

Bringing Events into Video Deblurring with Non-consecutively Blurry Frames

Supplementary Material

In this supplementary file, we provide the details of network architecture, more results of ablation study, more results on benchmark datasets and real-world blurry videos.

1. Network Architecture

1.1. Architectures of D²Nets

D²Nets includes BiLSTM Detector in \mathcal{F}_{DET} to detect blurry frames in a video, \mathcal{F}_{BRN} to restore blurry frame guided by detected NSFs and \mathcal{F}_{TCE} to enhance the temporal consistency of restored video. EFM is proposed to better utilize rich boundaries in events for facilitating deblurring, which can be embedded in \mathcal{F}_{BRN} and \mathcal{F}_{TCE} . We first formulate EFM, and then detail the network architecture.

1.1.1 Event Fusion Module (EFM)

To reduce the computational cost of EFM, we suggest to incorporate EFM into the bottleneck of encoder-decoder, which can be formulated as

$$\begin{aligned}
 z &= \mathcal{F}_{\text{Encoder}}(\mathbf{I}_i^+, \mathbf{B}_i, \mathbf{I}_i^-), \\
 e &= \mathcal{F}_{\text{Event}}(\mathbf{E}_i), \\
 \mathbf{m} &= \text{SoftMax} \left((\mathcal{R}(\mathbf{W}_e e))^T \otimes \mathcal{R}(\mathbf{W}_z z) \right), \\
 \mathbf{y} &= z + \mathcal{R}'(\mathcal{R}(\mathbf{W}_v z) \otimes \mathbf{m}),
 \end{aligned} \tag{1}$$

where \otimes is matrix multiplication, $\mathcal{F}_{\text{Encoder}}$ is convolutional Encoder, $\mathcal{F}_{\text{Event}}$ is convolutional module that can guarantee e has the same spatial dimension with z , \mathcal{R} is the reshape operation that changes dimension from (c, h, w) to $(c, h * w)$, and \mathcal{R}' is the reshape operation that changes dimension from $(c, h * w)$ to (c, h, w) . z is features of frames from latent space of encoder, and e is the feature of events \mathbf{E}_i extracted using CNN along with downsampling. \mathbf{m} is a reweighting map which can be used to facilitate deblurring by matrix multiplication with the features of frames. \mathbf{W}_e , \mathbf{W}_z and \mathbf{W}_v are learnable weight matrices. Finally, \mathbf{y} is the input of decoder to recover latent sharp frame.

1.1.2 Architecture details

D²Nets includes three parts: \mathcal{F}_{DET} , \mathcal{F}_{BRN} and \mathcal{F}_{TCE} . Since \mathcal{F}_{TCE} shares the same architecture with \mathcal{F}_{BRN} , we only present architecture details of \mathcal{F}_{BRN} . BiLSTM detector in \mathcal{F}_{DET} is detailed in Table s1, \mathcal{F}_{BRN} is detailed in Table s3. $\mathcal{F}_{\text{Event}}$ is detailed in Table s4, which can be embedded into the bottleneck of \mathcal{F}_{BRN} or \mathcal{F}_{TCE} .

Table s1. The architecture of BiLSTM detector in \mathcal{F}_{DET} . Convolution is with the form Conv(input channel, output channel, kernel size, stride, padding size) and fully connected layer is with the form Linear(input channel, output channel).

<i>Input</i>	A sequence of frames (indim = 3) or events (indim = 20);
<i>Layer f_{in}</i>	Conv(indim,3, 3, 1, 1); LeakyReLU;
<i>Layer $f_{\text{extractor}}$</i>	ResNet-152;
<i>Layer f_{c1}</i>	Linear(2048, 512);
<i>Layer f_{rec}</i>	BiLSTM;
<i>Layer f_{c2}</i>	Linear(1024, 1); Sigmoid;
<i>Output</i>	A sequence of probabilities of inputs being sharp frames;

Table s2. The basic component ResBlock. Convolution is with the form Conv(input channel, output channel, kernel size, stride, padding size). $feats$ is the number of channels from previous layer.

<i>Input</i>	Feature maps from the previous layer;
<i>Layer 1</i>	Conv($feats, feats, 3, 1, 1$); LeakyReLU;
<i>Layer 2</i>	Conv($feats, feats, 3, 1, 1$); Add(<i>Input</i> , <i>Layer2</i>);

Table s3. The architecture of f_{rec} in \mathcal{F}_{BRN} , while we omit the architecture of f_{flow} due to its exactly same architecture with PWC-Net [5]. Convolution is with the form Conv(input channel, output channel, kernel size, stride, padding size) and deconvolution is with the form ConvTrans(input channel, output channel, kernel size, stride, padding size).

Input: Front NFS G_i^- , blurry frame B_i , near NFS G_i^+	
Frames aligning: Warping NFSs aligned to blurry frame B_i using f_{flow} .	
Output: the warped frames I_i^- and I_i^+	
Concatenate(I_i^-, B_i, I_i^+)	
Encoder:	
<i>Layer f_{in}</i>	Conv(9, 32, 3, 1, 1); LeakyReLU; 3 ResBlocks;
<i>Layer f_{en1}</i>	Conv(32, 64, 3, 2, 1); LeakyReLU; 3 ResBlocks;
<i>Layer f_{en2}</i>	Conv(64, 128, 3, 2, 1); LeakyReLU; 3 ResBlocks;
Bottleneck:	
<i>Layer f_{rec}</i>	<i>LSTM</i> ;
Decoder:	
<i>Layer f_{de2}</i>	3 ResBlocks; ConvTrans(128, 64, 3, 2, 1); LeakyReLU;
<i>Layer f_{de1}</i>	3 ResBlocks; ConvTrans(64, 32, 3, 2, 1); LeakyReLU;
<i>Layer f_{out}</i>	3 ResBlocks; Conv(32, 3, 3, 1, 1);
Output: The latent sharp frame \hat{I}_i	

Table s4. The architecture of \mathcal{F}_{Event} . Convolution is with the form Con.(input channel, output channel, kernel size, stride, padding size). UnPS(factor) is a reverse operation of pixel-shuffle to perform downsampling, and we use factor = 2 to guarantee that the feature size is same with z .

<i>Input</i>	Events E_i ;
<i>Layer $f_{event.in1}$</i>	Conv(20,32,1,1,0); LeakyReLU;
<i>Layer $f_{event.in2}$</i>	Conv(32,32,3,2,1); LeakyReLU;
<i>Layer f_{ups}</i>	UnPS(2);
<i>Output</i>	Events feature e ;

2. Ablation Study

2.1. Effects of Components in D²Nets

Here, we present the deblurring results by D²Nets and its variants. One can see that full D²Nets* obtain visually favorable results as shown in Fig. s1.



Figure s1. Visual comparison of component analysis on GoPro dataset. The first row is blurry frames, and (b)~(e) correspond to the results of 1 ~ 4 rows of Table 4 in the manuscript.

2.2. Effectiveness of Event Fusion Module

In Fig. s2, we visualize the feature map in bottleneck with and without our EFM, respectively.

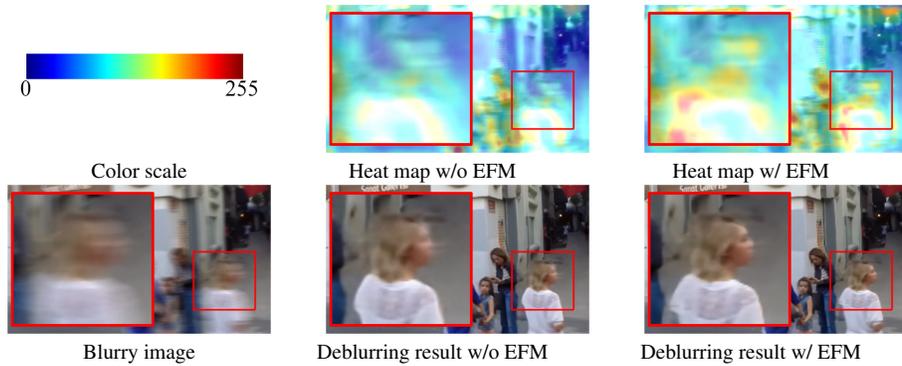


Figure s2. Effect of EFM. The feature map reweighted by our EFM pays more attention on blurry boundaries, which benefits the sharp frame restoration.

3. More Results on Benchmark Datasets

3.1. More Results on GoPro Dataset [2]

The results by DMPHN [6] and CDVD-TSP [3] cannot fully remove severe blur, while our results are more visually plausible.



Figure s3. More results on GoPro dataset. From 1 to 4 rows are blurry frames, DMPHN [6], CDVD-TSP [3] and D²Net*, respectively.

3.2. More Results on Blur-DVS Dataset [1]

The results by DMPHN [6] and CDVD-TSP [3] still have severe blur.



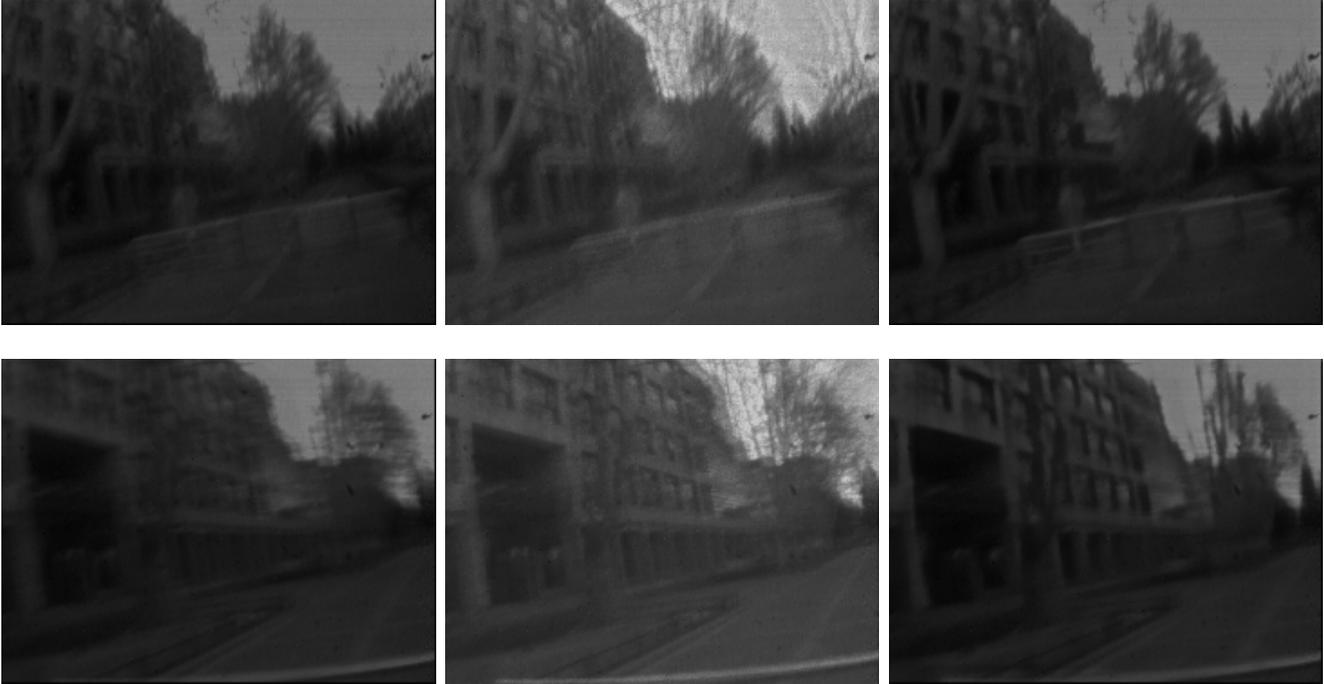
Figure s4. More results comparison on Blur-DVS dataset. From 1 to 4 rows are Blurry frames, DMPHN [6], CDVD-TSP [3] and D²Net*, respectively.

3.3. More Results on Severe Blurry Frames

We present more results on severe blurry frames with events from fast-motion subset of Blur-DVS dataset.

The left is blurry frames, the middle is the results by BHA [4], and the right is the results by our D²Nets*





References

- [1] Zhe Jiang, Yu Zhang, Dongqing Zou, Jimmy Ren, Jiancheng Lv, and Yebin Liu. Learning event-based motion deblurring. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3320–3329, 2020. 5
- [2] Seungjun Nah, Tae Hyun Kim, and Kyoung Mu Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3883–3891, 2017. 4
- [3] Jinshan Pan, Haoran Bai, and Jinhui Tang. Cascaded deep video deblurring using temporal sharpness prior. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3043–3051, 2020. 4, 5
- [4] Liyuan Pan, Cedric Scheerlinck, Xin Yu, Richard Hartley, Miaomiao Liu, and Yuchao Dai. Bringing a blurry frame alive at high frame-rate with an event camera. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 6820–6829, 2019. 6
- [5] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *IEEE Conference on computer vision and pattern recognition*, pages 8934–8943, 2018. 2
- [6] Hongguang Zhang, Yuchao Dai, Hongdong Li, and Piotr Koniusz. Deep stacked hierarchical multi-patch network for image deblurring. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5978–5986, 2019. 4, 5