# Spatially-Attentive Patch-Hierarchical Network
# for Adaptive Motion Deblurring

Maitreya Suin*       Kuldeep Purohit*       A. N. Rajagopalan
Indian Institute of Technology Madras, India
maitreyasuin21@gmail.com, kuldeeppurohit3@gmail.com, raju@ee.iitm.ac.in

## Abstract

*This paper tackles the problem of motion deblurring of dynamic scenes. Although end-to-end fully convolutional designs have recently advanced the state-of-the-art in non-uniform motion deblurring, their performance-complexity trade-off is still sub-optimal. Existing approaches achieve a large receptive field by increasing the number of generic convolution layers and kernel-size, but this comes at the expense of of the increase in model size and inference speed. In this work, we propose an efficient pixel adaptive and feature attentive design for handling large blur variations across different spatial locations and process each test image adaptively. We also propose an effective content-aware global-local filtering module that significantly improves performance by considering not only global dependencies but also by dynamically exploiting neighboring pixel information. We use a patch-hierarchical attentive architecture composed of the above module that implicitly discovers the spatial variations in the blur present in the input image and in turn, performs local and global modulation of intermediate features. Extensive qualitative and quantitative comparisons with prior art on deblurring benchmarks demonstrate that our design offers significant improvements over the state-of-the-art in accuracy as well as speed.*

## 1. Introduction

Motion-blurred images form due to relative motion during sensor exposure and are favored by photographers and artists in many cases for aesthetic purpose, but seldom by computer vision researchers, as many standard vision tools including detectors, trackers, and feature extractors struggle to deal with blur. Blind motion deblurring is an ill-posed problem that aims to recover a sharp image from a given image degraded due to motion-induced smearing of texture and high-frequency details. Due to its diverse applications in surveillance, remote sensing, and cameras mounted
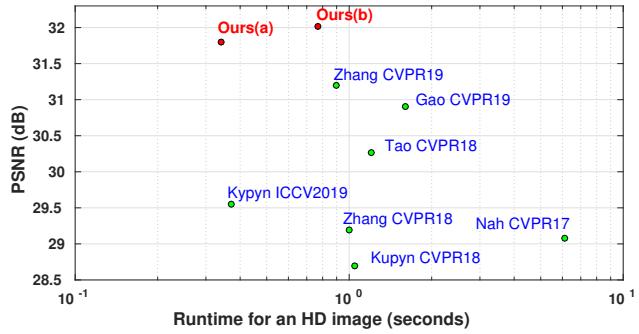


Figure 1. Comparison of different methods in terms of accuracy and inference time. Our approach outperforms all previous methods.

on hand-held and vehicle-mounted cameras, deblurring has gathered substantial attention from computer vision and image processing communities in the past two decades.

Majority of traditional deblurring approaches are based on variational model, whose key component is the regularization term. The restoration quality depends on the selection of the prior, its weight, as well as tuning of other parameters involving highly non-convex optimization setups[14]. Non-uniform blind deblurring for general dynamic scenes is a challenging computer vision problem as blurs arise from various sources including moving objects, camera shake and depth variations, causing different pixels to capture different motion trajectories. Such hand-crafted priors struggle while generalizing across different types of real-world examples, where blur is far more complex than modeled [3].

Recent works based on deep convolutional neural networks (CNN) have studied the benefits of replacing the image formation model with a parametric model that can be trained to emulate the non-linear relationship between blurred-sharp image pairs. Such works [13] directly regress to deblurred image intensities and overcome the limited representative capability of variational methods in describing dynamic scenes. These methods can handle combined effects of camera motion and dynamic object motion and

---

*Equal contribution.

1

achieve state-of-the-art results on single image deblurring task. They have reached a respectable reduction in model size, but still lack in accuracy and are not real-time.

Existing CNN-based methods have two major limitations: a) Weights of the CNN are fixed and spatially invariant which may not be optimal for different pixels in a dynamically blurred scene (e.g., sky vs. moving car pixels). This issue is generally tackled by learning a highly non-linear mapping by stacking a large number of filters. But this drastically increases the computational cost and memory consumption. b) A geometrically uniform receptive field is sub-optimal for the task of deblurring. Large image regions tend to be used to increase the receptive field even though the blur is small. This inevitably leads to a network with a large number of layers and a high computation footprint which slows down the convergence of the network.

Reaching a trade-off between the inference-speed, receptive field and the accuracy of a network is a non-trivial task (see Fig. 1). Our work focuses on the design of efficient and interpretable filtering modules that offer a better accuracy-speed trade-off as compared to simple cascade of convolutional layers. We investigate motion-dependent adaptability within a CNN to directly address the challenges in single image deblurring. Since motion blur is inherently directional and different for each image instance, a deblurring network can benefit from adapting to the blur present in each input test image. We deploy content-aware modules which adjust the filter to be applied and the receptive field at each pixel. Our analysis shows that the benefits of these dynamic modules for the deblurring task are two-fold: i) Cascade of such layers provides a large and dynamically adaptive receptive field. Directional nature of blur requires a directional receptive field, which a normal CNN cannot achieve within a small number of layers. ii) It efficiently enables spatially varying restoration, since changes in filters and features occur according to the blur in the local region. No previous work has investigated incorporating awareness of blur-variation within an end-to-end single image deblurring model.

Following the state of the art in deblurring, we adopt a multi-patch hierarchical design to directly estimate the restored sharp image. Instead of cascading along the depth, we introduce content-aware feature and filter transformation capability through a global-local attentive module and residual attention across layers to improve performance. These modules learn to exploit the similarity in the motion between different pixels within an image and are also sensitive to position-specific local context.

The efficiency of our architecture is demonstrated through a comprehensive evaluation on two benchmarks and comparisons with the state-of-the-art deblurring approaches. Our model achieves superior performance while
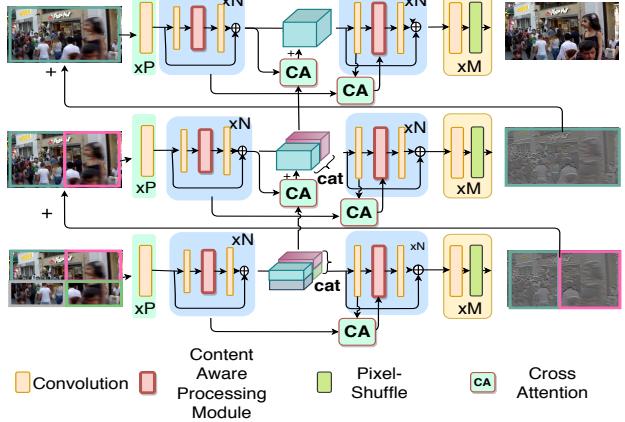


Figure 2. Overall architecture of our proposed network. CA block represents cross attention between different levels of encoder-decoder and different levels. All the resblock contains one content aware processing module. Symbol '+' denotes elementwise summation.

being computationally more efficient. The major contributions of this work are:

- We propose an efficient deblurring design built on new convolutional modules that learn the transformation of features using global attention and adaptive local filters. We show that these two branches complement each other and result in superior deblurring performance. Moreover, the efficient design of attention-module enables us to use it throughout the network without the need for explicit downsampling.

- We further demonstrate the efficacy of learning cross-attention between encode-decoder as well as different levels in our design.

- We provide extensive analysis and evaluations on dynamic scene deblurring benchmarks, demonstrating that our approach yields state-of-the-art results while being $3\times$ faster than the nearest competitor [26].

## 2. Proposed Architecture

To date, the driving force behind performance improvement in deblurring has been the use of a large number of layers and larger filters which assist in increasing the "static" receptive field and the generalization capability of a CNN. However, these techniques offer suboptimal design, since network performance does not always scale with network depth, as the effective receptive field of deep CNNs is much smaller than the theoretical value (investigated in [12]).

We claim that a superior alternative is a dynamic framework wherein the filtering and the receptive field change

across spatial locations and also across different images. Our experiments show that this approach is a considerably better choice due to its task-specific efficacy and utility for computationally limited environments. It delivers consistent performance across diverse magnitudes of blur.

Although previous multi-scale and scale-recurrent methods have shown good performance in removing non-uniform blur, they suffer from expensive inference time and performance bottleneck while simply increasing model depth. Instead, inspired by [26], we adopt multi-patch hierarchical structure as our base-model, which compared to multi-scale approach has the added advantage of residual-like architecture that leads to efficient learning and faster processing speed. The overall architecture of our proposed network is shown in Fig. 2. We divide the network into 3 levels instead of 4 as described in [26]. We found that the relative performance gain due to the inclusion of level 4 is negligible compared to the increase in inference time and number of parameters. At the bottom level input sliced into 4 non-overlapping patches for processing, and as we gradually move towards higher levels, the number of patches decrease and lower level features are adaptively fused using attention module as shown in Fig. 2. The output of level 1 is the final deblurred image. Note that unlike [26], we also avoid cascading of our network along depth, as that adds severe computational burden. Instead, we advocate the use of content-aware processing modules which yield significant performance improvements over even the deepest stacked versions of original DMPHN [26]. Major changes incorporated in our design are described next.

Each level of our network consists of an encoder and a decoder. Both the encoder and the decoder are made of standard convolutional layer and residual blocks where each of these residual blocks contains 1 convolution layer followed by a content-aware processing module and another convolutional layer. The content-aware processing module comprises two branches for global and local level feature processing which are dynamically fused at the end. The residual blocks of decoder and encoder are identical except for the use of cross attention in decoder. We have also designed cross-level attention for effective propagation of lower level features throughout the network. We begin with describing content-aware processing module, then proceed towards the detailed description of the two branches and finally how these branches are adaptively fused at the end.

## 3. Content-Aware Processing Module

In contrast to high-level problems such as classification and detection [22], which can obtain large receptive field by successively down-sampling the feature map with pooling or strided convolution, restoration tasks like deblurring need finer pixel details that can not be achieved from highly downsampled features. Most of the previous deblurring ap-

proaches uses standard convolutional layers for local filtering and stack those layers together to increase the receptive field. [1] uses self-attention and standard convolution on parallel branch and shows that best results are obtained when both features are combined together compared to using each feature separately. Inspired by this approach, we design a content-aware "global-local" processing module which depending on the input, deploys two parallel branches to fuse global and local features. The "global" branch is made of attention module. For decoder, this includes both self and cross-encoder-decoder attention whereas for encoder only self-attention is used. For local branch we design a pixel-dependent filtering module which determines the weight and the local neighbourhood to apply the filter adaptively. We describe in detail these two branches and their adaptive fusion strategy in the following sections.

### 3.1. Attention

Following the recent success of transformer architecture [21] in natural language processing domain, it has been introduced in image processing tasks as well [15, 11]. The main building block of this architecture is self-attention which as the name suggests calculates the response at a position in a sequence by attending to all positions within the same sequence. Given an input tensor of shape $(C, H, W)$ it is flattened to a matrix $z \in \mathbb{R}^{HW \times C}$ and projected to $d_a$ and $d_c$ dimensional spaces using embedding matrices $W_a, W_b \in \mathbb{R}^{C \times d_a}$ and $W_c \in \mathbb{R}^{C \times d_c}$. Embedded matrices $A, B \in \mathbb{R}^{HW \times d_a}$ and $C \in \mathbb{R}^{HW \times d_c}$ are known as query, key and value, respectively. The output of the self-attention mechanism for a single head can be expressed as

$$O = \text{softmax}\left(\frac{AB^T}{\sqrt{d_a}}\right) C \qquad (1)$$

The main drawback of this approach is very high memory requirement due to the matrix multiplication $AB^T$ which requires storing a high dimensional matrix of dimension $(HW, HW)$ for image domain. This requires a large downsampling operation before applying attention. [15] and [17] use a local memory block instead of global all-to-all for making it practically usable. [1] uses attention only from the layer with the smallest spatial dimension until it hits memory constraints. Also, these works typically resort to smaller batch size and sometimes additionally downsampling the inputs to self-attention layers. Although self attention is implemented in recent video super-resolution work [25], to reduce memory requirement it resorts to pixel-shuffling. This process is sub-optimal for spatial attention as pixels are transferred to channel domain to reduce the size.

Different from others, we resort to an attention mechanism which is lightweight and fast. If we consider Eq. (1)
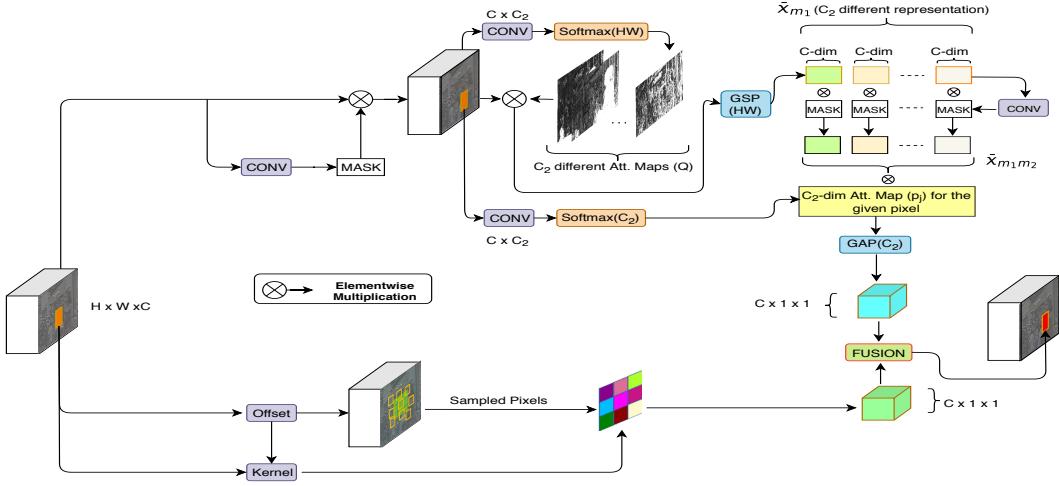
Figure 3. Illustration of our content-aware processing module. The upper and lower branch show self-attention (Sec. 3.1.1) and PDF module (Sec. 3.2). The fusion module is described in Eqs. 12 and 13.

without the softmax and scaling factor for simplicity, we first do a $(HW, d_a) \times (d_a, HW)$ matrix multiplication and then another $(HW, HW) \times (HW, d_c)$ matrix multiplication which is responsible for the high memory requirement and has a complexity of $\mathcal{O}(d_a(HW)^2)$. Instead, if we look into this equation differently and first compute $B^T C$ which is an $(d_a, HW) \times (HW, d_c)$ matrix multiplication followed by $A(B^T C)$ which is an $(HW, d_a) \times (d_a, d_c)$ matrix multiplication, this whole process becomes lightweight with a complexity of $\mathcal{O}(d_a d_c HW)$. We suitably introduce softmax operation at two places which makes this approach intuitively different from standard self-attention but still efficiently gathers global information for each pixel. Empirically we show that it performs better than standard self-attention as discussed in ablation studies. Also due to the light-weight nature, it not only enables us to use this in all the encoder and decoder blocks across levels for self-attention but also across different layers of encoder-decoder and levels for cross attention which results in a significant increase of accuracy.

### 3.1.1 Self-Attention (SA)

We start with generating a spatial attention mask $M_1$ describing which spatial features to emphasize or suppress for better motion understanding. Given the input feature map $x \in \mathbb{R}^{C \times H \times W}$ we generate $M_1$ as

$$M_1 = f_{m_1}(x; \theta_{m_1}) \qquad (2)$$

where $M_1 \in \mathbb{R}^{H \times W}$, $f_{m_1}(\cdot)$ is convolution followed by a sigmoid operation to generate a valid attention map. We generate the enhanced feature map $x_{m_1}$ by element-wise

multiplication as

$$x_{m_1} = x \odot M_1 \qquad (3)$$

where $x_m \in \mathbb{R}^{C \times H \times W}$ and $M$ is broadcast along channel dimension accordingly. Next, we distribute these informative features to all the pixels adaptively which is similar to standard self-attention operation.

Given $x_m$, we generate three attention maps $P \in \mathbb{R}^{C_2 \times HW}$, $Q \in \mathbb{R}^{C_2 \times HW}$ and $M_2 \in \mathbb{R}^C$ using convolutional operations $f_p(\cdot)$, $f_q(\cdot)$ and $f_{M_2}(\cdot)$ where global-average-pooling is used for the last case to get $C$ dimensional representation. We take the first cluster of attention map $Q$ and split it into $C_2$ different maps $Q = \{q_1, q_2, ..., q_{C_2}\}$, $q_i \in \mathbb{R}^{HW}$ and these represent $C_2$ different spatial attention-weights. A single attention reflects one aspect of the blurred image. However, there are multiple pertinent properties like edges, textures etc. in the image that together helps removing the blur. Therefore, we deploy a cluster of attention maps to effectively gather $C_2$ different key features. Each attention map is element-wise multiplied with the input feature map $x_{m_1}$ to generate $C_2$ part feature maps as

$$x_{m_1}^k = q_k \odot x_{m_1} \quad , \text{with} \sum_{i=1}^{HW} q_{ki} = 1 \qquad (k = 1, 2, ..., N) \qquad (4)$$

where $x_m^k \in \mathbb{R}^{C \times HW}$. We further extract descriptive global feature by global-sum-pooling (GSP) along $HW$ dimension to obtain $k^{th}$ feature representation as

$$\bar{x}_{m_1}^k = GSP_{HW}(x_{m_1}^k) \qquad (k = 1, 2, ..., N) \qquad (5)$$

where $\bar{x}_m^k \in \mathbb{R}^C$. Now we have $\bar{x}_{m_1} = \{\bar{x}_{m_1}^1, \bar{x}_{m_1}^2, ..., \bar{x}_{m_1}^{C_2}\}$ which are obtained from $C_2$ different

attention-weighted average of the input $x_m$. Each of these $C_2$ representations is expressed by an $C$-dimensional vector which is a feature descriptor for the $C$ channels. Similar to the first step (Eq.(3)), we further enhance these $C$ dimensional vectors by emphasizing the important feature-embeddings as

$$\bar{x}^k_{m_1 m_2} = M_2 \odot \bar{x}^k_{m_1} \tag{6}$$

where $M_2$ can be expressed as

$$M_2 = f_{m_2}(\bar{x}_{m_1}; \theta_{m_2}) \in \mathbb{R}^C \tag{7}$$

Eq.(3) and Eq.(6) can be intuitively compared to [4], where similar gated-enhancement technique is used to refine the result by elementwise-multiplication with an attention mask that helps in propagating only the relevant information. Next we take the set of attention maps $P = \{p_1, p_2, ..., p_{HW}\}$ where $p_i \in \mathbb{R}^{C_2}$ is represents attention map for $i^{th}$ pixel. Intuitively, $p_i$ shows the relative importance of $C_2$ different attention-weighted average ($\bar{x}_{m_1 m_2}$) for the current pixel and it allows the pixel to adaptively select the weighted average of all the pixels. For each output pixel $j$, we element-wise multiply these $C_2$ feature representations $\bar{x}^k_{m_1 m_2}$ with the corresponding attention map $p_j$, to get

$$y^j = p_j \odot \bar{x}_{m_1 m_2} \text{ with } \sum_{i=1}^{C_2} p_{ji} = 1 \ , (j = 1, 2, ..., HW) \tag{8}$$

where $y^j \in \mathbb{R}^{C \times C_2}$. We again apply global-average-pooling on $y^j$ along $C_2$ to get $C$ dimensional feature representation for each pixel as

$$\bar{y}^j = GAP_{C_2}(y^j) \tag{9}$$

where $\bar{y}^j \in \mathbb{R}^C$ represent the accumulated global feature for the $j^{th}$ pixel. Thus, each pixel flexibly selects features that are complementary to the current one and accumulates a global information. This whole sequence of operations can be expressed by efficient matrix-operations as

$$y^{att} = C \odot \left[(A)\text{softmax}(B)^T\right] \text{softmax}(D) \tag{10}$$

where $A$, $B$, $C$, $D$ are given by

$$C = \sigma(f_{M_2}(x_{m_1})) \in \mathbb{R}^C, A = \sigma(f_{M_1}(x)) \in \mathbb{R}^{C \times HW}$$
$$B = f_Q(x_{m_1}) \in \mathbb{R}^{HW \times C_2}, D = f_P(x_{m_1}) \in \mathbb{R}^{C_2 \times HW}$$

This efficient and simple matrix multiplication makes this attention module very fast whereas the order of operation (first computing $[(A)\text{softmax}(B)^T]$) results in low memory footprint. Note that, $C$ is broadcast along $HW$ dimension appropriately. We utilize this attention block in both encoder and decoder at each level for self-attention.

### 3.1.2 Cross-Attention (CA)

Inspired from the use of cross-attention in [21], we implement cross encoder-decoder and cross level attention in our model. For cross encoder-decoder attention, we deploy similar attention module where the information to be attended is from different encoder layers and all the attention maps are generated by the decoder. Similarly for cross-level, the attended feature is from a lower level and the attention decisions are made by features from a higher level. We have observed that this helps in the propagation of information across layers and levels compared to simply passing the whole input or doing elementwise sum as done in [26].

### 3.2. Pixel-Dependent Filtering Module (PDF)

In contrast to [1], for the local branch, we use Pixel-Dependent Filtering Module to handle spatially-varying dynamic motion blur effectively. Previous works like [6] generate sample-specific parameters on-the-fly using a filter generation network for image classification. [10] uses input text to construct the motion-generating filter weights for video generation task. [28] uses an adaptive convolutional layer where the convolution filter weights are the outputs of a separate filter-manifold network for crowd counting task. Our work is based on [19] as we use a meta-layer to generate pixel dependent spatially varying kernel to implement spatially variant convolution operation. Along with that, the local pixels where the filter is to be applied, are also determined at runtime as we adjust the offsets of these filters adaptively. Given the input feature map $x \in \mathbb{R}^{C \times H \times W}$, we apply a kernel generation function to generate a spatially varying kernel $V$ and do the convolution operation for pixel $j$ as

$$y^{dyn}_{j,c} = \sum_{k=1}^{K} V_{j,j_k} W_c[j_k] x[j + j_k + \Delta j_k] \tag{11}$$

where $y^{dyn}_j \in \mathbb{R}^C$, $K$ is the kernel size, $j_k \in \{(-(K-1)/2, -(K-1)/2), ..., ((K-1)/2, (K-1)/2)\}$ defines position of the convolutional kernel of dilation 1, $V_{j,j_k} \in \mathbb{R}^{K^2 \times H \times W}$ is the pixel dependent kernel generated, $W_c \in \mathbb{R}^{C \times C \times K \times K}$ is the fixed weight and $\Delta j_k$ are the learnable offsets. We set a maximum threshold $\Delta_{\max}$ for the offsets to enforce efficient local processing which is important for low level tasks like deblurring. Note that the kernels ($V$) and offsets vary from one pixel to another, but are constant for all the channels, promoting efficiency. Standard spatial convolution can be seen as a special case of the above with adapting kernel being constant $V_{j,j_k} = 1$ and $\Delta j_k = 0$. In contrast to [1], which simply concatenates the output of these two branches, we design attentive fusion between these two branches so that the network can adaptively adjust the importance of each branch for each pixel at runtime. Empirically we observed that it performs better than simple

| (a) Blurred Image | (b) Blurred patch | (c) MS-CNN | (d) DelurGAN | (e) SRN | (f) DelurGAN-V2 | (g) Stack(4)-DMPHN | (h) Ours (a) |

Figure 4. Visual comparisons of deblurring results on images from the GoPro test set [13]. Key blurred patches are shown in (b), while zoomed-in patches from the deblurred results are shown in (c)-(h).

addition or concatenation. Also, as discussed in visualization section, it gives an insight into the specific requirement for different levels of blur. Given the original input $x$ to this content-aware module, we generate a fusion mask as

$$M_{fus} = sigmoid(f_{fus}(x)) \qquad (12)$$

where $M_{fus} \in \mathbb{R}^{H \times W}$, $f_{fus}$ is a single convolution layer generating single channel output. Then we fuse the two branches as

$$y^{GL} = M_{fus} \odot y^{att} + (1 - M_{fus}) \odot y^{dyn} \qquad (13)$$

The fused output $y^{GL}$ contains global as well as local information distributed adaptively along pixels which helps in handling spatially-varying motion blur effectively.

## 4. Experiments

### 4.1. Implementation Details

**Datasets:** We follow the configuration of [26, 9, 20, 8, 13], which train on 2103 images from the GoPro dataset [13]. For testing, we use two benchmarks: GoPro [13] (1103 HD images), and HIDE [18] (2025 HD images).

**Training settings and implementation details:** All the convolutional layers within our proposed modules contain 128 filters. The hyper-parameters for our encoder-decoder backbone are $N = 3$, $M = 2$, and $P = 2$, and filter size in PDF modules is $5 \times 5$. Following [26], we use batch-size of 6 and patch-size of $256 \times 256$. Adam optimizer [7] was used with initial leaning rate $10^{-4}$, halved after every $2 \times 10^5$ iterations. We use PyTorch [16] library and Titan Xp GPU.

### 4.2. Performance comparisons

The main application of our work is efficient deblurring of general dynamic scenes. Due to the complexity of the blur present in such images, conventional image formation model based deblurring approaches struggle to per-

form well. Hence, we compare with only two conventional methods [23, 24] (which are selected as representative traditional methods for non-uniform deblurring, with publicly available implementations). We provide extensive comparisons with state-of-the-art learning-based methods, namely MS-CNN[13], DeblurGAN[8], DeblurGAN-v2[9], SRN[20], and Stack(4)-DMPHN[26]. We use official implementation from the authors with default parameters.

**Quantitative Evaluation** We show performance comparisons on two different benchmark datasets. The quantitative results on GoPro testing set and HIDE Dataset [18] are listed in Table 1 and 2. We evaluate two variants of our model with(b) and without(a) learnable offsets as shown in Table 1.

The average PSNR and SSIM measures obtained on the GoPro test split is provided in Table 1. It can be observed from the quantitative measures that our method performs better compared to previous state-of-the-art. The results shown in Figure 4. shows the large dynamic blur handling capability of our model while preserving sharpness. We further evaluate the run-time of all the methods on a single GPU with images of resolution $720 \times 1280$. The standard-deviation of the PSNR, SSIM, and run-time scores on the GoPro test set are 1.78, 0.018, and 0.0379, respectively. As reported in Table 1, our method takes significantly less time compared to other methods.

We also evaluate our method on the recent HIDE Dataset [18]. Both of GoPro and HIDE datasets contain dominant foreground object motion along with camera motion. We compare against all existing models trained on GoPro train-set for fair comparisons. As shown in Table 2, our approach outperforms all methods including [18], without requiring any human bounding box supervision. The superiority of our model is owed to the robustness of the proposed adaptive modules.

**Qualitative Evaluation:** Visual comparisons on different dynamic and 3D scenes are shown in Figs. 4 and 5. Visual comparisons are given in Fig. 4. We observe that the

|  (a) Blurred Image | (b) Blurred patch | (c) DelurGAN | (d) SRN | (e) DelurGANv2 | (f) Stack(4)-DMPHN | (g) Ours |

Figure 5. Visual comparisons of deblurring results on images from the HIDE test set [18]. Key blurred patches are shown in (b), while zoomed-in patches from the deblurred results are shown in (c)-(g).

Table 1. Performance comparisons with existing algorithms on 1103 images from the deblurring benchmark GoPro [13].

| Method | [24] | [23] | [5] | [3] | [13] | [8] | [20] | [27] | [2] | [26] | [9] | Ours(a) | Ours(b) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PSNR (dB) | 21 | 24.6 | 23.64 | 26.4 | 29.08 | 28.7 | 30.26 | 29.19 | 30.90 | 31.20 | 29.55 | 31.85 | **32.02** |
| SSIM | 0.741 | 0.846 | 0.824 | 0.863 | 0.914 | 0.858 | 0.934 | 0.931 | 0.935 | 0.940 | 0.934 | 0.948 | **0.953** |
| Time (s) | 3800 | 700 | 3600 | 1200 | 6 | 1 | 1.2 | 1 | 1.0 | 0.98 | 0.48 | **0.34** | 0.77 |



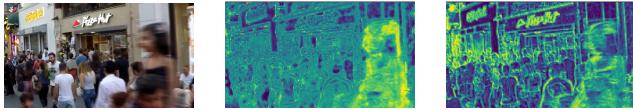(a) Input Image · (b) Fusion $M_{fus}$ · (c) Mask $M_1$

Figure 6. Visualization of intermediate results on images from the GoPro test set [18].

Table 2. Performance comparisons with existing algorithms on 2025 images from the deblurring benchmark HIDE [18].

| Method | [8] | [9] | [20] | [18][1] | [26] | Ours |
|---|---|---|---|---|---|---|
| PSNR | 24.51 | 26.61 | 28.36 | 28.89 | 29.09 | **29.98** |
| SSIM | 0.871 | 0.875 | 0.915 | 0.930 | 0.924 | **0.930** |

Table 3. Quantitative comparison of different ablations of our network on GoPro testset.

| Design | $SA$ | $CA$ | $CLA$ | $Kernel$ | $Offset$ | PSNR |
|---|---|---|---|---|---|---|
| Net1 | ✗ | ✗ | ✗ | ✗ | ✗ | 30.25 |
| Net2 | ✗ | ✗ | ✗ | ✓ | ✗ | 30.81 |
| Net3 | ✓ | ✗ | ✗ | ✗ | ✗ | 30.76 |
| Net4 | ✓ | ✓ | ✗ | ✗ | ✗ | 30.93 |
| Net5 | ✓ | ✗ | ✓ | ✗ | ✗ | 31.12 |
| Net6 | ✓ | ✓ | ✗ | ✓ | ✗ | 31.44 |
| Net7 | ✓ | ✓ | ✓ | ✓ | ✗ | 31.85 |
| Net8 | ✓ | ✓ | ✓ | ✓ | ✓ | 32.02 |

results of prior works suffer from incomplete deblurring or artifacts. In contrast, our network is able to restore scene details more faithfully which are noticeable in the regions containing text, edges, etc. An additional advantage over [5, 23] is that our model waives-off the requirement of parameter tuning during test phase.

On both the datasets, the proposed method achieves consistently better PSNR, SSIM and visual results with lower inference-time than DMPHN [26] and a comparable number of parameters.

### 4.3. Ablation studies

In Table 3, we analyse the effect of individual modules on our network's performance, using 1103 test images from GoPro dataset [13]. As shown in Figure 2, the proposed resblock contains one content-aware processing module and two standard convolutional layers. To find the optimal number of resblock in encoder and decoder we trained different versions of our network with varying number of resblocks. Although, the training performance as well as the quantitative results got better with the increase in number of blocks, beyond 3 the improvement was marginal. This led us to the choice of using 3 resblocks in each encoder and decoder and serves as a good balance between efficiency and performance as well.

As the use of local convolution and global attention together [1] or replacing local convolution with attention [17] is explored recently for image recognition tasks, we further analyze it for image restoration tasks like deblurring. As shown in Table 3, we observe that the advantages of SA and PDF modules are complimentary and their union leads to better performance (Net4 vs Net6). For better information flow between different layers of encoder-decoder and also between different levels we used CA, where the advantage of this attentive information flow rather than simple addition can be observed by comparing the performance of Net4 and Net5 compared to Net3. We also analyze the role of both adaptive weights and the adaptive local-neighborhood for PDF module. As shown quantitatively in Table 3 (Net7 and Net8) and visualized in Figure 7, adaptiveness of the offsets along with the weights perform better as it satisfies the need of directional local filters. We have also showed comparisons of the convergence plots of these models in supplementary. We also try to incorporate the attention mechanism used in [1] in our model for fair comparison. Due
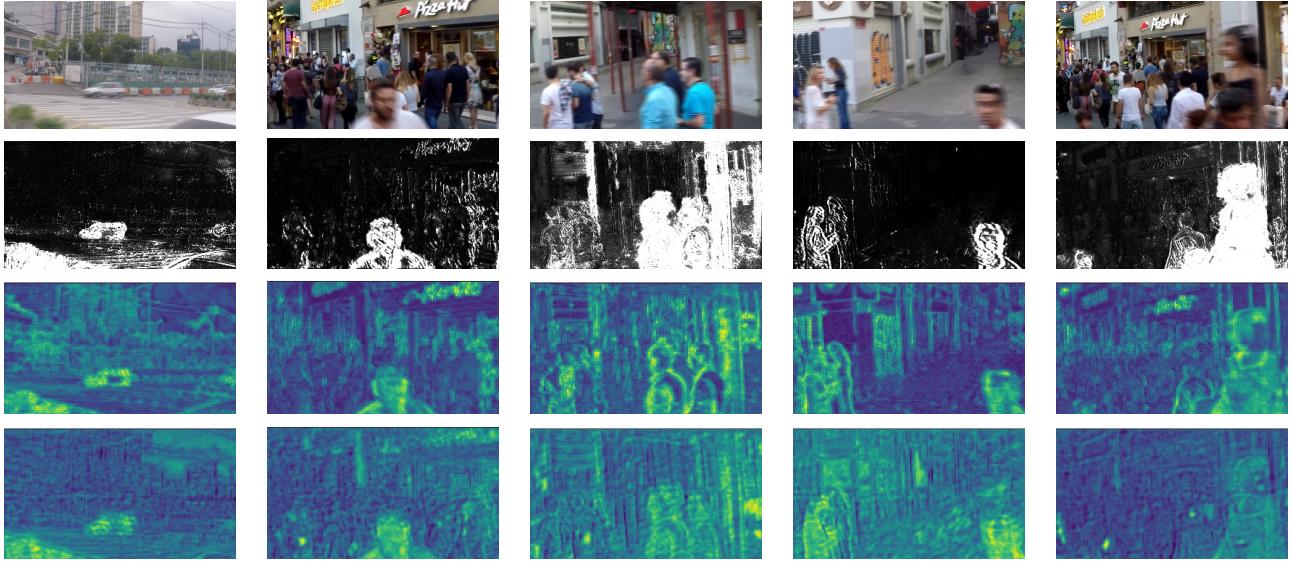
Figure 7. The second row shows one of the spatial attention map for each image. The third row shows the spatial distribution of the horizontal-offset values for the filter. Fourth row shows the variance of the predicted kernel values.

to high memory requirement, we were only able to use one attention module in the decoder in each level. The resultant PSNR was 30.52 compared to 30.76 of Net3. But, as it already occupied full GPU memory, we were unable to introduce more blocks, or cross attention.

### 4.4. Visualization and Analysis

The first row of Fig. 7 contains images from the testing datasets which suffer from complex blur due to large camera and object motion. In the subsequent rows, we visualize the output of different modules of our network and analyze the behavior change while handling different levels of blur due to camera motion, varying depth, moving objects, etc. The second row of Fig. 7 shows one of the attention-maps ($q_i$, $i \in 1, 2, ...C_2$) corresponding to each image. We can observe the high correlation between estimated attention weights and the dominant motion blurred regions present in the image. This adaptive ability of the network to focus on relevant parts of the image can be considered crucial to the observed performance improvement. The third and fourth rows of Fig. 7 show the spatially-varying nature of filter weights and offsets. Observe that a large horizontal offset is estimated in the regions with high horizontal blur so that the filter shape can spread along the direction of motion. Although the estimated filter wights are not directly interpretable, it can be seen that the variance of the filters correlates with the magnitude of blur. We further visualize the behavior of the fusion mask which adaptively weighs the outputs of the two branches for each pixel location. As shown in Fig. 6, PDF module output is more preferred in regions with moving foreground objects or blurred edges where most of the other regions give almost equal weight to both the branches. On the other hand, homogeneous regions

where the effect of blur is negligible, have shown a preference towards the attention branch. To further investigate this behavior, we have visualized the spatial mask ($M_1$). As we can observe in Fig. 6(c), the mask suppresses these homogeneous regions even before calculating self-attention for each pixel. This shows the robustness and interpretability of our attention module while handling any type of blur. **PDF Module:** We synthetically blurred 25 sharp images using synthetic linear PSFs oriented in 4 different directions ($0°,45°,90°,135°$). For these images, we recorded the dominant direction of filter offsets estimated by our PDF module. The values obtained ($11°,50°,81°,126°$) show high correlation between the offset orientations and the PSF angles.

## 5. Conclusions

We proposed a new content-adaptive architecture design for the challenging task of removing spatially-varying blur in images of dynamic scenes. Efficient self-attention is utilized in all the encoder-decoder to get better representation whereas cross-attention helps in efficient feature propagation across layers and levels. Proposed dynamic filtering module shows content-awareness for local filtering. The complimentary behaviour of the two branches are shown in Table 3 and Fig. 6. Different from existing deep learning-based methods for such applications, the proposed method is more interpretable which is one of its key strengths. Our experimental results demonstrated that the proposed method achieved better results than state-of-the-art methods on two benchmarks both qualitatively and quantitatively. We showed that the proposed content-adaptive approach achieves an optimal balance of memory, time and accuracy and can be applied to other image-processing tasks.

# References

[1] Irwan Bello, Barret Zoph, Ashish Vaswani, Jonathon Shlens, and Quoc V Le. Attention augmented convolutional networks. *arXiv preprint arXiv:1904.09925*, 2019.

[2] Hongyun Gao, Xin Tao, Xiaoyong Shen, and Jiaya Jia. Dynamic scene deblurring with parameter selective sharing and nested skip connections. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3848–3856, 2019.

[3] Dong Gong, Jie Yang, Lingqiao Liu, Yanning Zhang, Ian Reid, Chunhua Shen, AVD Hengel, and Qinfeng Shi. From motion blur to motion flow: a deep learning solution for removing heterogeneous motion blur. In *The IEEE conference on computer vision and pattern recognition (CVPR)*, 2017.

[4] Lun Huang, Wenmin Wang, Jie Chen, and Xiao-Yong Wei. Attention on attention for image captioning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4634–4643, 2019.

[5] Tae Hyun Kim, Byeongjoo Ahn, and Kyoung Mu Lee. Dynamic scene deblurring. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3160–3167, 2013.

[6] Xu Jia, Bert De Brabandere, Tinne Tuytelaars, and Luc V Gool. Dynamic filter networks. In *Advances in Neural Information Processing Systems*, pages 667–675, 2016.

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

[8] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiri Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. *arXiv preprint arXiv:1711.07064*, 2017.

[9] Orest Kupyn, Tetiana Martyniuk, Junru Wu, and Zhangyang Wang. Deblurgan-v2: Deblurring (orders-of-magnitude) faster and better. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 8878–8887, 2019.

[10] Yitong Li, Martin Renqiang Min, Dinghan Shen, David Carlson, and Lawrence Carin. Video generation from text. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.

[11] Ding Liu, Bihan Wen, Yuchen Fan, Chen Change Loy, and Thomas S Huang. Non-local recurrent network for image restoration. In *Advances in Neural Information Processing Systems*, pages 1673–1682, 2018.

[12] Wenjie Luo, Yujia Li, Raquel Urtasun, and Richard Zemel. Understanding the effective receptive field in deep convolutional neural networks. In *Advances in neural information processing systems*, pages 4898–4906, 2016.

[13] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, volume 1, page 3, 2017.

[14] TM Nimisha, Akash Kumar Singh, and AN Rajagopalan. Blur-invariant deep learning for blind-deblurring. In *Proceedings of the IEEE E International Conference on Computer Vision (ICCV)*, 2017.

[15] Niki Parmar, Ashish Vaswani, Jakob Uszkoreit, Łukasz Kaiser, Noam Shazeer, Alexander Ku, and Dustin Tran. Image transformer. *arXiv preprint arXiv:1802.05751*, 2018.

[16] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in pytorch. 2017.

[17] Prajit Ramachandran, Niki Parmar, Ashish Vaswani, Irwan Bello, Anselm Levskaya, and Jonathon Shlens. Stand-alone self-attention in vision models. *arXiv preprint arXiv:1906.05909*, 2019.

[18] Ziyi Shen, Wenguan Wang, Xiankai Lu, Jianbing Shen, Haibin Ling, Tingfa Xu, and Ling Shao. Human-aware motion deblurring. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5572–5581, 2019.

[19] Hang Su, Varun Jampani, Deqing Sun, Orazio Gallo, Erik Learned-Miller, and Jan Kautz. Pixel-adaptive convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11166–11175, 2019.

[20] Xin Tao, Hongyun Gao, Xiaoyong Shen, Jue Wang, and Jiaya Jia. Scale-recurrent network for deep image deblurring. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8174–8182, 2018.

[21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pages 5998–6008, 2017.

[22] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7794–7803, 2018.

[23] Oliver Whyte, Josef Sivic, Andrew Zisserman, and Jean Ponce. Non-uniform deblurring for shaken images. *International journal of computer vision*, 98(2):168–186, 2012.

[24] Li Xu, Shicheng Zheng, and Jiaya Jia. Unnatural l0 sparse representation for natural image deblurring. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1107–1114, 2013.

[25] Peng Yi, Zhongyuan Wang, Kui Jiang, Junjun Jiang, and Jiayi Ma. Progressive fusion video super-resolution network via exploiting non-local spatio-temporal correlations. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3106–3115, 2019.

[26] Hongguang Zhang, Yuchao Dai, Hongdong Li, and Piotr Koniusz. Deep stacked hierarchical multi-patch network for image deblurring. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5978–5986, 2019.

[27] Jiawei Zhang, Jinshan Pan, Jimmy Ren, Yibing Song, Linchao Bao, Rynson WH Lau, and Ming-Hsuan Yang. Dynamic scene deblurring using spatially variant recurrent neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2521–2529, 2018.

[28] Lu Zhang, Miaojing Shi, and Qiaobo Chen. Crowd counting via scale-adaptive convolutional neural network. In *2018*

*IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1113–1121. IEEE, 2018.