# Spatio-Temporal Filter Adaptive Network for Video Deblurring

Shangchen Zhou[1*]   Jiawei Zhang[1*]   Jinshan Pan[2†]   Haozhe Xie[1,3]   Wangmeng Zuo[3]   Jimmy Ren[1]

[1]SenseTime Research    [2]Nanjing University of Science and Technology, Nanjing, China
[3]Harbin Institute of Technology, Harbin, China

https://shangchenzhou.com/projects/stfan

## Abstract

*Video deblurring is a challenging task due to the spatially variant blur caused by camera shake, object motions, and depth variations, etc. Existing methods usually estimate optical flow in the blurry video to align consecutive frames or approximate blur kernels. However, they tend to generate artifacts or cannot effectively remove blur when the estimated optical flow is not accurate. To overcome the limitation of separate optical flow estimation, we propose a Spatio-Temporal Filter Adaptive Network (STFAN) for the alignment and deblurring in a unified framework. The proposed STFAN takes both blurry and restored images of the previous frame as well as blurry image of the current frame as input, and dynamically generates the spatially adaptive filters for the alignment and deblurring. We then propose the new Filter Adaptive Convolutional (FAC) layer to align the deblurred features of the previous frame with the current frame and remove the spatially variant blur from the features of the current frame. Finally, we develop a reconstruction network which takes the fusion of two transformed features to restore the clear frames. Both quantitative and qualitative evaluation results on the benchmark datasets and real-world videos demonstrate that the proposed algorithm performs favorably against state-of-the-art methods in terms of accuracy, speed as well as model size.*

## 1. Introduction

Recently, the hand-held and onboard video capturing devices have enjoyed widespread popularity, e.g., smartphone, action camera, unmanned aerial vehicle. The camera shake and high-speed movement in dynamic scenes often generate undesirable blur and result in blurry videos. The low-quality video not only leads to visually poor quality but also hampers some high-level vision tasks such as tracking [12, 21], video stabilization [20] and SLAM [18]. Thus, it is of great interest to develop an effective algorithm to deblur videos for above mentioned human perception and high-level vision tasks.



(a) Blurry frame   (b) SRN [38]   (c) GVD [9]   (d) OVD [10]

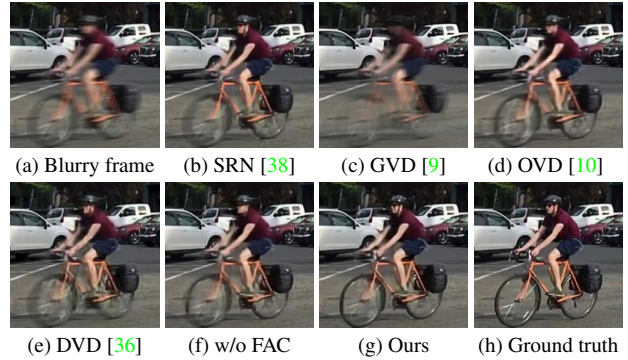(e) DVD [36]   (f) w/o FAC   (g) Ours   (h) Ground truth

Figure 1: One challenging example of video deblurring. Due to the large motion and spatially variant blur, the existing image (b) [38] and video deblurring (c, d, e) [9, 10, 36] methods are less effective. By using the proposed filter adaptive convolutional (FAC) layer for frame alignment and deblurring, our method generates a much clearer image. When the FAC layers are removed (f), our method cannot perform well anymore.

Unlike single-image deblurring, video deblurring methods can exploit additional information that exists across neighboring frames. Significant progress has been made due to the use of sharper regions from neighboring frames [20, 3] or the optical flow from consecutive frames [9, 32]. However, directly utilizing sharp regions of surrounding frames usually generates significant artifacts because the neighboring frames are not fully aligned. Although using the motion field from two adjacent frames, such as optical flow, is able to overcome the alignment problem or approximate the non-uniform blur kernels, the estimation of motion field from blurry adjacent frames is quite challenging.

Motivated by the success of the deep neural networks in low-level vision, several algorithms have been proposed to solve video deblurring [10, 36]. Kim *et al*. [10] concatenate the multi-frame features to restore the current image by a deep recurrent network. However, this method fails to make full use of the information of neighboring frames without explicitly considering alignment, and cannot perform well

---

*Equal contribution   †Corresponding author: sdluran@gmail.com.

when the videos contain large motion. Su *et al.* [36] align the consecutive frames to the reference frame. It shows that this method performs well when the input frames are not too blurry but are less effective for the frames containing severe blur. We also empirically find that both alignment and deblurring are crucial for deep networks to restore sharper frames from blurry videos.

Another group of methods [4, 37, 8, 9] use single or multiple images to estimate optical flow which is treated as the approximation of non-uniform blur kernels. With the estimated optical flow, these methods usually use the existing non-blind deblurring algorithms (*e.g.*, [46]) to reconstruct the sharp images. However, these methods highly depend on the accuracy of the optical flow field. In addition, these methods can only predict line-shaped blur kernel which is inaccurate under some scenarios. To handle non-uniform blur in dynamic scenes, Zhang et al. [44] develop the spatially variant recurrent neural network (RNN) [19] for image deblurring, whose pixel-wise weights are learned from a convolutional neural network (CNN). This algorithm does not need additional non-blind deblurring algorithms. However, it is limited to single image deblurring and cannot be directly extended to video deblurring.

To overcome the above limitations, we propose a Spatio-Temporal Filter Adaptive Network (STFAN) for video deblurring. Motivated by dynamic filter networks [11, 24, 22] which apply the generated filters to the input images, we propose the element-wise filter adaptive convolutional (FAC) layer. Compared with [11, 24, 22], FAC layer applies the generated spatially variant filters on down-sampled features, which allows it to obtain a larger receptive field using a smaller filter size. It also has stronger capability and flexibility due to different filters are dynamically estimated for different channels of the features. The proposed method formulates the alignment and deblurring as two element-wise filter adaptive convolution processes in a unified network. Specifically, given both blurry and restored images of the previous frame and blurry image of the current frame, STFAN dynamically generates corresponding alignment and deblurring filters for feature transformation. In contrast with estimating non-uniform blur kernels from a single blurry image [44, 4, 37, 8] or two adjacent blurry images [9], our method estimates the deblurring filters from a richer inputs: three images and the motion information of two adjacent frames obtained from alignment filters. By using FAC layer, STFAN adaptively aligns the features obtained at different time steps, without explicitly estimating optical flow and warping images, thereby leading to a tolerance of alignment accuracy. In addition, the FAC layers allow our network handle spatially variant blur better, with deblurring in the feature domain. An example in Figure 1 shows that our method generates a much sharper image (Figure 1(g)) than our baseline without FAC layers (Figure 1(f)) as well as the competing methods.

The main contributions are summarized as follows:

- We propose a filter adaptive convolutional (FAC) layer that applies the generated element-wise filters to feature transformation, which is utilized for two spatially variant tasks, i.e. alignment and deblurring in the feature domain.

- We propose a novel spatio-temporal filter adaptive network (STFAN) for video deblurring. It integrates the frame alignment and deblurring into a unified framework without explicit motion estimation and formulates them as two spatially variant convolution process based on the FAC layers.

- We quantitatively and qualitatively evaluate our network on benchmark dataset and show that it performs favorably against state-of-the-art algorithms in terms of accuracy, speed as well as model size.

## 2. Related Work

Our work formulates the neighboring frame alignment and non-uniform blur removal in video deblurring task as two element-wise filter adaptive convolution processes. The following is a review of relevant works on single-image deblurring, multi-image deblurring, and kernel prediction network, respectively.

**Single-Image Deblurring.** Numerous methods have been proposed for single-image deblurring. Early researchers assume a uniform blur kernel and design some natural image priors, such as $L_0$-regularized prior [43], dark channel prior [28], to compensate for the ill-posed blur removal process. However, it is hard for these methods to model spatially-varying blur under dynamic scenes. To model the non-uniform blur, the method [7] and [27] estimate different blur kernels for different segmented the image patches. Other works [4, 37, 8] estimate a dense motion field and a pixel-wise blur kernel.

With the development of deep learning, many CNN-based methods have been proposed to solve dynamic scene deblurring. Method [37] and [4] utilize CNNs to estimate the non-uniform blur kernels. However, the predicted kernels are line-shaped which are inaccurate in some scenarios, and time-consuming conventional non-blind deblurring [46] is generally required to restore the sharp image. More recently, many end-to-end CNN models [38, 44, 17, 23, 26] have also been proposed for image deblurring. To obtain a large receptive field for handling the large blur, the multi-scale strategy is used in [38, 23]. In order to deal with dynamic scene blur, Zhang *et al.* [44] use spatially variant RNNs [19] to remove blur in feature space with a generated RNN weights by a neural network. However, compared with the video-based method, the accuracy of RNN weights is highly limited to having only a single blurry image as input. To reduce the difficulty of restoration and ensures color consistency, Noroozi *et al.* [26] build skip connections between the input and output. The adversarial loss is used in [23, 17] to generate sharper images with more details.
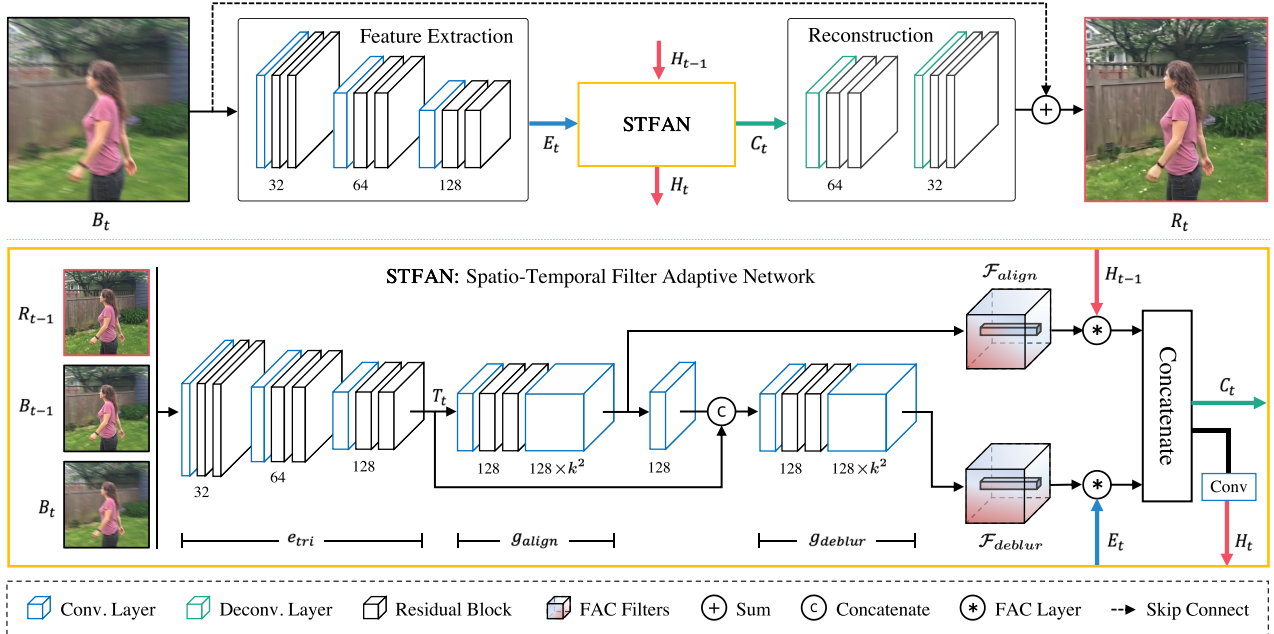
Figure 2: Proposed network structure. It contains three sub-networks: spatio-temporal filter adaptive network (STFAN), feature extraction network, and reconstruction network. Given the triplet images (blurry $B_{t-1}$ and restored $R_{t-1}$ image of the previous frame, and current input image $B_t$), the sub-network STFAN generates the alignment filters $\mathcal{F}_{align}$ and deblurring filters $\mathcal{F}_{deblur}$ in order. Then, using the proposed FAC layer $\circledast$, STFAN aligns deblurred features $H_{t-1}$ of the previous time step with the current time step and removes blur from the features $E_t$ extracted from the current blurry image by the feature extraction network. At last, the reconstruction network is utilized to restore the sharp image from the fused features $C_t$. $k$ denotes the filter size of FAC layer.

**Multi-Image Deblurring.** Many methods utilize multiple images to solve dynamic scene deblurring from video, burst or stereo images. The algorithms by [41] and [32] use the predicted optical flow to segment layers with different blur and estimate the blur layer-by-layer. In addition, Kim *et al.* [9] treat optical flow as a line-shaped approximation of blur kernels, which optimize optical flow and blur kernels iteratively. The stereo-based methods [42, 34, 29] estimate depth from stereo images, which is used to predict the pixel-wise blur kernels. Zhou *et al.* [45] propose a stereo deblurring network with depth awareness and view aggregation. To improve the generalization ability, Chen *et al.* [2] propose an optical flow based reblurring step to reconstruct the blurry input, which is employed to fine-tune deblurring network via self-supervised learning. Recently, several end-to-end CNN methods [36, 10, 15] have been proposed for video deblurring. After image alignment using optical flow, [36] and [15] aggregate information across the neighboring frames to restore the sharp images. Kim *et al.* [10] apply a temporal recurrent network to propagate the features from the previous time step into those of the current one. Despite the fact that motion can be the useful guidance for blur estimation, Aittala *et al.* [1] propose a burst deblurring network in an order-independent manner by repeatedly exchanging the information between the features of the burst images.

**Kernel Prediction Network.** Kernel (filter) prediction network (KPN) has recently witnessed rapid progress in low-level vision tasks. Jia *et al.* [11] first propose the dynamic filter network, which consists of a filter prediction network that predicts kernels conditioned on an input image, and a dynamic filtering layer that applies the generated kernels to another input. Their method shows the effectiveness on video and stereo prediction tasks. Niklaus *et al.* [24] apply kernel prediction network to video frame interpolation, which merges optical flow estimation and frame synthesis into a unified framework. To alleviate the demand for memories, they subsequently propose separable convolution [25] which estimates two separable 1D kernels to approximate 2D kernels. In [22], they utilize KPN for both burst frame alignment and denoising, using the same predicted kernels. [13] reconstructs high-resolution image from low-resolution input using generated dynamic upsampling filters. However, all the above methods directly apply the predicted kernels (filters) in the image domain. In addition, Wang *et al.* [39] propose a spatial feature transform (SFT) layer for image super-resolution. It generates transformation parameters for pixel-wise feature modulation, which can be considered as the KPN with a kernel size of $1 \times 1$ in the feature domain.

3

# 3. Proposed Algorithm

In this section, we first give an overview of our algorithm in Sec. 3.1. Then we introduce the proposed filter adaptive convolutional (FAC) layer in Sec. 3.2. Upon this layer, we show the structure of the proposed networks in Sec. 3.3. Finally, we present the loss functions that are used to constrain the network training in Sec. 3.4.

## 3.1. Overview

Different from the standard CNN-based video deblurring methods [36, 10, 15] that take five or three consecutive blurry frames as input to restore the sharp mid-frame, we propose a frame-recurrent method, which requires information of the previous frame and the current input. Due to the recurrent property, the proposed method is able to explore and utilize the information from a large number of previous frames without increasing the computational demands. As shown in Figure 2, the proposed STFAN generates the filters for alignment and deblurring from the triplet images (blurry and restored image of the previous time step $t-1$, and current input blurry image). Then, using FAC layers, STFAN aligns the deblurred features from the previous time step with the current one and removes blur from the features extracted from the current blurry image. Finally, a reconstruction network is applied to restore the sharp image by fusing the above two transformed features.

## 3.2. Filter Adaptive Convolutional Layer

Motivated by the Kernel Prediction Network (KPN) [11, 24, 22], which applies the generated spatially variant filters to the input image, we propose the filter adaptive convolutional (FAC) layer which applies generated element-wise convolutional filters to the features, as shown in Figure 3. The filters predicted in [11, 24, 22] are the same for RGB channels of each position. To be more capable and flexible for spatially variant tasks, the generated filters for FAC layer are different for each channel. Limited by large memory demand, we only consider the convolution within channels. In theory, the element-wise adaptive filters is five-dimensional ($h \times w \times c \times k \times k$). In practice, the dimension of the generated filter $\mathcal{F}$ is $h \times w \times ck^2$ and we reshape it into the five-dimensional filter. For each position $(x, y, c_i)$ of input feature $Q \in \mathbb{R}^{h \times w \times c}$, a specific local filter $\mathcal{F}_{x,y,c_i} \in \mathbb{R}^{k \times k}$ (reshape from $1 \times 1 \times k^2$) is applied to the region centered around $Q_{x,y,c_i}$ as follows:

$$
\begin{aligned}
\hat{Q}(x, y, c_i) &= \mathcal{F}_{x,y,c_i} * Q_{x,y,c_i} \\
&= \sum_{n=-r}^{r} \sum_{m=-r}^{r} \mathcal{F}(x, y, k^2 c_i + kn + m) \\
&\quad \times Q(x - n, y - m, c_i),
\end{aligned}
\tag{1}
$$

where $r = \frac{k-1}{2}$, $*$ donates convolution operation, $\mathcal{F}$ is the generated filter, $Q(x, y, c_i)$ and $\hat{Q}(x, y, c_i)$ denote the input
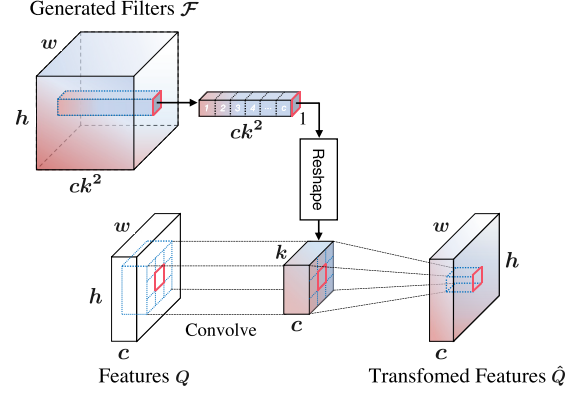


Figure 3: Filter Adaptive Convolutional Layer

features and transformed features, respectively. The proposed FAC layer is trainable and efficient, which is implemented and accelerated by CUDA.

A large receptive field is essential to handle large motions and blurs. The standard KPN methods [11, 24, 22] have to predict the filter much larger in size than motion blur for each pixel of the input image, which requires large computational cost and memory. In contrast, the proposed network does not require a large filter size due to the use of FAC layer on down-sampled features. The experimental results in Table 4 show a small filter size (e.g. 5) on intermediate feature layer is sufficient for deblurring.

## 3.3. Network Architecture

As shown in Figure 2, our network is composed of a spatio-temporal filter adaptive network (STFAN), a feature extraction network, and a reconstruction network.

**Feature Extraction Network.** This network extracts features $E_t$ from the blurry image $B_t$, which consists of three convolutional blocks and each of them has one convolutional layer with stride 2 and two residual blocks [6] with LeakyReLU (negative slope $\lambda = 0.1$) as the activation functions. The extracted features are feed into STFAN for deblurring using FAC layer.

**Spatio-Temporal Filter Adaptive Network.** The proposed STFAN consists of three modules: encoder $e_{tri}$ of triplet images, alignment filter generator $g_{align}$, and deblurring filter generator $g_{deblur}$.

Given the triplet input: the blurry image $B_{t-1}$ and restored image $R_{t-1}$ of the previous frame and the current blurry image $B_t$, STFAN extracts features $T_t$ by the encoder $e_{tri}$. The encoder consists of three convolutional blocks (kernel size 3) and each of them is composed of one convolutional layer with stride 2 and two residual blocks. The alignment filter generator $g_{align}$ takes the extracted features $T_t$ of triplet images as input to predict the adaptive filters for alignment, denoted as $\mathcal{F}_{align} \in \mathbb{R}^{h \times w \times ck^2}$:

$$
\mathcal{F}_{align} = g_{align}(e_{tri}(B_{t-1}, R_{t-1}, B_t)),
\tag{2}
$$

4

where generated $\mathcal{F}_{align}$ contains rich motion information, which is helpful to model the non-uniform blur in the dynamic scene. To make full use of it, the deblurring filter generator $g_{deblur}$ takes alignment filters $\mathcal{F}_{align}$ as well as the features $T$ of triplet images to generate the spatially variant filters for deblurring, denoted as $\mathcal{F}_{deblur} \in \mathbb{R}^{h \times w \times ck^2}$:

$$\mathcal{F}_{deblur} = g_{deblur}(e_{tri}(B_{t-1}, R_{t-1}, B_t), \mathcal{F}_{align}), \quad (3)$$

Both filter generators consist of one convolution layer and two residual blocks with kernel size $3 \times 3$, followed by a $1 \times 1$ convolution layer to expand the channels of output to $ck^2$.

With the two generated filters, two FAC layers are utilized to align the deblurred features $H_{t-1}$ from the previous time step with the current frame and remove the blur from the extracted features $E_t$ of current blurry frame in the feature domain. After that, we concatenate these two transformed features as $C_t$ and restore the sharp image by the reconstruction network. To propagate the deblurred information $H_t$ to the next time step, we pass the features $C_t$ to the next iteration through a convolutional layer.

It is worth noting that both the blurry $B_{t-1}, B_t$ and restored $R_{t-1}$ are required to learn the filters for alignment and deblurring, and thus are taken as the triplet input to STFAN. On the one hand, $B_{t-1}$ and $B_t$ are crucial to capture the motion information across frames and thus benefit alignment. On the other hand, the inclusion of $B_{t-1}$ and $R_{t-1}$ makes it possible to implicitly exploit the blur kernel at frame $t-1$ for improving the deblurring at frame $t$. Moreover, deblurring is assumed to be more difficult but can be benefited by alignment. Thus we stack $g_{deblur}$ upon $g_{align}$ in our implementation. We will analyze the effect of taking triplet images $B_{t-1}, R_{t-1}, B_t$ as input in Sec. 5.3.

**Reconstruction Network.** The reconstruction network is used to restore the sharp images by taking the fusion features from STFAN as input. It consists of scale convolutional blocks, each of which has one deconvolutional layer and two residual blocks as shown in Figure 2.

## 3.4. Loss Function

To effectively train the proposed network, we consider two kinds of loss functions. The first loss is the mean squared error (MSE) loss that measures the differences between the restored frame $R$ and its corresponding sharp ground truth $S$:

$$\mathcal{L}_{mse} = \frac{1}{CHW} ||R - S||^2, \quad (4)$$

where $C, H, W$ are dimensions of image, respectively; $R$ and $S$ respectively denote the restored image and the corresponding ground truth.

To generate more realistic images, we further use the perceptual loss proposed in [14], which is defined as the Euclidean distance between the VGG-19 [35] features of restored frame $R$ and ground truth $S$:

$$\mathcal{L}_{perceptual} = \frac{1}{\mathcal{C}_j \mathcal{H}_j \mathcal{W}_j} ||\Phi_j(R) - \Phi_j(S)||^2, \quad (5)$$

where $\Phi_j(\cdot)$ denotes the features from the $j$-th convolutional layer of the pretrained VGG-19 network and $\mathcal{C}_j, \mathcal{H}_j, \mathcal{W}_j$ are dimensions of features. In this paper, we use the features of conv3-3 ($j = 15$). The final loss function for the proposed network is defined as:

$$\mathcal{L}_{deblur} = \mathcal{L}_{mse} + \lambda \mathcal{L}_{perceptual}, \quad (6)$$

where the weight $\lambda$ is set as 0.01 in our experiments.

## 4. Experiments

### 4.1. Implementation Details

In our experiments, we train the proposed network using the video deblurring dataset from [36]. It contains 71 videos (6,708 blurry-sharp pairs), splitting into 61 training videos (5,708 pairs) and 10 testing videos (1,000 pairs).

**Data Augmentation.** We perform several data augmentations for training. We first divide each video into several sequences with length 20. To add motion diversity into the training data, we reverse the order of sequence randomly. For each sequence, we perform the same image transformations. It consists of chromatic transformations such as brightness, contrast as well as saturation, which are uniformly sampled from [0.8, 1.2] and geometric transformations including randomly flipping horizontally and vertically and randomly cropping to $256 \times 256$ patches. To make our network robust in real-world scenarios, a Gaussian random noise from $\mathcal{N}(0, 0.01)$ is added to the input images.

**Experimental Settings.** We initialize our neural network using the initialization method in [5], and train it using Adam [16] optimizer with $\beta_1 = 0.9$ and $\beta_2 = 0.999$. We set the initial learning rate to $10^{-4}$ and decayed by 0.1 every 400k iterations. The proposed network converges after 900k iterations. We quantitatively and qualitatively evaluate the proposed method on the video deblurring dataset [36]. For a fair comparison with the most related deep learning-based algorithms [23, 17, 44, 38], we finetune all these methods by the corresponding publicly released implementations on the video deblurring dataset [36]. In our experiments, we use both PSNR and SSIM as quantitative evaluation metrics for synthetic testing set. The training code, test model, and experimental results will be available to the public.

### 4.2. Experimental Results

**Quantitative Evaluations.** We compare the proposed algorithm with the state-of-the-art video deblurring methods including conventional optical flow-based pixel-wise kernel estimation [9] and CNN based methods [36, 10]. We also compare it with the state-of-the-art image deblurring methods including conventional non-uniform deblurring [40],

Table 1: Quantitative evaluation on the video deblurring dataset [36], in terms of PSNR, SSIM, running time (sec) and parameter numbers ($\times 10^6$) of different networks. All existing methods are evaluated using their publicly available code. '-' indicates that it is not available.

| Method | Whyte [40] | Sun [37] | Gong [4] | Nah [23] | Kupyn [17] | Zhang [44] | Tao [38] | Kim [9] | Kim [10] | Su [36] | Ours |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Frame# | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 5 | 5 | 2 |
| PSNR | 25.29 | 27.24 | 28.22 | 29.51 | 26.78 | 30.05 | 29.97 | 27.01 | 29.95 | 30.05 | **31.24** |
| SSIM | 0.832 | 0.878 | 0.894 | 0.912 | 0.848 | 0.922 | 0.919 | 0.861 | 0.911 | 0.920 | **0.934** |
| Time (sec) | 700 | 1200 | 1500 | 4.78 | 0.22 | 1.40 | 2.52 | 880 | 0.13 | 6.88 | 0.15 |
| Params (M) | - | 7.26 | 10.29 | 11.71 | 11.38 | 9.22 | 8.06 | - | 0.92 | 16.67 | 5.37 |



(a) Blurry image
PSNR / SSIM

(b) Gong *et al.* [4]
22.72 / 0.7911

(c) Kupyn *et al.* [17]
21.22 / 0.7189

(d) Zhang *et al.* [44]
23.92 / 0.8321

(e) Tao *et al.* [38]
25.29 / 0.8533

(f) Kim and Lee [9]
20.97 / 0.7235

(g) Kim *et al.* [10]
23.21 / 0.8023

(h) Su *et al.* [36]
23.98 / 0.8291

(i) Ours
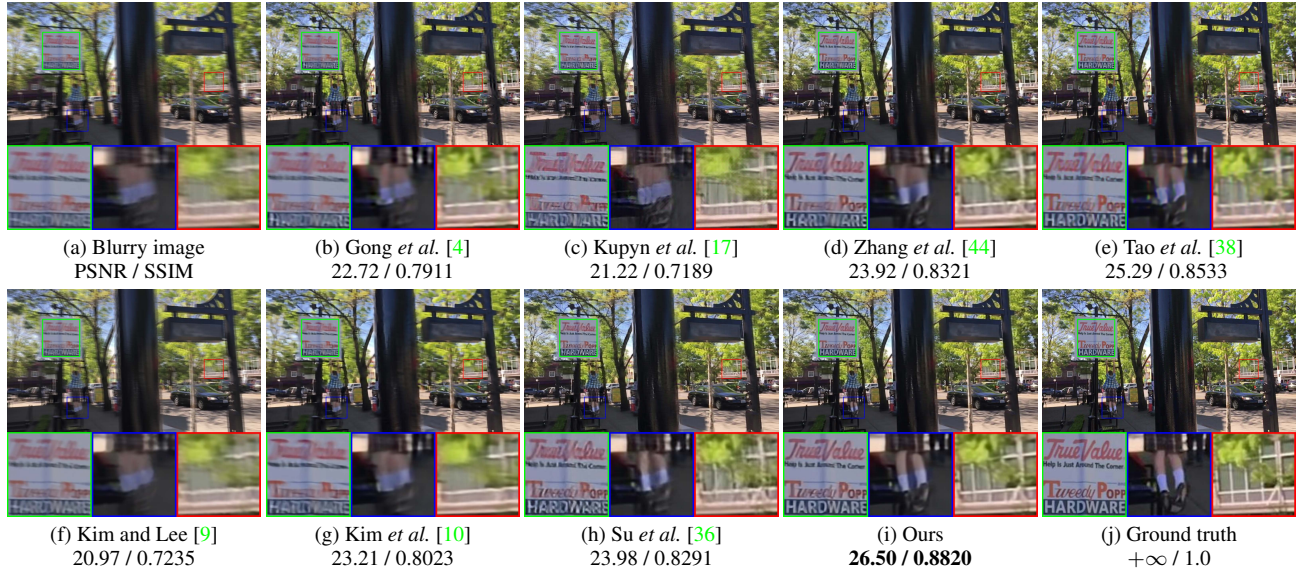**26.50 / 0.8820**

(j) Ground truth
$+\infty$ / 1.0

Figure 4: Qualitative evaluations on Video Deblurring Dataset [36]. The proposed method generates much sharper images with higher PSNR and SSIM.

CNN based spatially variant blur kernel estimation [37, 4], and end-to-end CNN methods [23, 17, 44, 38].

Table 1 shows that the proposed method performs favorably against the state-of-the-art algorithms on the testing set of dynamic scene video deblurring dataset [36].

Figure 4 shows some examples in the testing set from [36]. It shows that the existing methods cannot keep sharp details and remove the non-uniform blur well. With temporal alignment and spatially variant deblurring, our network performs the best and restores much clearer images with more details.

**Qualitative Evaluations.** To further validate the generalization ability of the proposed method, we also qualitatively compare the proposed network with other algorithms on real blurry images from [36]. As illustrated in Figure 5, the proposed method can restore shaper images with more image details than the state-of-the-art image and video deblurring methods. The comparison results show that our STFAN can robustly handle unknown real blur in dynamic scenes, which further demonstrates the superiority of the proposed framework.

### 4.3. Running Time and Model Size

We implement the proposed network using PyTorch platform [30]. To speed up, we implement the proposed FAC layer with CUDA. We evaluate the proposed method and state-of-the-art image or video deblurring methods on the same server with an Intel Xeon E5 CPU and an NVIDIA Titan Xp GPU. The traditional algorithms [40, 9] are time-consuming due to a complex optimization process. Therefore, [37] and [4] utilize the CNN to estimate non-uniform blur kernels based on motion flow. However, they are still time-consuming since the traditional non-blind deblurring algorithm [46] is used to restore the sharp images. DVD [36] uses CNN to restore sharp images from neighboring multiple blurry frames, but they use a traditional optical flow method [31] to align these input frames and is computationally expensive. With GPU implementation, the end-to-end CNN-based methods [23, 17, 44, 38, 10] are relatively efficient. To enlarge the receptive field, the networks in [23, 17, 44, 38] are very deep, which lead to a large model size as well as a long processing time. Even though spatially variant RNNs are used in [44] to enlarge

(a) Blurry image     (b) Gong *et al.* [4]     (c) Nah *et al.* [23]     (d) Kupyn *et al.* [17]     (e) Zhang *et al.* [44]

(f) Tao *et al.* [38]     (g) Kim and Lee [9]     (h) Kim *et al.* [10]     (i) Su *et al.* [36]     (j) Ours
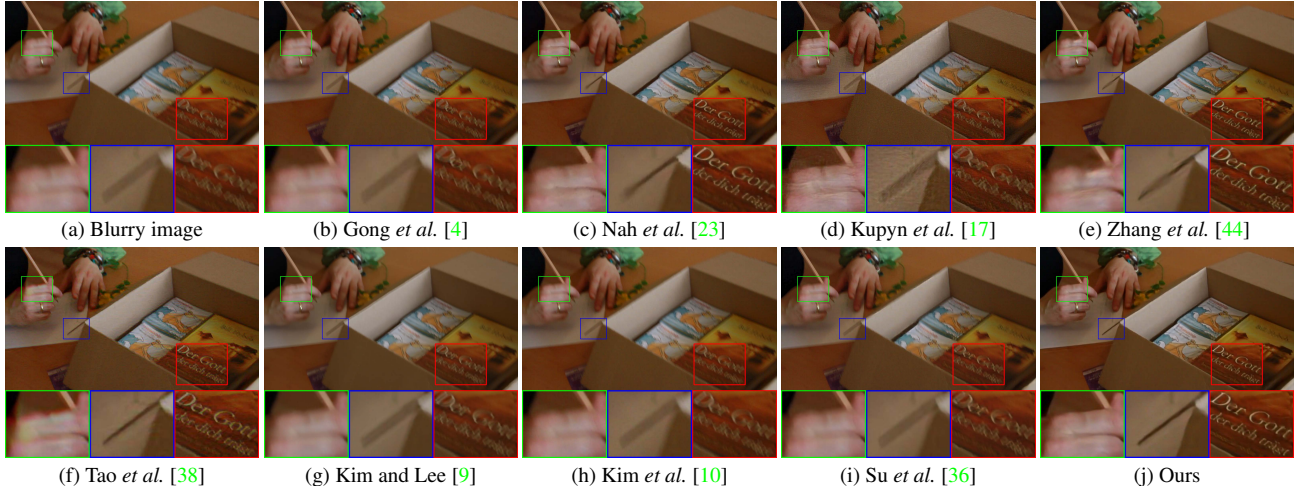
Figure 5: Qualitative evaluations on the real blurry videos [36]. The proposed method generates much clearer images.

the receptive field, they need a deep network to estimate the RNN weights and RNNs are also time-consuming. Our network uses the aligned deblurred features of the previous frame, which reduces the difficulty for the network to restore the sharp image of the current frame. In addition, the FAC layer is effective for spatially variant alignment and deblurring. Benefited from the above two merits, our networks are designed to be small and efficient. As shown in Table 1, the proposed network has less running time and smaller model size than the existing end-to-end CNN methods. Even though [10] runs slightly faster and has smaller model size, the proposed method performs better with the frame alignment and deblurring in the feature domain.

### 4.4. Temporal consistency

To enforce temporal consistency, we adopt the recurrent network to transfer previous feature maps over time, and propose the FAC layer for propagating information between consecutive frames via explicit alignment. Fig. 6 shows that our method not only restores sharper frames but also keeps better temporal consistency. In addition, the video results are given on our [project webpage].

## 5. Analysis and Discussions

We have shown that the proposed algorithm performs favorably against state-of-the-art methods. In this section, we conduct a number of comparative experiments for ablation study and analysis further.

### 5.1. Effectiveness of the FAC layers

The generated alignment filters and deblurring filters are visualized in Figure 7(c) and (h), respectively. According to the optical flow estimated by EpicFlow [33] in Figure 7(b), there is a vehicle moving in the video which is coherent with the alignment filters estimated by our network. Since removing different blur requires different operations and blur
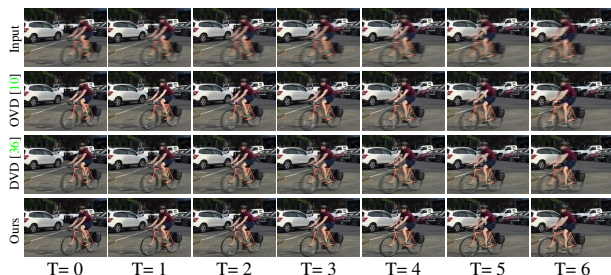


Figure 6: Temporal consistency evaluation on consecutive frames from a blurry video. (**zoom in for best view**).

is somehow related to the optical flow, our network estimates different deblurring filters for foreground vehicle and backgrounds.

To validate the effectiveness of the FAC layer for alignment and deblurring, some intermediate features are shown in Figure 7. According to Figure 7(d) and (i), the FAC layer for alignment can correctly warp the head of the vehicle from green line to purple line even without an image alignment constraint during training. As for the transformed features in Figure 7(j) for deblurring, they are sharper than those before the FAC layer in Figure 7(e), which means the deblurring branch can effectively remove blur in the feature domain.

We also conduct three experiments which replace one or both the FAC layers by concatenating the corresponding features directly, without features transformation by FAC layers. In Table 2, (w/o A, w/ D), (w/ A, w/o D) and (w/o A, w/o D) represent removing FAC layers for feature domain alignment only, feature domain deblurring only and both of them, respectively (refer to Figure 2 for clarification). It shows that the network performs worse without the help of the feature transformation by FAC layers. In addition, Figure 1 also shows that our method cannot restore such a sharp image without using FAC layers.

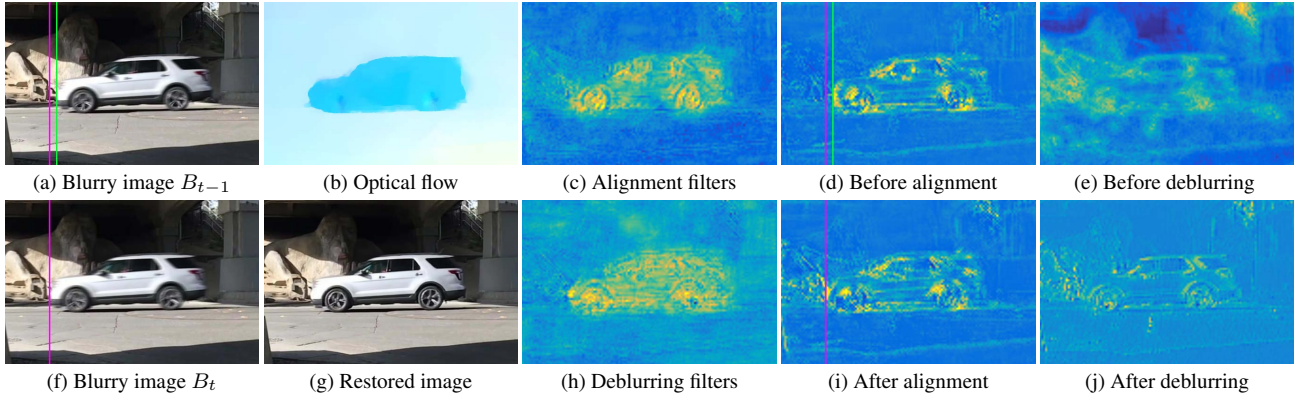| (a) Blurry image $B_{t-1}$ | (b) Optical flow | (c) Alignment filters | (d) Before alignment | (e) Before deblurring |
| (f) Blurry image $B_t$ | (g) Restored image | (h) Deblurring filters | (i) After alignment | (j) After deblurring |

Figure 7: Effectiveness of the adaptive filter generator and FAC layer. (b) is the optical flow from the adjacent input blurry frames (a) and (f) according to EpicFlow [33]. (c) and (h) are the visualization of the generated alignment and deblurring filters of FAC layers, respectively. (d) and (i) are selected feature maps before and after alignment using FAC layer. (e) and (j) are selected feature maps before and after deblurring using FAC layer.

Table 2: Results of different variants of structures. The (w/o A, w/ D), (w/ A, w/o D) and (w/o A, w/o D) represent removing FAC layers for alignment only, deblurring only and both of them, respectively. Unlike the above variants still considering nonalignment features, (-, w D) and (w A, -) denote removing the features of the alignment branch and removing the features of the deblurring branch.

| Structure | w/o A w/o D | w/o A w D | w A w/o D | - w D | w A - | Ours |
|---|---|---|---|---|---|---|
| PSNR | 29.91 | 30.92 | 30.59 | 30.80 | 30.29 | **31.24** |
| SSIM | 0.919 | 0.931 | 0.926 | 0.929 | 0.924 | **0.934** |

## 5.2. Effectiveness of the A and D Branches

To validate the effectiveness of both alignment (A) and deblurring (D) branches, we compare our network with two variant networks: removing the features of the alignment branch (-, w D) and removing the features of the deblurring branch (w A, -). According to Table 2, these two baseline networks do not generate satisfying deblurring results compared to our proposed method.

## 5.3. Effectiveness of the Triplet Input of STFAN

To generate adaptive alignment and deblurring filters, STFAN takes the triplet input (previous blurry image $B_{t-1}$, previous restored image $R_{t-1}$, and current blurry image $B_t$). Table 3 shows the results of two variants which take $(B_{t-1}, B_t)$ and $(R_{t-1}, B_t)$ as input, respectively. The triplet input leads to the best performance. As Sec. 3.3 discussed, the network can implicitly capture the motion and model dynamic scene blur better from the triplet input.

## 5.4. Effectiveness of the Size of Adaptive Filters

To further investigate the proposed network, we test different sizes of adaptive filters, shown in Table 4. The larger size of the adaptive filters leads to better performance. However, increasing the size of adaptive filters after $k = 5$ only

Table 3: Effectiveness of using triplet input of the STFAN. We replace the input of the STFAN by $(B_{t-1}, B_t)$ and $(R_{t-1}, B_t)$ as two variants of our network $(R_{t-1}, B_{t-1}, B_t)$, respectively.

| Input | $(B_{t-1}, B_t)$ | $(R_{t-1}, B_t)$ | $(R_{t-1}, B_{t-1}, B_t)$ |
|---|---|---|---|
| PSNR | 30.87 | 30.85 | **31.24** |
| SSIM | 0.930 | 0.930 | **0.934** |

Table 4: Results of different sizes of adaptive filters.

| Filter Size | $k = 3$ | $k = 5$ | $k = 7$ | $k = 9$ |
|---|---|---|---|---|
| PSNR | 30.95 | 31.24 | 31.27 | 31.30 |
| SSIM | 0.931 | 0.934 | 0.934 | 0.935 |
| Receptive Field | 79 | 87 | 95 | 103 |
| Params (M) | 4.58 | 5.37 | 6.56 | 8.14 |

has minor performance improvement. We empirically set $k = 5$ as a trade-off among the computational complexity, model size and performance.

## 6. Conclusion

We have proposed a novel spatio-temporal network for video deblurring based on filter adaptive convolutional (FAC) layers. The network dynamically generate element-wise alignment and deblurring filters in order. Using the generated filters and FAC layers, our network can perform temporal alignment and deblurring in the feature domain. We have shown that the formulation of two spatially variant problems in video deblurring (i.e., alignment and deblurring) as two filter adaptive convolution processes allows the proposed method to utilize features obtained at different time steps without explicit motion estimation (e.g., optical flow) and enables our method to handle spatially variant blur in dynamic scenes. The experimental results demonstrate the effectiveness of the proposed method in terms of accuracy, speed as well as model size.

# References

[1] M. Aittala and F. Durand. Burst image deblurring using permutation invariant convolutional neural networks. In *ECCV*, 2018.

[2] H. Chen, J. Gu, O. Gallo, M.-Y. Liu, A. Veeraraghavan, and J. Kautz. Reblur2deblur: Deblurring videos via self-supervised learning. In *ICCP*, 2018.

[3] S. Cho, J. Wang, and S. Lee. Video deblurring for hand-held cameras using patch-based synthesis. *TOG*, 31(4):64, 2012.

[4] D. Gong, J. Yang, L. Liu, Y. Zhang, I. D. Reid, C. Shen, A. Van Den Hengel, and Q. Shi. From motion blur to motion flow: A deep learning solution for removing heterogeneous motion blur. In *CVPR*, 2017.

[5] K. He, X. Zhang, S. Ren, and J. Sun. Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *ICCV*, 2015.

[6] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016.

[7] T. Hyun Kim, B. Ahn, and K. Mu Lee. Dynamic scene deblurring. In *ICCV*, 2013.

[8] T. Hyun Kim and K. Mu Lee. Segmentation-free dynamic scene deblurring. In *CVPR*, pages 2766–2773, 2014.

[9] T. Hyun Kim and K. Mu Lee. Generalized video deblurring for dynamic scenes. In *CVPR*, 2015.

[10] T. Hyun Kim, K. Mu Lee, B. Scholkopf, and M. Hirsch. Online video deblurring via dynamic temporal blending network. In *CVPR*, 2017.

[11] X. Jia, B. De Brabandere, T. Tuytelaars, and L. V. Gool. Dynamic filter networks. In *NIPS*, 2016.

[12] H. Jin, P. Favaro, and R. Cipolla. Visual tracking in the presence of motion blur. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, 2005.

[13] Y. Jo, S. Wug Oh, J. Kang, and S. Joo Kim. Deep video super-resolution network using dynamic upsampling filters without explicit motion compensation. In *CVPR*, 2018.

[14] J. Johnson, A. Alahi, and L. Fei-Fei. Perceptual losses for real-time style transfer and super-resolution. In *ECCV*, 2016.

[15] T. H. Kim, M. S. Sajjadi, M. Hirsch, and B. Schölkopf. Spatio-temporal transformer network for video restoration. In *ECCV*, 2018.

[16] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.

[17] O. Kupyn, V. Budzan, M. Mykhailych, D. Mishkin, and J. Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *CVPR*, 2018.

[18] H. S. Lee, J. Kwon, and K. M. Lee. Simultaneous localization, mapping and deblurring. In *ICCV*, 2011.

[19] S. Liu, J. Pan, and M.-H. Yang. Learning recursive filters for low-level vision via a hybrid neural network. In *ECCV*, 2016.

[20] Y. Matsushita, E. Ofek, W. Ge, X. Tang, and H.-Y. Shum. Full-frame video stabilization with motion inpainting. *TPAMI*, 28(7):1150–1163, 2006.

[21] C. Mei and I. Reid. Modeling and generating complex motion blur for real-time tracking. In *CVPR*, 2008.

[22] B. Mildenhall, J. T. Barron, J. Chen, D. Sharlet, R. Ng, and R. Carroll. Burst denoising with kernel prediction networks. In *CVPR*, 2018.

[23] S. Nah, T. H. Kim, and K. M. Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, 2017.

[24] S. Niklaus, L. Mai, and F. Liu. Video frame interpolation via adaptive convolution. In *ICCV*, 2017.

[25] S. Niklaus, L. Mai, and F. Liu. Video frame interpolation via adaptive separable convolution. In *CVPR*, 2017.

[26] M. Noroozi, P. Chandramouli, and P. Favaro. Motion deblurring in the wild. In *GCPR*, 2017.

[27] J. Pan, Z. Hu, Z. Su, H.-Y. Lee, and M.-H. Yang. Soft-segmentation guided object motion deblurring. In *CVPR*, 2016.

[28] J. Pan, D. Sun, H. Pfister, and M.-H. Yang. Blind image deblurring using dark channel prior. In *CVPR*, 2016.

[29] L. Pan, Y. Dai, M. Liu, and F. Porikli. Simultaneous stereo video deblurring and scene flow estimation. In *CVPR*, 2017.

[30] A. Paszke, S. Gross, S. Chintala, G. Chanan, E. Yang, Z. DeVito, Z. Lin, A. Desmaison, L. Antiga, and A. Lerer. Automatic differentiation in pytorch. In *NIPS Workshops*, 2017.

[31] J. S. Pérez, E. Meinhardt-Llopis, and G. Facciolo. Tv-l1 optical flow estimation. *Image Processing On Line*, 2013:137–150, 2013.

[32] W. Ren, J. Pan, X. Cao, and M.-H. Yang. Video deblurring via semantic segmentation and pixel-wise non-linear kernel. In *ICCV*, 2017.

[33] J. Revaud, P. Weinzaepfel, Z. Harchaoui, and C. Schmid. Epicflow: Edge-preserving interpolation of correspondences for optical flow. In *CVPR*, 2015.

[34] A. Sellent, C. Rother, and S. Roth. Stereo video deblurring. In *ECCV*, 2016.

[35] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[36] S. Su, M. Delbracio, J. Wang, G. Sapiro, W. Heidrich, and O. Wang. Deep video deblurring for hand-held cameras. In *CVPR*, 2017.

[37] J. Sun, W. Cao, Z. Xu, and J. Ponce. Learning a convolutional neural network for non-uniform motion blur removal. In *CVPR*, 2015.

[38] X. Tao, H. Gao, X. Shen, J. Wang, and J. Jia. Scale-recurrent network for deep image deblurring. In *CVPR*, 2018.

[39] X. Wang, K. Yu, C. Dong, and C. Change Loy. Recovering realistic texture in image super-resolution by deep spatial feature transform. In *CVPR*, 2018.

[40] O. Whyte, J. Sivic, A. Zisserman, and J. Ponce. Non-uniform deblurring for shaken images. *IJCV*, 98(2):168–186, 2012.

[41] J. Wulff and M. J. Black. Modeling blurred video with layers. In *ECCV*, 2014.

[42] L. Xu and J. Jia. Depth-aware motion deblurring. In *ICCP*, 2012.

[43] L. Xu, S. Zheng, and J. Jia. Unnatural l0 sparse representation for natural image deblurring. In *CVPR*, 2013.

[44] J. Zhang, J. Pan, J. Ren, Y. Song, L. Bao, R. W. Lau, and M.-H. Yang. Dynamic scene deblurring using spatially variant recurrent neural networks. In *CVPR*, 2018.

[45] S. Zhou, J. Zhang, W. Zuo, H. Xie, J. Pan, and J. S. Ren. Davanet: Stereo deblurring with view aggregation. In *CVPR*, 2019.

[46] D. Zoran and Y. Weiss. From learning models of natural image patches to whole image restoration. In *ICCV*, 2011.