

# Joint Video Multi-Frame Interpolation and Deblurring under Unknown Exposure Time *Supplementary Material*

Wei Shang<sup>1</sup>, Dongwei Ren<sup>1\*</sup>, Yi Yang<sup>1</sup>, Hongzhi Zhang<sup>1</sup>, Kede Ma<sup>2</sup>, Wangmeng Zuo<sup>1,3</sup>

<sup>1</sup>School of Computer Science and Technology, Harbin Institute of Technology

<sup>2</sup>City University of Hong Kong <sup>3</sup>Peng Cheng Laboratory, Shenzhen

In this supplementary file, we provide details of the network architecture of the proposed VIDUE method and more results on benchmark datasets. Please refer to the link for more visual results at <https://onedrive.live.com>.

## 1. Network Architecture

As mentioned in the main manuscript, the proposed VIDUE consists of an exposure-aware feature extractor  $g_e$ , an intra- and inter-motion analyzer  $g_a$ , and a video reconstruction network  $f$ .

### 1.1. Exposure-Aware Feature Extractor

We illustrate the network structure of the exposure-aware feature extractor  $g_e$  in Table S1, in which ResBlock is detailed in Table S2.

### 1.2. Intra- and Inter-Motion Analyzer

We choose to analyze both intra-motion within each video frame and inter-motion between frames. Our motion analyzer  $g_a : \mathbb{R}^{(T \times 3) \times H \times W} \mapsto \mathbb{R}^{T \times 1 \times H \times W} \times \mathbb{R}^{(S \times T) \times 2 \times H \times W}$ , computes, from an input video sequence,  $T$  intra-motion maps and  $(S \times T) \times 2$  inter-motion maps of the same size, respectively. We show the network structure for intra-motion analysis in Table S3. The network structure for inter-motion analysis is similar to RefineNet proposed in [3], which is detailed in Table S4.

### 1.3. Video Reconstruction Network

The video reconstruction network  $f : \mathbb{R}^{T \times 3 \times H \times W} \mapsto \mathbb{R}^{(S \times T) \times 3 \times H \times W}$  is built upon the exposure-aware representation  $u$  and motion-aware representations  $m$  and  $n$  by progressive exposure-adaptive convolution and motion refinement. The network structure is illustrated in Table S5.

## 2. More Results on Benchmark Datasets

In this section, we provide more visual results of VIDUE in comparison to existing video interpolation and deblurring methods.

### 2.1. More Results on the GoPro Dataset

On the GoPro dataset with mild blur, VIDUE generally achieves more visually appealing results with sharper boundaries and fewer artifacts, as shown in Figs. S1 and S2.

### 2.2. More Results on the Adobe Dataset

On the Adobe dataset with heavy blur, VIDUE still performs favorably as shown in Figs. S3 and S4.

### 2.3. More Results on the RealBlur Dataset

On the RealBlur dataset, VIDUE generates the least artifacts and sharper textures compared with other methods, as shown in Figs. S5 and S6.

### 2.4. More Results on Video $\times 16$ Interpolation and Deblurring

Although video  $\times 16$  interpolation and deblurring is extremely difficult (without the help of extra auxiliary information), VIDUE is better at it than others, as shown in Figs. S7 and S8.

## References

- [1] Seungjun Nah, Tae H Kim, and Kyoung M Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *CVPR*, pages 3883–3891, 2017. 2, 4, 6, 7
- [2] Shuo Chen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *CVPR*, pages 1279–1288, 2017. 4, 5, 6
- [3] Xiangyu Xu, Siyao Li, Wenxiu Sun, Qian Yin, and Ming-Hsuan Yang. Quadratic video interpolation. In *NeurIPS*, pages 1–10, 2019. 1

\*Corresponding author: rendongweihi@gmail.com.

<i>Input</i>	A sequence of blurred frames (indim = $T \times 3$ )
<i>Layer 1</i>	Conv(indim, 64, 7, 1, 3); BN; LeakyReLU; ResBlock(64, 64, 3, 1, 1, BN = True) $\times$ 2
<i>Layer 2</i>	Conv(64, 128, 3, 2, 1); BN; LeakyReLU; ResBlock(128, 128, 3, 1, 1, BN = True) $\times$ 2
<i>Layer 3</i>	Conv(128, 256, 3, 2, 1); BN; LeakyReLU; ResBlock(256, 256, 3, 1, 1, BN = True) $\times$ 2
<i>Layer 4</i>	Conv(256, 256, 3, 2, 1); BN; LeakyReLU; ResBlock(256, 256, 3, 1, 1, BN = True) $\times$ 2
<i>Layer 5</i>	GlobalAvgPool
<i>Layer 6</i>	Linear(256, 1024); LeakyReLU; Linear(1024, 256)
<i>Output</i>	Exposure-aware feature representation $\mathbf{u}$

Table S1. The structure of the exposure-aware feature extractor  $g_e$ .  $T$  is the length of the input video sequence. BN stands for batch normalization. Convolution is in the form of (input channel, output channel, kernel size, stride, padding). The linear (*i.e.*, fully connected) layer is in the form of (input channel, output channel).

<i>Layer 1</i>	Conv( $D$ , $D$ , 3, 1, 1); BN; LeakyReLU
<i>Layer 2</i>	Conv( $D$ , $D$ , 3, 1, 1); BN
	Residual connection between the input and the output of <i>Layer 2</i> via addition

Table S2. The structure of ResBlock. We may remove BN by setting the corresponding indicator variable to False.  $D$  is the channel dimension of the feature maps from the previous layer.



Figure S1. More visual results on GoPro [1] ( $\times 8$ ). Blue and red patches are from deblurred and interpolated frames, respectively. TSP is short for CDVD-TSP, and MIMO is short for MIMOUNetPlus.

<i>Input</i>	One blurred frame (indim = 3)
<i>Layer 1</i>	PixelUnshuffle(2)
<i>Layer 2</i>	Conv(indim $\times$ 4, 16, 5, 1, 2); ResBlock(16, 16, 5, 1, 2, BN = False) $\times$ 3
<i>Layer 3</i>	Conv(16, 32, 5, 2, (1, 2, 1, 2)); ResBlock(32, 32, 5, 1, 2, BN = False) $\times$ 3
<i>Layer 4</i>	Conv(32, 64, 5, 2, (1, 2, 1, 2)); ResBlock(64, 64, 5, 1, 2, BN = False) $\times$ 3
<i>Layer 5</i>	Conv(64, 128, 1, 1, 0); LeakyReLU; Conv(128, 64, 3, 1, 1); LeakyReLU
<i>Layer 6</i>	SE(256, 64, 4); TransConv(64, 32, 4, 2, 1); LeakyReLU
	Concatenation of outputs from <i>Layer 6</i> and <i>Layer 3</i>
<i>Layer 7</i>	Conv(64, 128, 1, 1, 0); LeakyReLU; Conv(128, 64, 3, 1, 1); LeakyReLU
<i>Layer 8</i>	SE(256, 64, 4); TransConv(64, 16, 4, 2, 1); LeakyReLU
	Concatenation of outputs from <i>Layer 8</i> and <i>Layer 2</i>
<i>Layer 9</i>	Conv(32, 64, 1, 1, 0); LeakyReLU; Conv(64, 32, 3, 1, 1); LeakyReLU
<i>Layer 10</i>	SE(256, 32, 4); TransConv(32, 32, 4, 2, 1); LeakyReLU
<i>Layer 11</i>	Conv(32, 4, 5, 1, 2)
<i>Output</i>	Motion offsets

Table S3. The structure of intra-motion analysis. Adapting to  $\mathbf{u}$  via gain tuning is in the form of squeeze-and-excitation (SE) (vec channel, input channel, reduction). Transposed convolution is in the same format as convolution.

<i>Input</i>	A sequence of motion offsets (indim = $T \times 4$ )
<i>Layer 1</i>	Conv(indim, 32, 7, 1, 3); LeakyReLU; Conv(32, 32, 7, 1, 3); LeakyReLU
<i>Layer 2</i>	AvgPool(2); Conv(32, 64, 5, 1, 2); LeakyReLU; Conv(64, 64, 5, 1, 2); LeakyReLU
<i>Layer 3</i>	AvgPool(2); Conv(64, 128, 3, 1, 1); LeakyReLU; Conv(128, 128, 3, 1, 1); LeakyReLU
<i>Layer 4</i>	AvgPool(2); Conv(128, 256, 3, 1, 1); LeakyReLU; Conv(256, 256, 3, 1, 1); LeakyReLU
<i>Layer 5</i>	AvgPool(2); Conv(256, 512, 3, 1, 1); LeakyReLU; Conv(512, 512, 3, 1, 1); LeakyReLU
<i>Layer 6</i>	AvgPool(2); Conv(512, 512, 3, 1, 1); LeakyReLU; Conv(512, 512, 3, 1, 1); LeakyReLU
<i>Layer 7</i>	Up(2); Conv(512, 512, 3, 1, 1); LeakyReLU; Conv(512, 512, 3, 1, 1); LeakyReLU; SE(256, 512, 4)
<i>Layer 8</i>	Up(2); Conv(512, 256, 3, 1, 1); LeakyReLU; Conv(256, 256, 3, 1, 1); LeakyReLU; SE(256, 256, 4)
<i>Layer 9</i>	Up(2); Conv(256, 128, 3, 1, 1); LeakyReLU; Conv(256, 128, 3, 1, 1); LeakyReLU; SE(256, 128, 4)
<i>Layer 10</i>	Up(2); Conv(128, 64, 3, 1, 1); LeakyReLU; Conv(64, 64, 3, 1, 1); LeakyReLU; SE(256, 64, 4)
<i>Layer 11</i>	Up(2); Conv(64, 32, 3, 1, 1); LeakyReLU; Conv(32, 32, 3, 1, 1); LeakyReLU; SE(256, 32, 4)
<i>Layer 12</i>	Conv(32, $S \times T \times 2$ , 3, 1, 1)
<i>Output</i>	Inter-motion map $n$

Table S4. The structure of inter-motion analysis. Up( $\cdot$ ) stands for bilinear interpolation.

<i>Input</i>	A sequence of blurred frames (indim = $T \times 3$ )
<i>Layer 1</i>	Conv(indim, 64, 7, 1, 3); LeakyReLU; ResBlock(64, 64, 3, 1, 1, BN = False) $\times$ 3
<i>Layer 2</i>	Conv(64, 128, 3, 2, 1); LeakyReLU; ResBlock(128, 128, 3, 1, 1, BN = False) $\times$ 6
<i>Layer 3</i>	Conv(128, 256, 3, 2, 1); LeakyReLU; ResBlock(256, 256, 3, 1, 1, BN = False) $\times$ 6
<i>Layer 4</i>	Conv(256, 256, 3, 2, 1); LeakyReLU; ResBlock(256, 256, 3, 1, 1, BN = False) $\times$ 6
<i>Layer 5</i>	TransConv(256, 256, 3, 2, 1); LeakyReLU; ResBlock(256, 256, 3, 1, 1, BN = False) $\times$ 6
<i>Layer 6</i>	E_Conv(256, 256, 4)
<i>Layer 7</i>	Motion_Refine(256, 3, $S \times T$ ); Up(2) Addition of outputs from <i>Layer 6</i> and <i>Layer 3</i>
<i>Layer 8</i>	TransConv(256, 128, 3, 2, 1); LeakyReLU; ResBlock(128, 128, 3, 1, 1, BN = False) $\times$ 6
<i>Layer 9</i>	E_Conv(256, 128, 4)
<i>Layer 10</i>	Motion_Refine(128, 3, $S \times T$ ) Addition of outputs from <i>Layer 10</i> and <i>Layer 7</i>
<i>Layer 11</i>	Up(2) Addition of outputs from <i>Layer 9</i> and <i>Layer 2</i>
<i>Layer 12</i>	TransConv(128, 64, 3, 2, 1); LeakyReLU; ResBlock(64, 64, 3, 1, 1, BN = False) $\times$ 6
<i>Layer 13</i>	E_Conv(256, 64, 4)
<i>Layer 14</i>	Motion_Refine(64, 3, $S \times T$ )
<i>Layer 15</i>	Addition of outputs from <i>Layer 14</i> and <i>Layer 11</i> Addition of outputs from <i>Layer 13</i> and <i>Layer 1</i>
<i>Layer 16</i>	ResBlock(64, 64, 3, 1, 1, BN = False) $\times$ 3; Conv(64, $S \times T \times 3$ , 3, 1, 1); Split(3)
<i>Layer 17</i>	Addition of outputs from <i>Layer 16</i> and <i>Layer 15</i>
<i>Output</i>	Final sharp video sequence $\hat{x}$

Table S5. The architecture of video reconstruction network. Adapting to  $u$  via exposure-aware convolution is in the form E\_Conv (vec channel, input channel, reduction). Motion refinement is in the form of Motion\_Refine (input channel, kernel size, output sequence). Split( $\cdot$ ) stands for the channel splitting operation.

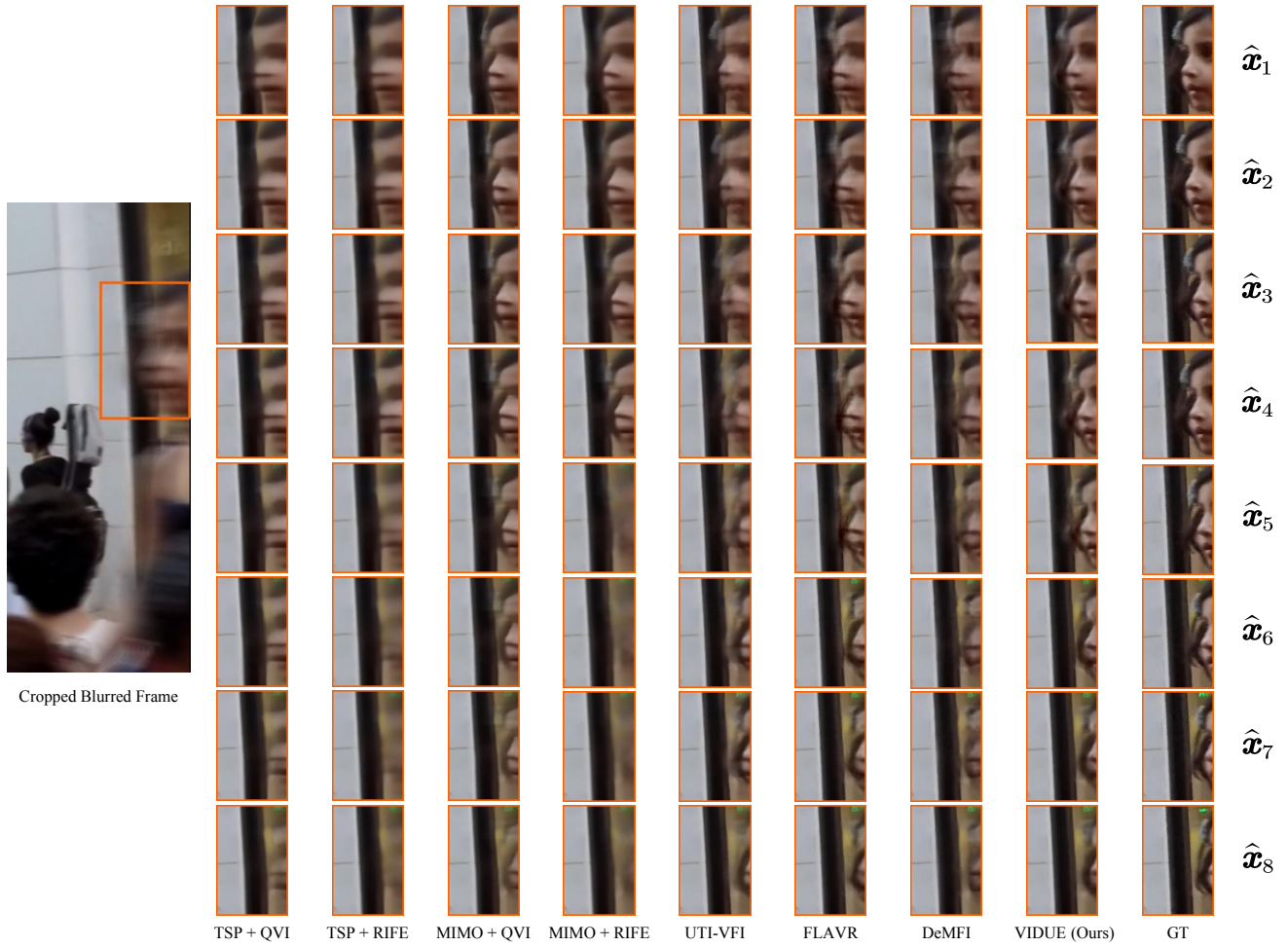


Figure S2. More visual results on GoPro [1] ( $\times 8$ ).



Figure S3. More visual results on Adobe [2] ( $\times 8$ ). Blue and red patches are from deblurred and interpolated frames, respectively.

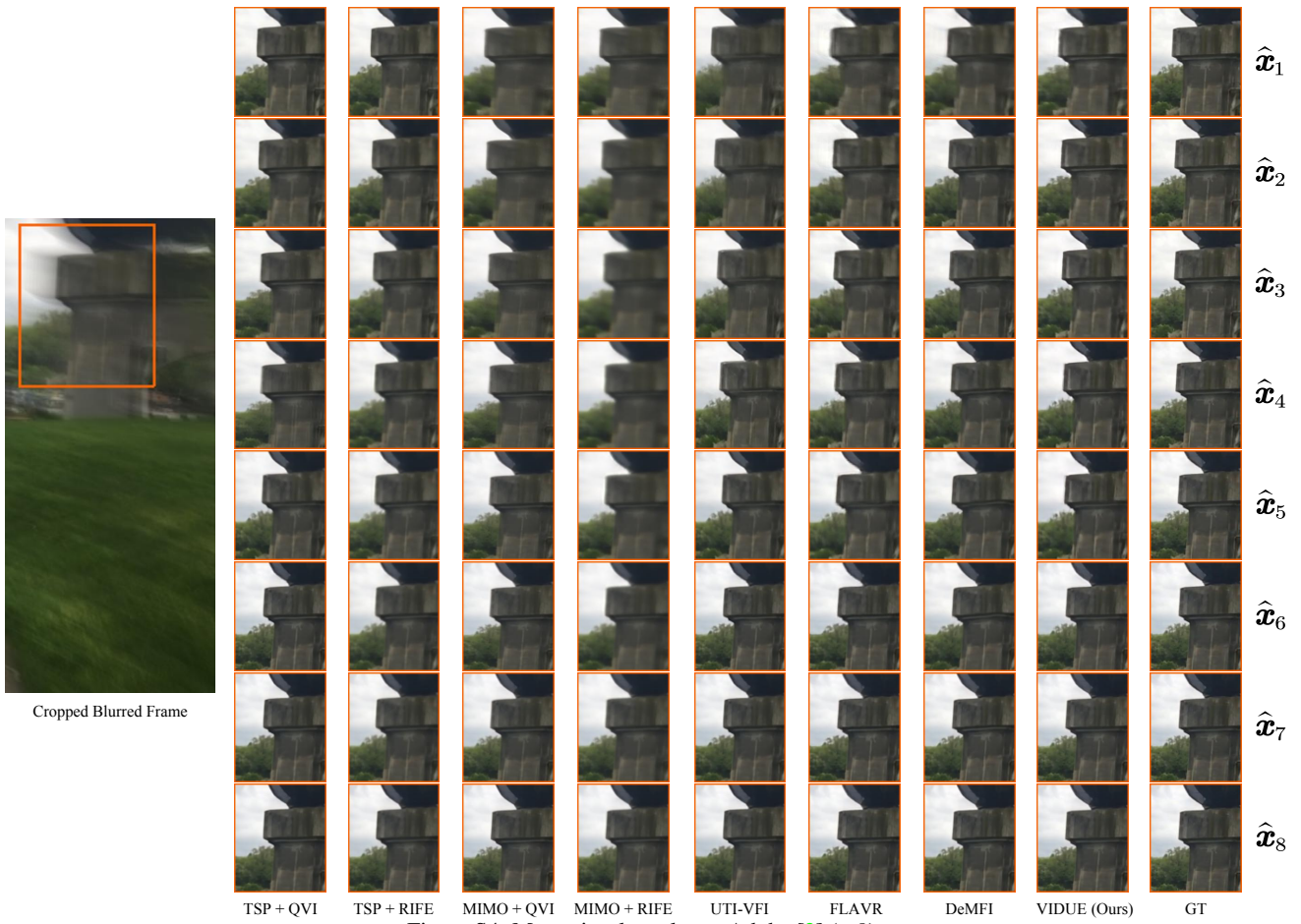


Figure S4. More visual results on Adobe [2] ( $\times 8$ ).

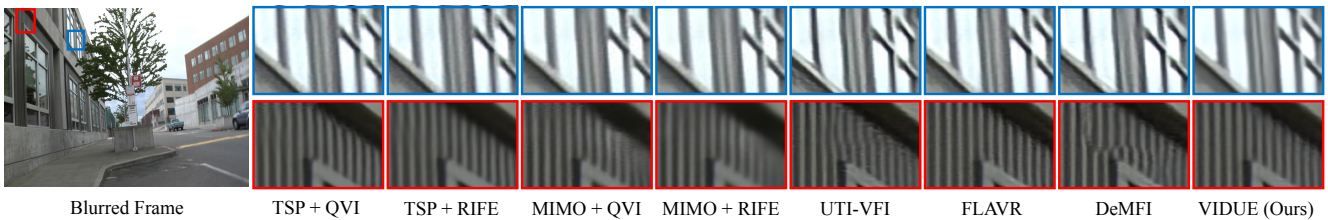


Figure S5. More visual results on RealBlur [2] ( $\times 8$ ). Blue and red patches are from deblurred and interpolated frames, respectively.



Figure S6. More visual results on RealBlur [2] ( $\times 8$ ).

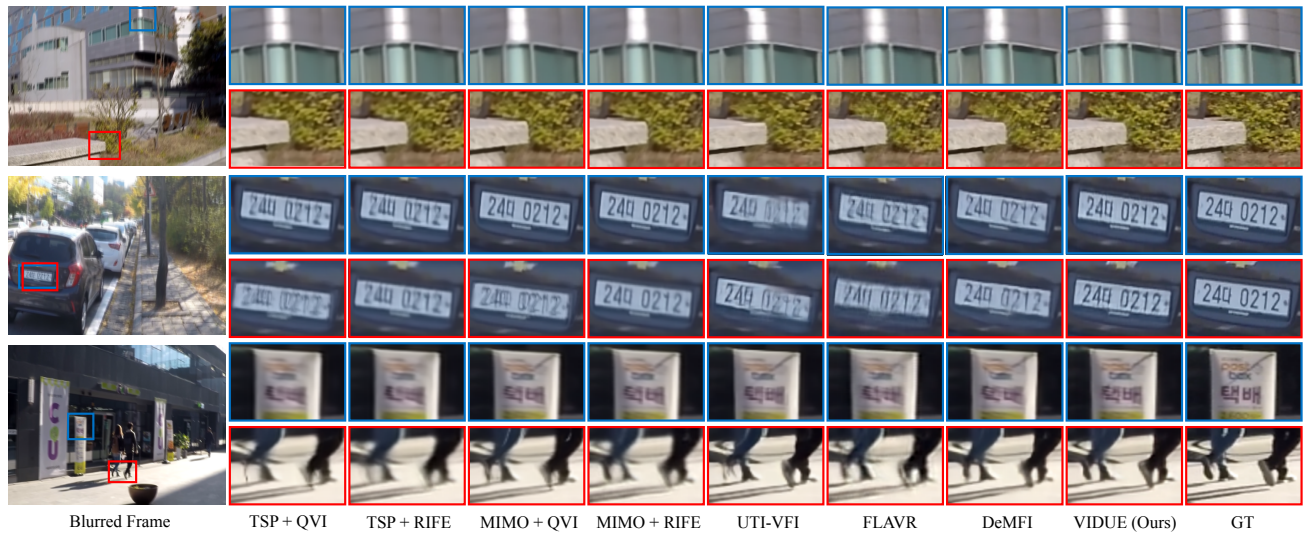


Figure S7. More visual results on GoPro [1] ( $\times 16$ ). Blue and red patches are from deblurred and interpolated frames, respectively.



Cropped Blurred Frame



Figure S8. More visual results on GoPro [1] ( $\times 16$ ).