - | 33.76 | - | - | - |
| SwinIR [50] | 2× | 11.75 | 38.42 | 0.9623 | 34.46 | 0.9250 | 32.53 | 0.9041 | 33.81 | 0.9427 | 39.92 | 0.9797 |
| CAT-A [15] | 2× | 16.46 | 38.51 | 0.9626 | 34.78 | 0.9265 | 32.59 | 0.9047 | 34.26 | 0.9440 | 40.10 | 0.9805 |
| ART [99] | 2× | 16.40 | 38.56 | 0.9629 | 34.59 | 0.9267 | 32.58 | 0.9048 | 34.30 | 0.9452 | 40.24 | 0.9808 |
| EDT [47] | 2× | 11.48 | 38.63 | 0.9632 | 34.80 | 0.9273 | 32.62 | 0.9052 | 34.27 | 0.9456 | 40.37 | 0.9811 |
| SemanIR-S (Ours) | 2× | 11.87 | 38.57 | 0.9651 | 34.99 | 0.9300 | 32.65 | 0.9078 | 34.86 | 0.9472 | 40.45 | 0.9824 |
| SemanIR-B (Ours) | 2× | 19.90 | 38.61 | 0.9654 | 35.08 | 0.9304 | 32.69 | 0.9084 | 34.99 | 0.9455 | 40.59 | 0.9830 |
| RCAN [109] | 4× | 15.59 | 32.63 | 0.9002 | 28.87 | 0.7889 | 27.77 | 0.7436 | 26.82 | 0.8087 | 31.22 | 0.9173 |
| SAN [18] | 4× | 15.86 | 32.64 | 0.9003 | 28.92 | 0.7888 | 27.78 | 0.7436 | 26.79 | 0.8068 | 31.18 | 0.9169 |
| HAN [63] | 4× | 64.20 | 32.64 | 0.9002 | 28.90 | 0.7890 | 27.80 | 0.7442 | 26.85 | 0.8094 | 31.42 | 0.9177 |
| IPT [10] | 4× | 115.63 | 32.64 | - | 29.01 | - | 27.82 | - | 27.26 | - | - | - |
| SwinIR [50] | 4× | 11.90 | 32.92 | 0.9044 | 29.09 | 0.7950 | 27.92 | 0.7489 | 27.45 | 0.8254 | 32.03 | 0.9260 |
| CAT-A [15] | 4× | 16.60 | 33.08 | 0.9052 | 29.18 | 0.7960 | 27.99 | 0.7510 | 27.89 | 0.8339 | 32.39 | 0.9285 |
| ART [99] | 4× | 16.55 | 33.04 | 0.9051 | 29.16 | 0.7958 | 27.97 | 0.751 | 27.77 | 0.8321 | 32.31 | 0.9283 |
| EDT [47] | 4× | 11.63 | 33.06 | 0.9055 | 29.23 | 0.7971 | 27.99 | 0.7510 | 27.75 | 0.8317 | 32.39 | 0.9283 |
| SemanIR-S (Ours) | 4× | 12.02 | 33.02 | 0.9082 | 29.29 | 0.8026 | 27.96 | 0.7582 | 28.34 | 0.8467 | 32.48 | 0.9322 |
| SemanIR-B (Ours) | 4× | 20.04 | 33.08 | 0.9090 | 29.34 | 0.8037 | 27.98 | 0.7599 | 28.51 | 0.8467 | 32.56 | 0.9335 |



Figure 6: Visual comparison of classical image SR (4×) on Urban100. Best viewed by zooming.

**Evaluation on SR.** For the classical image SR, we compared our SemanIR with both recent lightweight and accurate SR models, and the quantitative results are shown in Tab. 10. Compared to EDT, SemanIR-base achieves significant improvements on Urban100 (*i.e.*, 0.72 dB and 0.76dB for 2× and 4× SR) and Manga109 datasets (*i.e.*, 0.22dB and 0.17 dB for 2× and 4× SR). Even the SemanIR-small consistently ranks as the runner-up across the majority of test datasets, all while maintaining a reduced number of trainable parameters. Visual results in both Fig. 6 and *Appx.* E also validate the effectiveness of the proposed SemanIR. Specifically, it is clear from the zoomed part in Fig. 6 that SemanIR can restore more details and structural content compared to other methods.

## 5 Conclusion

In this paper, we propose a novel approach, SemanIR, for ViTs-based image restoration, which experimentally validated that global cues are essential to restore degraded images well, but the most semantically related global cures play the major role. Specifically, to capture the key semantics, we propose to construct a semantic dictionary (*i.e.*, naively by self-similarity is enough) for storing only the most related $k$ semantic information and then use it as a reference for guiding the attention operation for making the attention operation pay more attention only to these key semantics. Furthermore, we share the key-semantic dictionary with all the upcoming transformer layers within the same stage since each stage of the transformer is typically at the same semantic level. This strategy significantly reduces the computational cost for IR and functions as loop optimization, continuously restoring degraded patches with their most semantically related patches, which share similar texture or structural information. Extensive experiments on 6 IR tasks validated the effectiveness of SemanIR, demonstrating that our method achieves new state-of-the-art performance.

## Acknowledgments and Disclosure of Funding

## References

[1] Eirikur Agustsson and Radu Timofte. NTIRE 2017 challenge on single image super-resolution: Dataset and study. In *CVPRW*, pages 126–135, 2017. 17

[2] Pablo Arbelaez, Michael Maire, Charless Fowlkes, and Jitendra Malik. Contour detection and hierarchical image segmentation. *IEEE TPAMI*, 33(5):898–916, 2010. 17

[3] Mark R Banham and Aggelos K Katsaggelos. Digital image restoration. *IEEE Signal Processing Magazine*, 14(2):24–41, 1997. 3

[4] Marco Bevilacqua, Aline Roumy, Christine Guillemot, and Marie Line Alberi-Morel. Low-complexity single-image super-resolution based on nonnegative neighbor embedding. In *BMVC*, 2012. 17

[5] Christopher M Bishop and Nasser M Nasrabadi. *Pattern Recognition and Machine Learning*, volume 4. Springer, 2006. 2

[6] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your ViT but faster. In *ICLR*, 2023. 2

[7] Antoni Buades, Bartomeu Coll, and J-M Morel. A non-local algorithm for image denoising. In *CVPR*, volume 2, pages 60–65, 2005. 3

[8] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, pages 213–229, 2020. 3

[9] Dexiong Chen, Leslie O'Bray, and Karsten Borgwardt. Structure-aware transformer for graph representation learning. In *ICML*, pages 3469–3489, 2022. 5

[10] Hanting Chen, Yunhe Wang, Tianyu Guo, Chang Xu, Yiping Deng, Zhenhua Liu, Siwei Ma, Chunjing Xu, Chao Xu, and Wen Gao. Pre-trained image processing transformer. In *CVPR*, pages 12299–12310, 2021. 3, 6, 10

[11] Liangyu Chen, Xiaojie Chu, Xiangyu Zhang, and Jian Sun. Simple baselines for image restoration. In *ECCV*, pages 17–33, 2022. 6, 18, 19

[12] Wei-Ting Chen, Hao-Yu Fang, Jian-Jiun Ding, Cheng-Che Tsai, and Sy-Yen Kuo. JSTASR: Joint size and transparency-aware snow removal algorithm based on modified partial convolution and veiling effect removal. In *ECCV*, pages 754–770, 2020. 9, 17

[13] Xiang Chen, Hao Li, Mingqiang Li, and Jinshan Pan. Learning a sparse transformer network for effective image deraining. In *CVPR*, pages 5896–5905, 2023. 3

[14] Xiangyu Chen, Xintao Wang, Jiantao Zhou, Yu Qiao, and Chao Dong. Activating more pixels in image super-resolution transformer. In *CVPR*, pages 22367–22377, 2023. 3, 7, 19

[15] Zheng Chen, Yulun Zhang, Jinjin Gu, Yongbing Zhang, Linghe Kong, and Xin Yuan. Cross aggregation transformer for image restoration. *NeurIPS*, 35:25478–25490, 2022. 3, 7, 10

[16] Sung-Jin Cho, Seo-Won Ji, Jun-Pyo Hong, Seung-Won Jung, and Sung-Jea Ko. Rethinking coarse-to-fine approach in single image deblurring. In *ICCV*, 2021. 6

[17] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Image denoising by sparse 3-D transform-domain collaborative filtering. *IEEE TIP*, 16(8):2080–2095, 2007. 3

[18] Tao Dai, Jianrui Cai, Yongbing Zhang, Shu-Tao Xia, and Lei Zhang. Second-order attention network for single image super-resolution. In *CVPR*, pages 11065–11074, 2019. 10

[19] Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. FlashAttention: Fast and memory-efficient exact attention with IO-awareness. In *NeurIPS*, 2022. 2, 6

[20] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2020. 2, 3

[21] Max Ehrlich, Larry Davis, Ser-Nam Lim, and Abhinav Shrivastava. Quantization guided JPEG artifact correction. In *ECCV*, pages 293–309, 2020. 8

[22] Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. Pointwise shape-adaptive DCT for high-quality denoising and deblocking of grayscale and color images. *IEEE TIP*, 16(5):1395–1411, 2007. 17

[23] Rich Franzen. Kodak lossless true color image suite. *source: http://r0k. us/graphics/kodak*, 4(2), 1999. 17

[24] Hongyun Gao, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Dynamic scene deblurring with parameter selective sharing and nested skip connections. In *CVPR*, 2019. 6

[25] Sicheng Gao, Xuhui Liu, Bohan Zeng, Sheng Xu, Yanjing Li, Xiaoyan Luo, Jianzhuang Liu, Xiantong Zhen, and Baochang Zhang. Implicit diffusion models for continuous super-resolution. In *CVPR*, pages 10021–10030, 2023. 3

[26] Michaël Gharbi, Gaurav Chaurasia, Sylvain Paris, and Frédo Durand. Deep joint demosaicking and denoising. *ACM TOG*, 35(6):1–12, 2016. 9

[27] Gene H Golub, Per Christian Hansen, and Dianne P O'Leary. Tikhonov regularization and total least squares. *SIAM Journal on Matrix Analysis and Applications*, 21(1):185–194, 1999. 3

[28] Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph domains. In *IJCNN*, volume 2, pages 729–734, 2005. 3

[29] Yu Guo, Qiyu Jin, Gabriele Facciolo, Tieyong Zeng, and Jean-Michel Morel. Residual learning for effective joint demosaicing-denoising. *arXiv preprint arXiv:2009.06205*, 2020. 9

[30] Jia-Bin Huang, Abhishek Singh, and Narendra Ahuja. Single image super-resolution from transformed self-exemplars. In *CVPR*, pages 5197–5206, 2015. 17

[31] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *CVPR*, pages 1125–1134, 2017. 9, 17

[32] Bo Jiang, Yao Lu, Xiaosheng Chen, Xinhai Lu, and Guangming Lu. Graph attention in attention network for image denoising. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2023. 3

[33] Jiaxi Jiang, Kai Zhang, and Radu Timofte. Towards flexible blind JPEG artifacts removal. In *ICCV*, pages 4997–5006, 2021. 8

[34] Justin Johnson, Alexandre Alahi, and Li Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, pages 694–711, 2016. 19

[35] Salman Khan, Muzammal Naseer, Munawar Hayat, Syed Waqas Zamir, Fahad Shahbaz Khan, and Mubarak Shah. Transformers in vision: A survey. *ACM computing surveys*, 54(10s):1–41, 2022. 3

[36] Daisuke Kiku, Yusuke Monno, Masayuki Tanaka, and Masatoshi Okutomi. Beyond color difference: Residual interpolation for color image demosaicking. *IEEE TIP*, 25(3):1288–1300, 2016. 9

[37] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 19

[38] Filippos Kokkinos, Stamatios Lefkimmiatis, and B A. Iterative joint image demosaicking and denoising using a residual denoising network. *IEEE TIP*, 28(8):4177–4188, 2019. 9

[39] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. DeblurGAN: Blind motion deblurring using conditional adversarial networks. In *CVPR*, 2018. 6

[40] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. DeblurGAN-v2: Deblurring (orders-of-magnitude) faster and better. In *ICCV*, 2019. 6

[41] Wei-Sheng Lai, Jia-Bin Huang, Narendra Ahuja, and Ming-Hsuan Yang. Deep laplacian pyramid networks for fast and accurate super-resolution. In *CVPR*, pages 624–632, 2017. 3

[42] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. 2

[43] Hunsang Lee, Hyesong Choi, Kwanghoon Sohn, and Dongbo Min. Knn local attention for image restoration. In *CVPR*, pages 2139–2149, 2022. 3, 6, 18

[44] Juho Lee, Yoonho Lee, Jungtaek Kim, Adam Kosiorek, Seungjin Choi, and Yee Whye Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *ICML*, pages 3744–3753, 2019. 5

[45] Ruoteng Li, Loong-Fah Cheong, and Robby T Tan. Heavy rain image restoration: Integrating physics model and conditional adversarial learning. In *CVPR*, pages 1633–1642, 2019. 9, 17, 24

[46] Ruoteng Li, Robby T Tan, and Loong-Fah Cheong. All in one bad weather removal using architectural search. In *CVPR*, pages 3175–3185, 2020. 9, 17, 24

[47] Wenbo Li, Xin Lu, Jiangbo Lu, Xiangyu Zhang, and Jiaya Jia. On efficient transformer and image pre-training for low-level vision. *arXiv preprint arXiv:2112.10175*, 2021. 9, 10

[48] Yawei Li, Yuchen Fan, Xiaoyu Xiang, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. Efficient and explicit modelling of image hierarchies for image restoration. In *CVPR*, pages 18278–18289, 2023. 1, 2, 3, 7, 8, 9

[49] Yawei Li, Kai Zhang, Jingyun Liang, Jiezhang Cao, Ce Liu, Rui Gong, Yulun Zhang, Hao Tang, Yun Liu, Denis Demandolx, Rakesh Ranjan, Radu Timofte, and Luc Van Gool. LSDIR: A large scale dataset for image restoration. In *CVPRW*, pages 1775–1787, 2023. 3

[50] Jingyun Liang, Jiezhang Cao, Guolei Sun, Kai Zhang, Luc Van Gool, and Radu Timofte. SwinIR: Image restoration using Swin transformer. In *ICCVW*, pages 1833–1844, 2021. 2, 3, 7, 8, 9, 10, 17

[51] Bee Lim, Sanghyun Son, Heewon Kim, Seungjun Nah, and Kyoung Mu Lee. Enhanced deep residual networks for single image super-resolution. In *CVPRW*, pages 1132–1140, 2017. 3, 17

[52] Yun-Fu Liu, Da-Wei Jaw, Shih-Chia Huang, and Jenq-Neng Hwang. DesnowNet: Context-aware deep network for snow removal. *IEEE TIP*, 27(6):3064–3073, 2018. 9, 17

[53] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021. 2

[54] Zhilei Liu, Le Li, Yunpeng Wu, and Cuicui Zhang. Facial expression restoration based on improved graph convolutional networks. In *MMM*, pages 527–539. Springer, 2020. 3

[55] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2018. 19

[56] Ziwei Luo, Fredrik K Gustafsson, Zheng Zhao, Jens Sjölund, and Thomas B Schön. Image restoration with mean-reverting stochastic differential equations. *arXiv preprint arXiv:2301.11699*, 2023. 3

[57] Kede Ma, Zhengfang Duanmu, Qingbo Wu, Zhou Wang, Hongwei Yong, Hongliang Li, and Lei Zhang. Waterloo exploration database: New challenges for image quality assessment models. *IEEE TIP*, 26(2):1004–1016, 2016. 17

[58] David Martin, Charless Fowlkes, Doron Tal, and Jitendra Malik. A database of human segmented natural images and its application to evaluating segmentation algorithms and measuring ecological statistics. In *ICCV*, pages 416–423, 2001. 17

[59] Yusuke Matsui, Kota Ito, Yuji Aramaki, Azuma Fujimoto, Toru Ogawa, Toshihiko Yamasaki, and Kiyoharu Aizawa. Sketch-based manga retrieval using manga109 dataset. *Multimedia Tools and Applications*, 76(20):21811–21838, 2017. 17

[60] Rafael Molina, Jorge Núñez, Francisco J Cortijo, and Javier Mateos. Image restoration in astronomy: a Bayesian perspective. *IEEE Signal Processing Magazine*, 18(2):11–29, 2001. 1

[61] Chong Mou, Jian Zhang, and Zhuoyuan Wu. Dynamic attentive graph learning for image restoration. In *ICCV*, pages 4328–4337, 2021. 3

[62] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, pages 3883–3891, 2017. 6, 9, 17

[63] Ben Niu, Weilei Wen, Wenqi Ren, Xiangde Zhang, Lianping Yang, Shuzhen Wang, Kaihao Zhang, Xiaochun Cao, and Haifeng Shen. Single image super-resolution via a holistic attention network. In *ECCV*, pages 191–207, 2020. 10

[64] Dongwon Park, Dong Un Kang, Jisoo Kim, and Se Young Chun. Multi-temporal recurrent neural networks for progressive non-uniform single image deblurring with incremental temporal training. In *ECCV*, 2020. 6

[65] Kuldeep Purohit, Maitreya Suin, AN Rajagopalan, and Vishnu Naresh Boddeti. Spatially-adaptive image restoration using distortion-guided networks. In *ICCV*, 2021. 6

[66] Rui Qian, Robby T Tan, Wenhan Yang, Jiajun Su, and Jiaying Liu. Attentive generative adversarial network for raindrop removal from a single image. In *CVPR*, pages 2482–2491, 2018. 9, 17

[67] Ruijie Quan, Xin Yu, Yuanzhi Liang, and Yi Yang. Removing raindrops and rain streaks in one go. In *CVPR*, pages 9147–9156, 2021. 9, 17

[68] Yuhui Quan, Shijie Deng, Yixin Chen, and Hui Ji. Deep learning for seeing through window with raindrops. In *ICCV*, pages 2463–2471, 2019. 9, 17

[69] Yongming Rao, Wenliang Zhao, Benlin Liu, Jiwen Lu, Jie Zhou, and Cho-Jui Hsieh. DynamicViT: Efficient vision transformers with dynamic token sparsification. *NeurIPS*, 34:13937–13949, 2021. 2

[70] Bin Ren, Yawei Li, Nancy Mehta, Radu Timofte, Hongyuan Yu, Cheng Wan, Yuxin Hong, Bingnan Han, Zhuoyuan Wu, Yajun Zou, et al. The ninth ntire 2024 efficient super-resolution challenge report. In *CVPRW*, pages 6595–6631, 2024. 1

[71] Bin Ren, Yahui Liu, Yue Song, Wei Bi, Rita Cucchiara, Nicu Sebe, and Wei Wang. Masked jigsaw puzzle: A versatile position embedding for vision transformers. In *CVPR*, pages 20382–20391, 2023. 2

[72] Mengwei Ren, Mauricio Delbracio, Hossein Talebi, Guido Gerig, and Peyman Milanfar. Multiscale structure guided diffusion for image deblurring. In *ICCV*, pages 10721–10733, 2023. 6

14

[73] William Hadley Richardson. Bayesian-based iterative method of image restoration. *Journal of the Optical Society of America*, 62(1):55–59, 1972. 3

[74] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE Transactions on Neural Networks*, 20(1):61–80, 2008. 3

[75] MI Sezan and Henry Stark. Image restoration by the method of convex projections: Part 2-applications and numerical results. *IEEE TMI*, 1(2):95–101, 1982. 1

[76] HR Sheikh. Live image quality assessment database release 2. *http://live. ece. utexas. edu/research/quality*, 2005. 17

[77] Ziyi Shen, Wenguan Wang, Xiankai Lu, Jianbing Shen, Haibin Ling, Tingfa Xu, and Ling Shao. Human-aware motion deblurring. In *ICCV*, pages 5572–5581, 2019. 9, 17

[78] Martin Simonovsky and Nikos Komodakis. Dynamic edge-conditioned filters in convolutional neural networks on graphs. In *CVPR*, pages 3693–3702, 2017. 3

[79] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015. 19

[80] Maitreya Suin, Kuldeep Purohit, and A. N. Rajagopalan. Spatially-attentive patch-hierarchical network for adaptive motion deblurring. In *CVPR*, 2020. 6

[81] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Jiaya Jia. Scale-recurrent network for deep image deblurring. In *CVPR*, 2018. 6

[82] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Herve Jegou. Training data-efficient image transformers & distillation through attention. In *ICML*, 2021. 3

[83] Zhengzhong Tu, Hossein Talebi, Han Zhang, Feng Yang, Peyman Milanfar, Alan Bovik, and Yinxiao Li. MAXIM: Multi-axis mlp for image processing. In *CVPR*, pages 5769–5780, 2022. 2

[84] Jeya Maria Jose Valanarasu, Rajeev Yasarla, and Vishal M Patel. TransWeather: Transformer-based restoration of images degraded by adverse weather conditions. In *CVPR*, pages 2353–2363, 2022. 9, 17, 20

[85] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, 2017. 2, 3, 5

[86] Shashanka Venkataramanan, Amir Ghodrati, Yuki M Asano, Fatih Porikli, and Amir Habibian. Skip-attention: Improving vision transformers by paying less attention. In *ICLR*, 2024. 2

[87] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, pages 7794–7803, 2018. 3

[88] Yinhuai Wang, Jiwen Yu, and Jian Zhang. Zero-shot image restoration using denoising diffusion null-space model. *ICLR*, 2023. 3

[89] Jiqing Wu, Radu Timofte, and Luc Van Gool. Demosaicing based on directional difference regression and efficient regression priors. *IEEE TIP*, 25(8):3862–3874, 2016. 9

[90] Zhuofan Xia, Xuran Pan, Shiji Song, Li Erran Li, and Gao Huang. Vision transformer with deformable attention. In *CVPR*, pages 4794–4803, 2022. 2, 5

[91] Linwei Ye, Mrigank Rochan, Zhi Liu, and Yang Wang. Cross-modal self-attention network for referring image segmentation. In *CVPR*, pages 10502–10511, 2019. 3

[92] Hongxu Yin, Arash Vahdat, Jose M Alvarez, Arun Mallya, Jan Kautz, and Pavlo Molchanov. A-ViT: Adaptive tokens for efficient vision transformer. In *CVPR*, pages 10809–10818, 2022. 2

[93] Zongsheng Yue, Jianyi Wang, and Chen Change Loy. ResShift: Efficient diffusion model for image super-resolution by residual shifting. *arXiv preprint arXiv:2307.12348*, 2023. 3

[94] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, and Ming-Hsuan Yang. Restormer: Efficient transformer for high-resolution image restoration. In *CVPR*, pages 5728–5739, 2022. 3, 6, 7, 9, 18

[95] Syed Waqas Zamir, Aditya Arora, Salman Khan, Munawar Hayat, Fahad Shahbaz Khan, Ming-Hsuan Yang, and Ling Shao. Multi-stage progressive image restoration. In *CVPR*, pages 14821–14831, 2021. 2, 6, 9

[96] Roman Zeyde, Michael Elad, and Matan Protter. On single image scale-up using sparse-representations. In *Curves and Surfaces*, pages 711–730. Springer, 2010. 17

[97] Hongguang Zhang, Yuchao Dai, Hongdong Li, and Piotr Koniusz. Deep stacked hierarchical multi-patch network for image deblurring. In *CVPR*, 2019. 6

[98] Jiale Zhang, Yulun Zhang, Jinjin Gu, Jiahua Dong, Linghe Kong, and Xiaokang Yang. Xformer: Hybrid X-shaped transformer for image denoising. *arXiv preprint arXiv:2303.06440*, 2023. 9

[99] Jiale Zhang, Yulun Zhang, Jinjin Gu, Yongbing Zhang, Linghe Kong, and Xin Yuan. Accurate image restoration with attention retractable transformer. In *ICLR*, 2023. 3, 5, 10

[100] Kai Zhang, Yawei Li, Wangmeng Zuo, Lei Zhang, Luc Van Gool, and Radu Timofte. Plug-and-play image restoration with deep denoiser prior. *IEEE TPAMI*, 2021. 8, 9

[101] Kai Zhang, Wangmeng Zuo, Yunjin Chen, Deyu Meng, and Lei Zhang. Beyond a gaussian denoiser: Residual learning of deep cnn for image denoising. *IEEE TIP*, 26(7):3142–3155, 2017. 17

[102] Kai Zhang, Wangmeng Zuo, and Lei Zhang. FFDNet: Toward a fast and flexible solution for cnn-based image denoising. *IEEE TIP*, 27(9):4608–4622, 2018. 9

[103] Kaihao Zhang, Rongqing Li, Yanjiang Yu, Wenhan Luo, and Changsheng Li. Deep dense multi-scale network for snow removal using semantic and depth priors. *IEEE TIP*, 30:7419–7431, 2021. 9, 17

[104] Kaihao Zhang, Wenhan Luo, Yiran Zhong, Lin Ma, Bjorn Stenger, Wei Liu, and Hongdong Li. Deblurring by realistic blurring. In *CVPR*, 2020. 6

[105] Leheng Zhang, Yawei Li, Xingyu Zhou, Xiaorui Zhao, and Shuhang Gu. Transcending the limit of local window: Advanced super-resolution transformer with adaptive token dictionary. *arXiv preprint arXiv:2401.08209*, 2024. 3

[106] Lei Zhang, Xiaolin Wu, Antoni Buades, and Xin Li. Color demosaicking by local directional interpolation and nonlocal adaptive thresholding. *Journal of Electronic Imaging*, 20(2):023016, 2011. 17

[107] Qiming Zhang, Jing Zhang, Yufei Xu, and Dacheng Tao. Vision transformer with quadrangle attention. *IEEE TPAMI*, 2024. 2, 5

[108] Shuoxi Zhang, Hanpeng Liu, Stephen Lin, and Kun He. You only need less attention at each stage in vision transformers. In *CVPR*, pages 6057–6066, 2024. 2, 3

[109] Yulun Zhang, Kunpeng Li, Kai Li, Lichen Wang, Bineng Zhong, and Yun Fu. Image super-resolution using very deep residual channel attention networks. In *ECCV*, pages 286–301, 2018. 10

[110] Yulun Zhang, Kunpeng Li, Kai Li, Bineng Zhong, and Yun Fu. Residual non-local attention networks for image restoration. *arXiv preprint arXiv:1903.10082*, 2019. 3, 9

[111] Haiyu Zhao, Yuanbiao Gou, Boyun Li, Dezhong Peng, Jiancheng Lv, and Xi Peng. Comprehensive and delicate: An efficient transformer for image restoration. In *CVPR*, pages 14122–14132, 2023. 6

[112] Maria Zontak and Michal Irani. Internal statistics of a single natural image. In *CVPR*, pages 977–984, 2011. 2

16

# A Experimental Protocals

## A.1 Training/Testing Datasets

**JPEG compression artifact removal.** For JPEG compression artifact removal, the JPEG image is compressed by the `cv2` JPEG compression function. The compression function is characterized by the quality factor. We investigated four compression quality factors including 10, 20, 30, and 40. The smaller the quality factor, the more the image is compressed, meaning a lower quality.

- The training datasets: DIV2K [1], Flickr2K [51], and WED [57].
- The test datasets: Classic5 [22], LIVE1 [76], Urban100 [30], BSD500 [2].

**Image Denoising.** For image denoising, we conduct experiments on both color and grayscale image denoising. During training and testing, noisy images are generated by adding independent additive white Gaussian noise (AWGN) to the original images. The noise levels are set to $\sigma = 15, 25, 50$. We train individual networks at different noise levels. The network takes the noisy images as input and tries to predict noise-free images.

- The training datasets: DIV2K [1], Flickr2K [51], WED [57], and BSD400 [58].
- The test datasets for color image: CBSD68 [58], Kodak24 [23], McMaster [106], and Urban100 [30].
- The test datasets for grayscale image: Set12 [101], BSD68 [58], and Urban100 [30].

**Image Demosaicking.** For image demosaicking, the mosaic image is generated by applying a Bayer filter on the ground-truth image. Then the network try to restore high-quality image. The mosaic image is first processed by the default `Matlab` demosaic function and then passed to the network as input.

- The training datasets: DIV2K [1] and Flickr2K [51].
- The test datasets: Kodak [23], McMaster [106].

**IR in Adverse Weather Conditions.** For IR in adverse weather conditions, the model is trained on a combination of images degraded by a variety of adverse weather conditions. The same training and test dataset is used as in Transweather [84]. The training data comprises 9,000 images sampled from Snow100K [52], 1,069 images from Raindrop [66], and 9,000 images from Outdoor-Rain [45]. Snow100K includes synthetic images degraded by snow, Raindrop consists of real raindrop images, and Outdoor-Rain contains synthetic images degraded by both fog and rain streaks. The proposed method is tested on both synthetic and real-world datasets.

- The comparison methods in Tab. 6 of our main manuscript: pix2pix [31], HRGAN [45], SwinIR [50], All-in-One [46], Transweather [84], DesnowNet [52], JSTASR [12], DDM-SNET [103], Attn. GAN [66], [68], and CCGAN [67].
- The test datasets: test1 dataset [45, 46], the RainDrop test dataset [66], and the Snow100k-L test.

**Image SR.** For image SR, the LR image is synthesized by `Matlab` bicubic downsampling function before the training. We investigated the upscalingg factors $\times 2$, $\times 3$, and $\times 4$.

- The training datasets: DIV2K [1] and Flickr2K [51].
- The test datasets: Set5 [4], Set14 [96], BSD100 [58], Urban100 [30], and Manga109 [59].

**Image Deblurring.** For single-image motion deblurring,

- The training datasets: GoPro [62] .
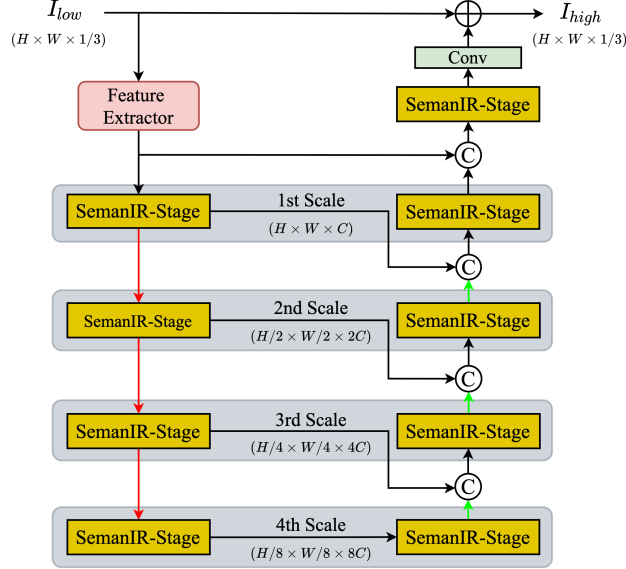- The test datasets: GoPro [62] and HIDE [77].

Figure 7: The U-shaped hierarchical architecture (Archi-V2) of the proposed SemanIR for Image Restoration. Note that this U-shaped one is used for image JPEG CAR, image denoising, image demosaicking, IR in AWC, and image deblurring. Symbols $\oplus$ and ⓒ denote the element-wise addition and channel-wise concatenation. The downsample and upsample operations are denoted by red and green arrows.

Table 11: The details of the transformer stages and layers per stage of SemanIR for both architectures.

|  | Archi-V1 (Columnar-shape) | | Archi-V2 (U-shape) | | |
|---|---|---|---|---|---|
|  | SemanIR-small | SemanIR-base | Down Stages | Up Stages | Final Stage |
| Num. of SemanIR Stages | 6 | 8 | 4 | 4 | 1 |
| Num. of SemanIR layer per stage | 6 | 8 | 6 | 6 | 6 |

## A.2 Model Architecture

In the proposed SemanIR, we adopt two kinds of base architecture *i.e.*, the widely used multi-stage one shown in Fig.1 of our main manuscript (Archi-V1) and a U-shaped hierarchical one shown in Fig. 7 (Archi-V2) for taking patterns of various scales into account (Note that 1/3 of $I_{low}$ and $I_{high}$ in Fig. 7 denotes the grayscale/color image cases). This is consistent with previous methods such as Restormer [94], KiT [43], and NAFNet [11].

Note that The feature extractor for both architectures is implemented as a simple convolution and converts the input image into feature maps. The image reconstructor for Archi-V1 takes the rich features calculated by the previous operations and estimates a recovered image.

In addition to introducing the two base architectures of the proposed SemanIR, we have provided comprehensive details of its structure in Table 11. This table outlines the number of SemanIR stages and the distribution of layers within each stage, offering a thorough understanding of our model's architecture.

## A.3 Efficiency Analysis

We provide a complexity comparison among the standard multi-head self-attention (MSA), the Window-wise MSA (W-MSA), and the proposed KeySemanIR MSA (SemanIR-MSA) in the Tab. 12. $(H, W, C)$ indicate the feature size, $M$ represents the window size, and $h$ denotes the number of heads. It is commonly demonstrated and proven that the complexity of the W-MSA is much lower than that of the standard MSA, *i.e.*,

$$\mathcal{O}(4HWC^2 + 2(M)^2HWC) < \mathcal{O}(4HWC^2 + 2(HW)^2C) \tag{3}$$

18

Table 12: The Computation Complexity Comparison.

|  | Time Complexity | Space Complexity |
|---|---|---|
| MSA | $\mathcal{O}(4HWC^2 + 2(HW)^2C)$ | $\mathcal{O}(4HWC^2 + 2h(HW)^2C)$ |
| W-MSA | $\mathcal{O}(4HWC^2 + 2(M)^2HWC)$ | $\mathcal{O}(4HWC^2 + 2h(M)^2HWC)$ |
| SemanIR MSA (Ours) | $\mathcal{O}(4HWC^2 + 2kHWC)$ | $\mathcal{O}(4HWC^2 + 2hkHWC)$ |

To better understand the efficiency of the proposed method, it should be considered together with all transformer layers within a certain stage. Specifically, take one stage, which contains 6 transformer Layers, as an example (To simplify the illustration, we omit the convolution operation at the end of each stage).

First, the total complexity of W-MSA within each stage can be calculated as:

$$\mathcal{O}(6 \times [4HWC^2 + 2(M)^2HWC]) \tag{4}$$

Second, similarly, the complexity of the proposed SemanIR-MSA can be calculated as follows. ($\mathcal{O}(HWC)$ indicates the complexity of the key-semantic dictionary construction at the start of each transformer stage. All the other layers then share it, hence it is calculated only once).

$$\mathcal{O}(6 \times [4HWC^2 + 2kHWC] + (HW)^2C) \tag{5}$$

Third, a simple subtraction can be done as follows to validate that the proposed method is more efficient compared to the W-MSA within each Transformer stage:

$$\mathcal{O}(6 \times [4HWC^2 + 2(M)^2HWC] - (6 \times [4HWC^2 + 2kHWC] + (HW)^2C)) \\ = \mathcal{O}((12M^2 - 12k - HW)HWC) \tag{6}$$

In the last equation, to provide further clarity, let's consider a common setting where the window size $M = 7$ and the patch size is 16. The height $H$ and the width $W$ of the feature map are 64. The number of pixels within the window is approximately $(7 \times 7) \times (16 \times 16)$. In the proposed SemanIR, the k value is set to 512 or randomly sampled from [64, 128, 256, 384, 512]. To this end, we have:

$$\mathcal{O}((12M^2 - 12k - HW)HWC) = \mathcal{O}(12 \times (7 \times 7) \times (16 \times 16) - 12 \times 512 - 64 \times 64) \\ = \mathcal{O}(150528 - 6144 - 4096) >> 0 \tag{7}$$

This shows that the complexity is significantly greater than zero.

Based on the above analysis, it can be concluded that, together with the proposed transformer layer, constructing the key-semantic dictionary at the start of each stage leads to greater efficiency.

### A.4 Training Details

Our method explores **6** IR tasks, and the training settings vary slightly for each task. These differences encompass the architecture of the proposed SemanIR, variations in the choice of the optimizer, and loss functions. Each experiments are conducted on 4 NVIDIA Tesla V100 32G GPUs.

**Architecture.** We use the columnar multi-stage architecture (without changing the feature map resolution and number of channels) for image SR and the U-shaped architecture for the other tasks including image denoising, image deblurring, and other tasks. The strategy of using multiple architectures is also explored by the previous method [11, 14].

**Optimizer.** We adopt the same optimizer as all other comparison methods, *i.e.*, Adam [37], for IR in AWC, and AdamW [55] for the rest IR tasks.

**Loss Function.** We adopt the same loss function as all other comparison methods, *i.e.*, smooth L1 loss and VGG loss [34, 79] for IR in AWC, the Charbonnier loss for Deblurring, and L1 loss for the rest IR tasks.

**Batch Size and Patch Size.** We keep the similar batch size as other comparison methods, *i.e.*, (Batch size = 16, Patch Size = 64) for JPEG CAR, denoising, demosaicking, and SR. (Batch Size = 32, Patch Size = 16) for IR in AWC. (Batch Size = 8, Patch Size = 192) for deblurring.
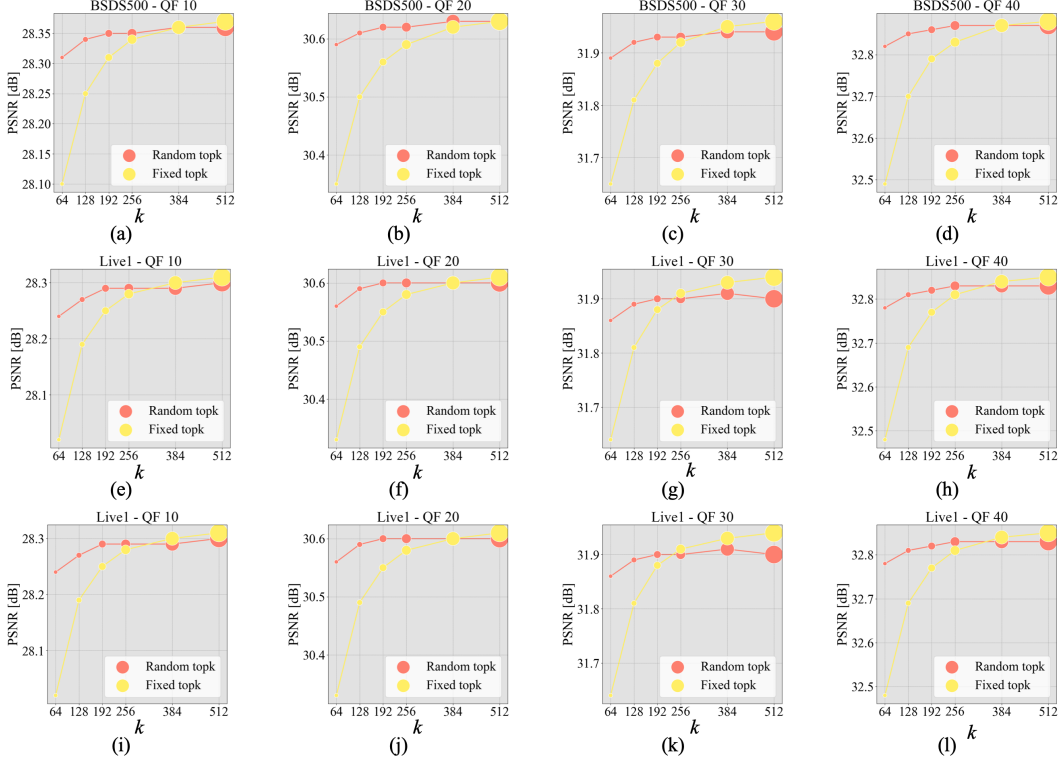
Figure 8: Ablation study on the impact of $k$ for **Color JPEG CAR** on CBSD68 (a-d), Live (e-h), and Urban100 (i-l) datasets with $QF = \{10, 20, 30, 40\}$.

**Learning Rate Schedule.** For all the IR tasks, similar to other comparison methods, we set the initial learning rate to $2 \times 10^{-4}$, and then the half-decay is adopted during the training. Note that the training iteration for JPEG CAR, denoising, demosaicking, and SR is set to 1M. For IR in AWC and debluriing, it is set to 750K.

## A.5 Evaluation Introduction

Note that the results of all the comparison methods are reported from their original papers. The details of the evaluation metric (*i.e.*, SSIM, PSNR) are described as follows:

**JPEG compression artifact removal.** For color image JPEG compression artifact removal, the PSNR is reported on the RGB channels while for grayscale image JPEG compression artifact removal, the PSNR is reported on the Y channel.

**Image Denoising.** For color image denoising, the PSNR is reported on the RGB channels while for grayscale image denoising, the PSNR is reported on the Y channel.

**Image Demosaicking.** For the comparison between different methods, PSNR is reported on the RGB channels.

**IR in Adverse Weather Conditions.** We adopted the same PSNR evaluation metric used in Transweather [84].

**Image SR.** The PSNR is reported on the Y channel.

**Image Deblurring.** The PSNR and SSIM on the RGB channels are reported.

## B Limitations

This study faces a task-specific limitation: each image restoration task requires training a separate network. While efforts have been made to train models for varying degradation levels within specific
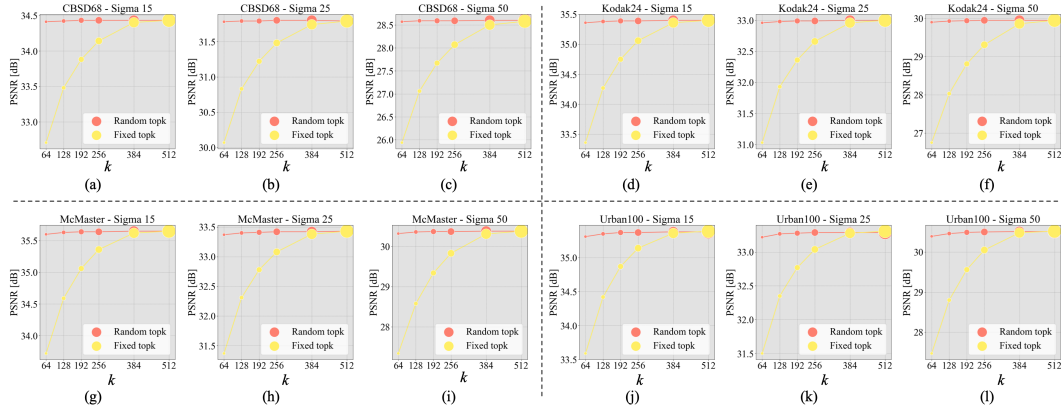
Figure 9: Ablation study on the impact of $k$ for **Color Image Denoising** on CBSD68 (a-c), Kodak24 (d-f), McMaster (g-i), and Urban100 (j-l) datasets with $\sigma = \{15, 25, 50\}$.
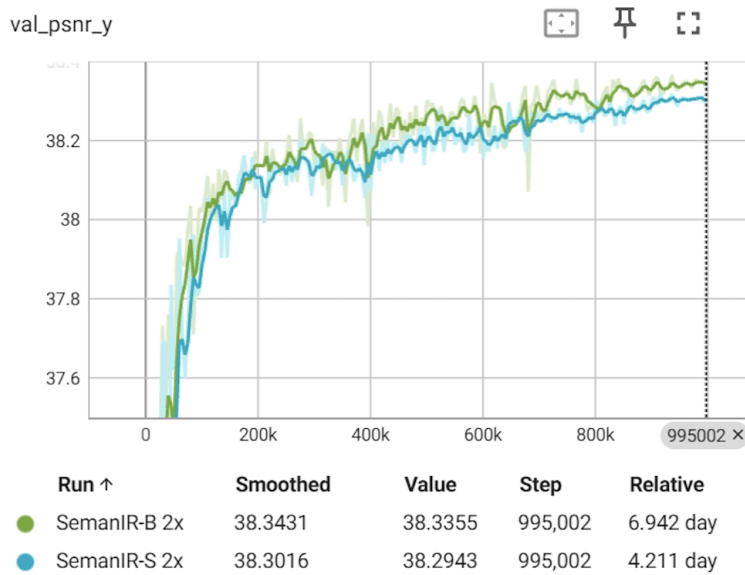


Figure 10: Training log shows the convergence of the proposed SemanIR during training. The upscaling factor is $2\times$.

types, such as image denoising and removal of JPEG compression artifacts, this approach still leads to inefficiencies in model training. It constrains the utility of the trained networks. A potential future enhancement involves developing a mechanism enabling a network to handle diverse image degradation types and levels. Another challenge is the substantial parameter requirement of the proposed SemanIR, which operates within a tens-of-millions parameter budget. Deploying such a large image IR network on handheld devices with limited resources is challenging, if not unfeasible. Therefore, a promising research direction is the creation of more efficient versions of SemanIR, integrating non-local context more effectively, to overcome these limitations.

## C   Impact Statement

This paper introduces a transformer-based approach that significantly enhances the efficiency and performance of image restoration tasks, including image deblurring, JPEG CAR, image denoising, IR in adverse weather conditions, demosaicking, and image super-resolution. The proposed SemanIR's notable efficiency improvement holds promise for resource-effective implementations in real-world applications. This elevated performance creates opportunities for enhanced image quality across diverse domains. While our primary contribution lies in the technical aspects of Machine Learning, we are cognizant of potential societal impacts, particularly in healthcare, surveillance, and digital
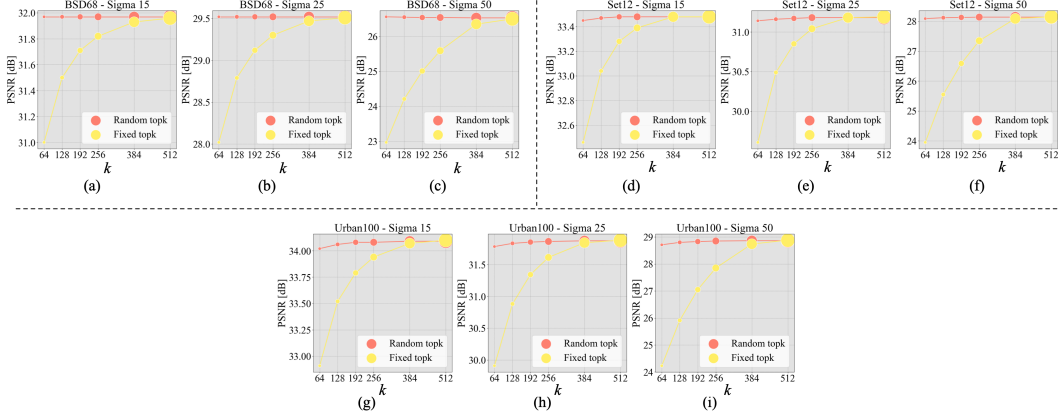
Figure 11: Ablation study on the impact of $k$ for **Grayscale Image Denoising** on BSD68 (a-c), Set12 (d-f), and Urban100 (g-i) datasets with $\sigma = \{15, 25, 50\}$.

imaging. As with any technology, ongoing vigilance in ethical considerations during deployment is essential, ensuring responsible use and proactively addressing any unintended consequences.

## D   More Ablation Analyses

Besides the ablation studies presented in our main manuscript, we further provide the following two analyses:

**Convergence Visualization.** The training log of the proposed SemanIR for image SR is shown in Fig. 10. The log is reported for the PSNR on the Set5 dataset during training. Two versions of the proposed method including SemanIR-S and SemanIR-B are shown in this figure. As shown in this figure, the proposed network converges gradually during the training.

**The Impact of the $k$ in Key-Semantic Dictionary Construction under Various IR Tasks.** To explore how the $k$ value of top-k will affect the IR performance of the proposed SemanIR. We conduct exhaustive experiments on JPEG compression artifact reduction for color images under different QF values (*i.e.*, QF = [10, 20, 30, 40]), image denoising for both color and grayscale images under different noise levels (*i.e.*, $\sigma$ = [15, 25, 50]), as well as image SR under different scales (*i.e.*, $2\times$, $3\times$, $4\times$) with the proposed SemanIR. Note that all the experiments for each IR task are conducted under two kinds of top-k settings, *i.e.*, (i) $k$ was randomly sampled from the range [64, 512] during the overall training phase, and (ii) $k$ was held constant at 512 throughout the training phase. For inference, $k$ was configured to the specified value for both settings.

The results of the JPEG CAR in terms of the hyper-parameters $k$ under different training settings during inference for color image are shown in Fig. 8. It is clear that for the color JPEG CAR task when $k$ is set to 64 during inference, there is a huge performance cat between the random top-k setting and the fixed top-k setting. In addition, the fixed top-k setting performs well or sometimes even a bit better than the random top-k setting only when $k$ is also set to the same number (*i.e.*, 512). With the decrease of $k$ during inference for the fixed top-k setting, the PSNR drops largely marginally for all the datasets under every kind of degraded QF factor.

The results of the image denoising in terms of hyper-parameters $k$ under different training settings during inference for both color image and grayscale image are shown in Fig. 9 and Fig. 11. All the experimental results on various datasets (*i.e.*, BSD68/CBSD68, Kodak24, McMaster, and Urban100) share a similar trend for image denoising compared to the JPEG CAR task. The random top-k setting can maintain a relatively stable PSNR score under different $k$ during inference compared to its fixed counterpart. In addition, a decent result can be obtained for the fixed top-k setting only when the $k$ is set to the same (*i.e.*, 512) during the inference.

The results of the image SR in terms of hyper-parameters $k$ under different training settings during inference for color images with different scale factors (*i.e.*, $2\times$, $3\times$, and $4\times$) are also provided in Fig. 12. All experiments are conducted in various datasets (*i.e.*, BSD100, Manga109, Set5, Set14, and Urban100). It shows that for datasets like BSD100, Set14, and Urban100, a similar trend can be also
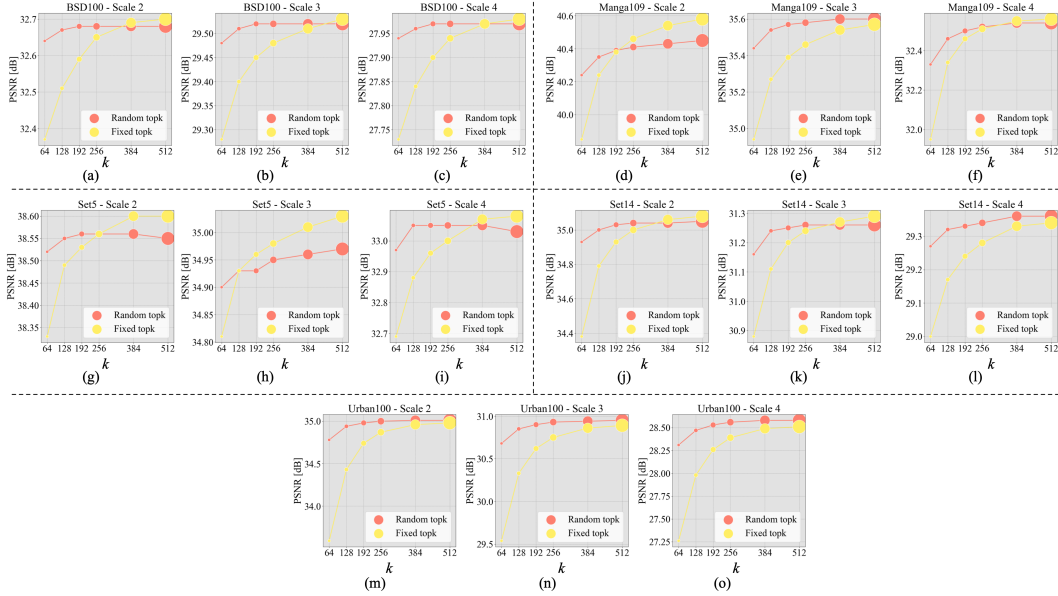
Figure 12: Ablation study on the impact of $k$ for **Image SR** with **SemanIR-B** on BSD100 (a-c), Manga109 (d-f), Set5 (g-i), Set14 (i-l) and Urban100 (m-o) datasets with $scale = (2\times, 3\times, \text{and } 4\times)$.

observed in Fig. 12 compared to JPEG CAR and image denoising tasks, *i.e.*, the random top-k setting performs more stable regardless the change of the K during the inference. However, for Manga109 and Set5 dataset. The best PSNR is obtained by the fixed top-k setting (in (d) - (i) in Fig. 12).

In general, based on all the experimental results mentioned above, we conclude that (1) the random top-k setting performs better than the fixed $k$ setting, and usually outperforms the latter by a large margin when $k$ is fixed to small values (*i.e.*, 64, 128, 192, or 256.). (2) For the fixed top-k setting, if $k$ is set to big enough (*i.e.*, 512) during inference, the fixed top-k setting can also achieve comparative performance or even better performance compared to the random top-k setting for several experiments (*e.g.*color JPEG CAR in Fig. 8 (g) and Fig. 8 (k)). However, it is not always possible that the large fixed $k$ setting can be generalized to limited computation resources, and the model trained with large fixed $k$ usually needs the same $k$ for inference to maintain the performance, which leads to heavy computation resources needed even for inference.

To this end, we propose to decouple the way to use $k$ between training and inference. *i.e.*, we can use the random sample $k$ during training while an optional fixed $k$ during inference without degenerating the overall performance. It makes it possible to deploy models that heavily rely on large GPU memory during training but to limited GPU resources while maintaining reliable performance during inference. This is also consistent with the way we implement the proposed attention block (*i.e.*, we adopt a *Torch-Mask* version that requires affordable large GPU memories during training compared to *Torch-Gather* while adopting the *Triton* version during inference) of SemanIR.

In addition, setting a predetermined $k$ value for each patch/pixel enhances computational efficiency. A fixed $k$ value facilitates parallel computation, particularly in attention operations. Conversely, making $k$ values learnable for each patch or pixel would significantly increase the complexity of the attention operation. Nonetheless, exploring the potential of learnable $k$ values for each patch or pixel represents an intriguing avenue for further investigation.

# E More visual Results

To further support the effectiveness of the proposed SemanIR intuitively. We provide more visual comparison in terms of image deblurring, JPEG CAR, image denoising, and image SR below.

**Image Deblurring** The visual results for single image motion deblurring are shown in Fig. 14. As shown in this figure, the proposed method can effectively remove the motion blur in the input images
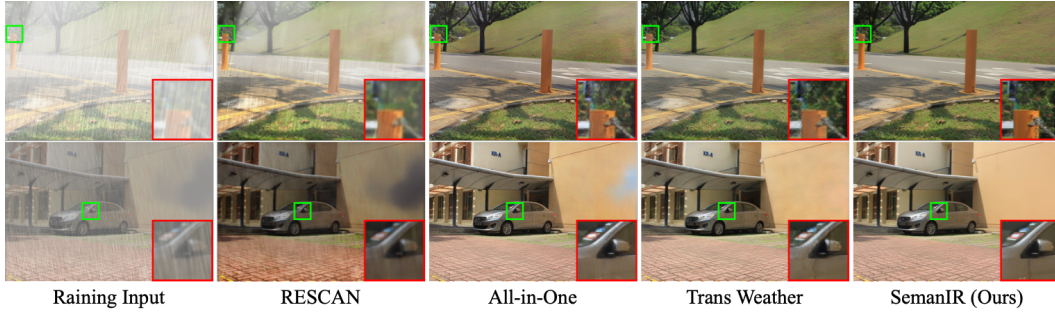
Figure 13: Visual comparison for restoring images in AWC. Best viewed by zooming.

and restore more details such as the facial contour, and the characters compared to MPRNet or Restormer.

**JPEG Compression Artifact Removal** For JPEG compression artifact removal, the visual results for color images on the Urban100 dataset are shown in Fig. 15. The proposed method achieves state-of-the-art performance in removing the blocking artifacts in the input images.

**Image Denoising** The qualitative results for image denoising on the BSD68 and the Urban100 dataset are shown in Fig. 16 (grayscale image) and Fig. 17 (color image). It is clear that for both the grayscale and color inputs, the proposed SemanIR can remove the noise in the noisy input images and recover more realistic textural details in the restored images.

**IR Adverse Weather Conditions.** The qualitative results for IR in AWC on the Test1 [45, 46] dataset are shown in Fig. 13. It shows a challenging case but our method can restore better structural content and clearer details.

**Image SR.** The comparison of visual results of different image SR methods is shown in Fig. 18 and Fig. 19. Fig. 18 shows the results on the Urban100 dataset, and Fig. 19 shows the results on the Manga109 dataset. The proposed SemanIR can restore more missing details in the LR images compared to other state-of-the-art methods like SwinIR, ART, CAT, and EDT.
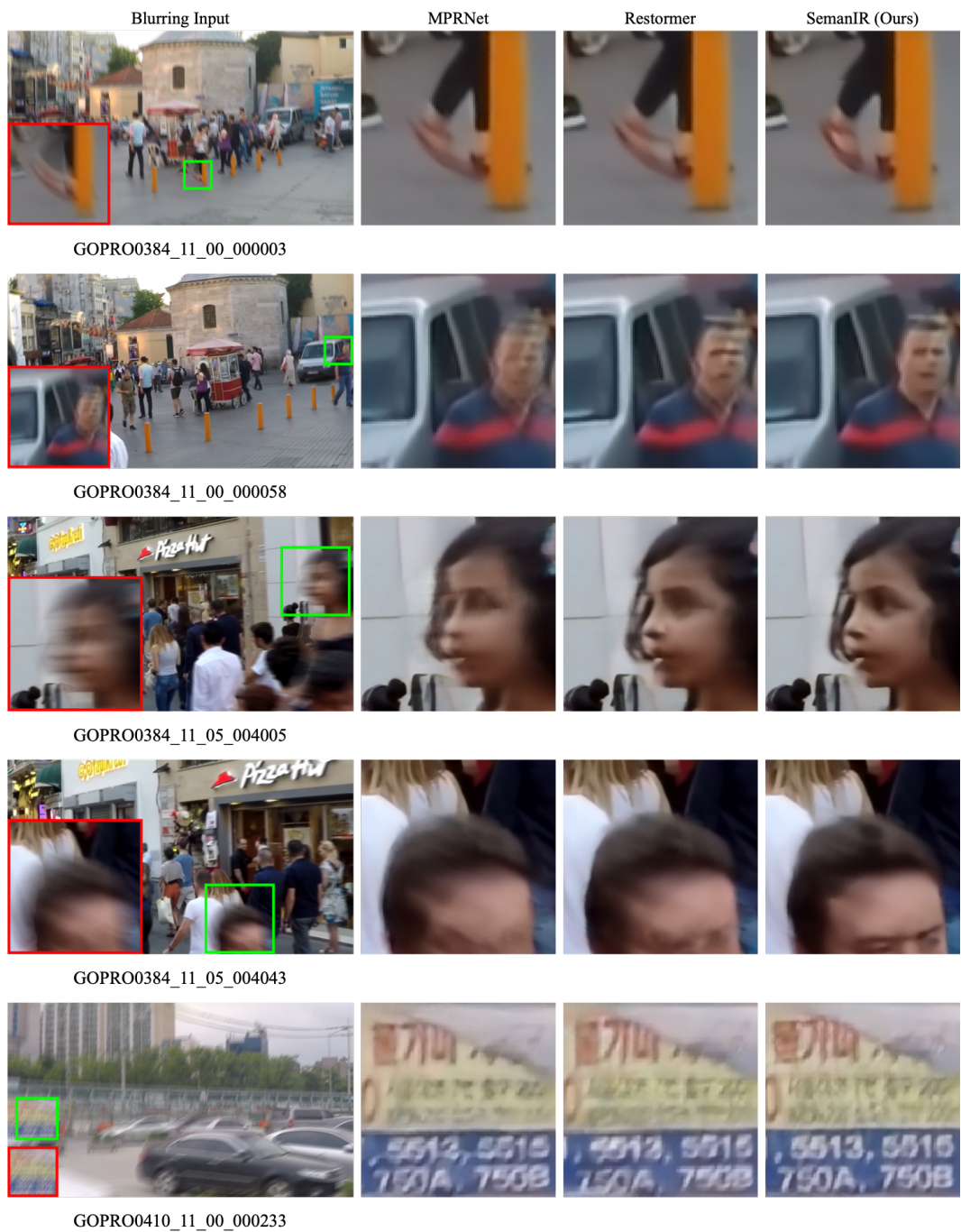
| Blurring Input | MPRNet | Restormer | SemanIR (Ours) |

GOPRO0384_11_00_000003

GOPRO0384_11_00_000058

GOPRO0384_11_05_004005

GOPRO0384_11_05_004043

GOPRO0410_11_00_000233

Figure 14: Visual comparison with single image motion deblurring on GoPro dataset. Best viewed by zooming.

| Input | DRUNet | SwinIR | GRL | SemanIR (Ours) |

Urban100: img_033 (q=10)

Urban100: img_034 (q=10)

Urban100: img_038 (q=10)

Urban100: img_039 (q=10)

Urban100: img_057 (q=10)

Urban100: img_089 (q=10)

Figure 15: Visual comparison of color JPEG CAR on Urban100 dataset. Best viewed by zooming.

BSD68: test001 (GT)    Input (Sigma=50)    SwinIR    Restormer    SemanIR (Ours)

BSD68: test024 (GT)    Input (Sigma=50)    SwinIR

Restormer    SemanIR (Ours)

BSD68: test038 (GT)    Input (Sigma=50)    SwinIR    Restormer    SemanIR (Ours)

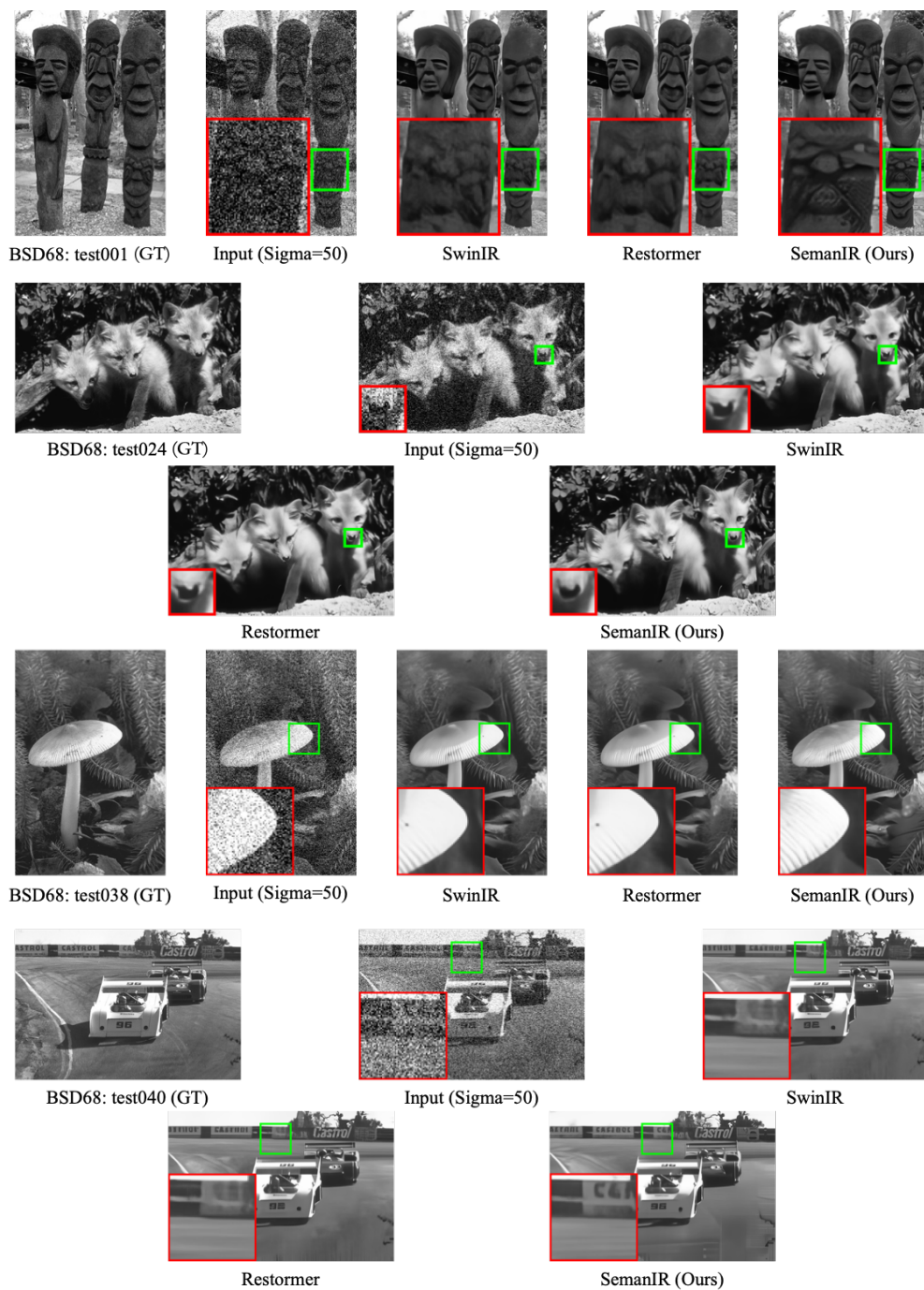BSD68: test040 (GT)    Input (Sigma=50)    SwinIR

Restormer    SemanIR (Ours)

Figure 16: Visual comparison with image denoising on BSD68 dataset. Best viewed by zooming.
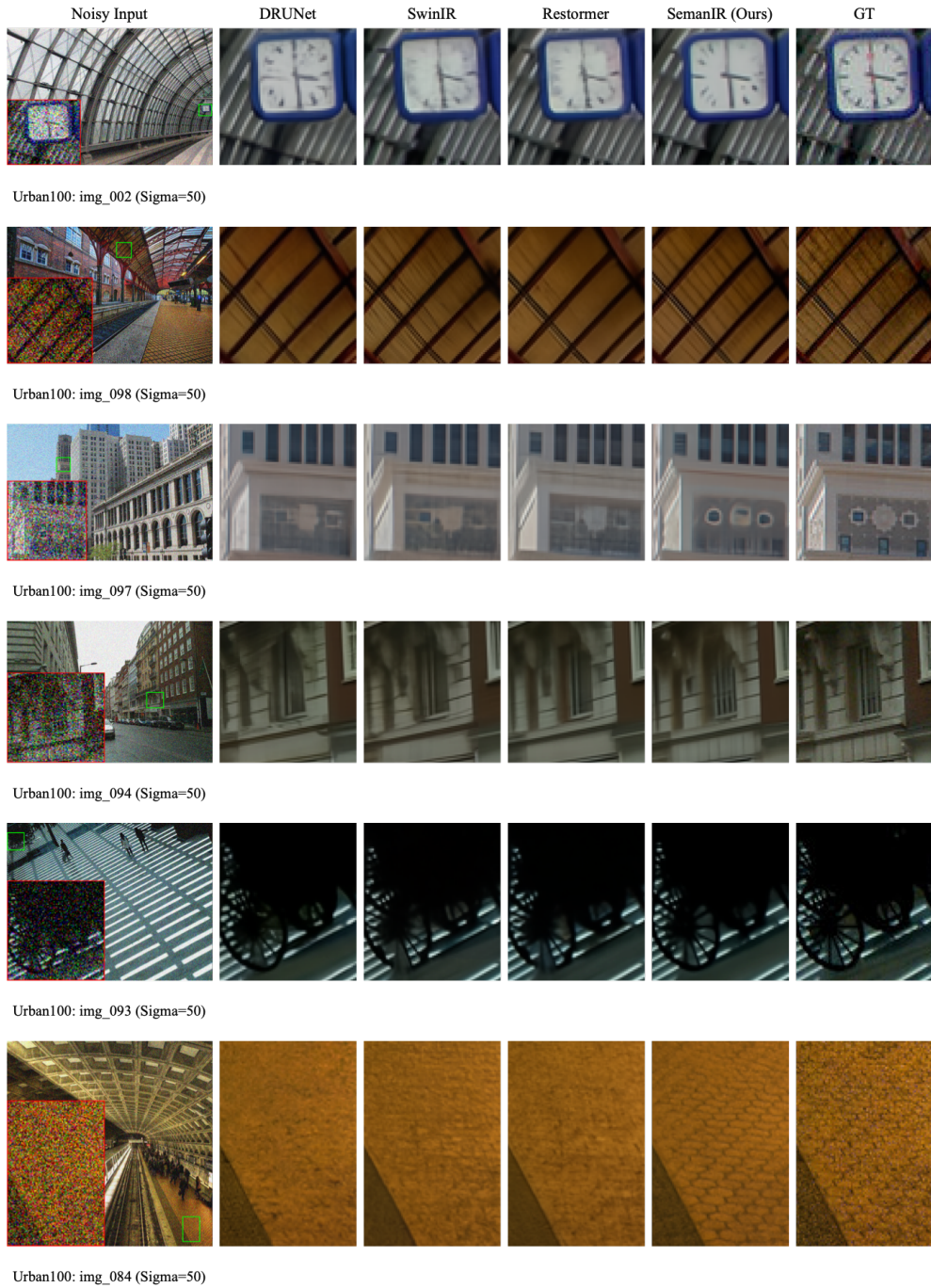
Figure 17: Visual comparison with image denoising on Urban100 dataset. Best viewed by zooming.
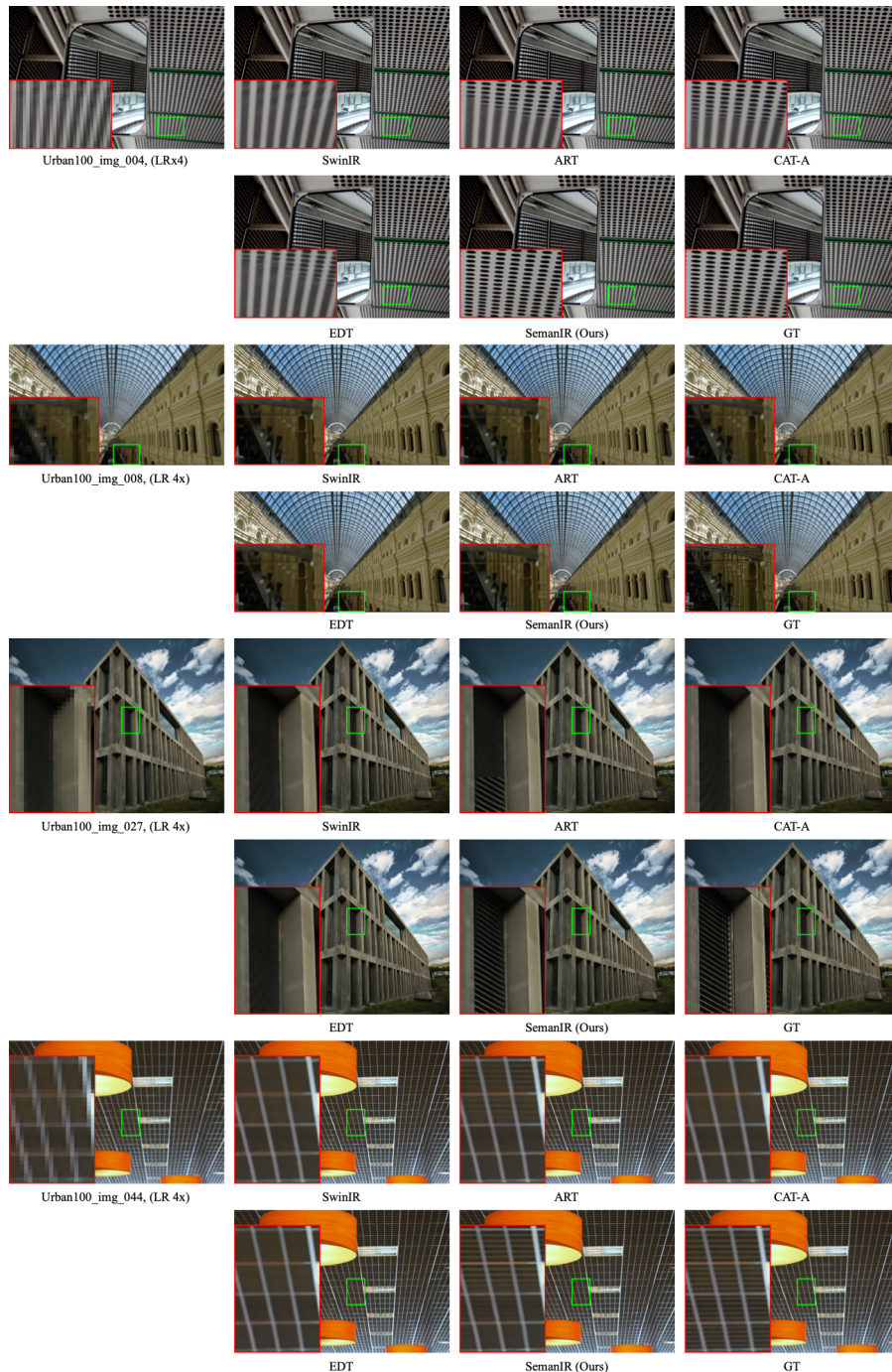
Figure 18: Visual comparison (4×) with image SR on Urban100 dataset. Best viewed by zooming.

Figure 19: Visual comparison (4×) with image SR on Manga109 dataset. Best viewed by zooming.