

Bringing a Blurry Frame Alive at High Frame-Rate with an Event Camera

Liyuan Pan^{1,2}, Cedric Scheerlinck^{1,2}, Xin Yu^{1,2}, Richard Hartley^{1,2}, Miaomiao Liu^{1,2}, and Yuchao Dai³

¹ Research School of Engineering, Australian National University, Canberra, Australia

² Australia Centre for Robotic Vision

³ School of Electronics and Information, Northwestern Polytechnical University, Xi'an, China

liyuan.pan@anu.edu.au

Abstract

*Event-based cameras can measure intensity changes (called ‘events’) with microsecond accuracy under high-speed motion and challenging lighting conditions. With the active pixel sensor (APS), the event camera allows simultaneous output of the intensity frames. However, the output images are captured at a relatively low frame-rate and often suffer from motion blur. A blurry image can be regarded as the integral of a sequence of latent images, while the events indicate the changes between the latent images. Therefore, we are able to model the blur-generation process by associating event data to a latent image. In this paper, we propose a simple and effective approach, the **Event-based Double Integral (EDI)** model, to reconstruct a high frame-rate, sharp video from a single blurry frame and its event data. The video generation is based on solving a simple non-convex optimization problem in a single scalar variable. Experimental results on both synthetic and real images demonstrate the superiority of our EDI model and optimization method in comparison to the state-of-the-art.*

1. Introduction

Event cameras (such as the Dynamic Vision Sensor (DVS) [17] and the Dynamic and Active-pixel Vision Sensor (DAVIS) [3]) are sensors that asynchronously measure the intensity changes at each pixel independently with microsecond temporal resolution¹. The event stream encodes the motion information by measuring the precise pixel-by-pixel intensity changes. Event cameras are more robust to low lighting and highly dynamic scenes than traditional cameras since they are not affected by under/over exposure or motion blur associated with a synchronous shutter.

Due to the inherent differences between event cameras and standard cameras, existing computer vision algorithms designed for standard cameras cannot be applied to event

cameras directly. Although the DAVIS [3] can provide the simultaneous output of the intensity frames and the event stream, there still exist major limitations with current event cameras:

- **Low frame-rate intensity images:** In contrast to the high temporal resolution of event data ($\geq 3\mu s$ latency), the current event cameras only output low frame-rate intensity images ($\geq 5ms$ latency).
- **Inherent blurry effects:** When recording highly dynamic scenes, motion blur is a common issue due to the relative motion between the camera and the scene. The output of the intensity image from the APS tends to be blurry.

To address these above challenges, various methods have been proposed by reconstructing high frame-rate videos. The existing methods can be in general categorized as **1**) Event data only solutions [1, 27, 2], where the results tend to lack the texture and consistency of natural videos, as they fail to use the complementary information contained in the low frame-rate intensity image; **2**) Low frame-rate intensity-image-only solutions [11], where an end-to-end learning framework has been proposed to learn regression between a single blurry image and a video sequence, whereas the rich event data are not used; and **3**) Jointly exploiting event data and intensity images [29, 31, 4], building upon the interaction between both sources of information. However, these methods fail to address the blur issue associated with the captured image frame. Therefore, the reconstructed high frame-rate videos can be degraded by blur.

Although blurry frames cause undesired image degradation, they also encode the relative motion between the camera and the observed scene. Taking full advantage of the encoded motion information would benefit the reconstruction of high frame-rate videos.

To this end, we propose an **Event-based Double Integral (EDI)** model to resolve the above problems by reconstructing a high frame-rate video from a single image (even

¹If nothing moves in the scene, no events are triggered.

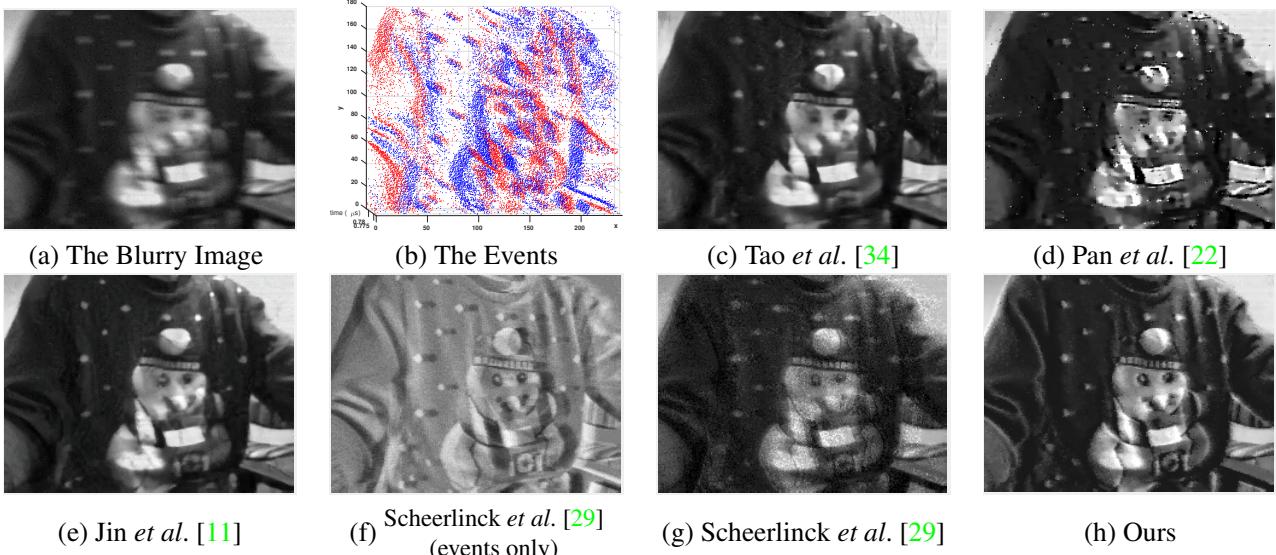


Figure 1. *Deblurring and reconstruction results of our method compared with the state-of-the-art methods on our real blurry event dataset.* (a) The input blurry image. (b) The corresponding event data. (c) Deblurring result of Tao et al. [34]. (d) Deblurring result of Pan et al. [22]. (e) Deblurring result of Jin et al. [11]. Jin uses video as training data to train a supervised model to perform deblur, where the video can also be considered as similar information as the event data. (f)-(g) Reconstruction results of Scheerlinck et al. [29], (f) from only events, (g) from combining events and frames. (h) Our reconstruction result. (Best viewed on screen).

blur) and its event sequence, where the blur effects have been reduced in each reconstructed frame. Our EDI model naturally relates the desired high frame-rate sharp video, the captured intensity frame and event data. Based on the EDI model, high frame-rate video generation is as simple as solving a non-convex optimization problem in a single scalar variable.

Our main contributions are summarized as follows.

- 1) We propose a simple and effective model, named the Event-based Double Integral (EDI) model, to restore a high frame-rate sharp video from a single image (even blur) and its corresponding event data.
- 2) Using our proposed formulation of EDI, we propose a stable and general method to generate a sharp video under various types of blur by solving a single variable non-convex optimization problem, especially in low lighting and complex dynamic conditions.
- 3) The frame rate of our reconstructed video can theoretically be as high as the event rate (200 times greater than the original frame rate in our experiments).

2. Related Work

Event cameras such as the DAVIS and DVS [3, 17] report log intensity changes, inspired by human vision. Although several works try to explore the advantages of the high temporal resolution provided by event cameras [39, 13, 26, 41, 40, 8, 15], how to make the best use of the event camera has not yet been fully investigated.

Event-based image reconstruction. Kim et al. [12] reconstruct high-quality images from an event camera under a strong assumption that the only movement is pure camera rotation, and later extend their work to handle 6-degree-of-freedom motion and depth estimation [13]. Bardow et al. [1] aim to simultaneously recover optical flow and intensity images. Reinbacher et al. [27] restore intensity images via manifold regularization. Barua et al. [2] generate image gradients by dictionary learning and obtain a logarithmic intensity image via Poisson reconstruction. However, the intensity images reconstructed by the previous approaches suffer from obvious artifacts as well as lack of texture due to the spatial sparsity of event data.

To achieve more image detail in the reconstructed images, several methods trying to combine events with intensity images have been proposed. The DAVIS [3] uses a shared photo-sensor array to simultaneously output events (DVS) and intensity images (APS). Scheerlinck et al. [29] propose an asynchronous event-driven complementary filter to combine APS intensity images with events, and obtain continuous-time image intensities. Brandli et al. [4] directly integrate events from a starting APS image frame, and each new frame resets the integration. Shedligeri et al. [31] first exploit two intensity images to estimate depth. Then, they use the event data only to reconstruct a pseudo-intensity sequence (using [27]) between the two intensity images and use the pseudo-intensity sequence to estimate ego-motion using visual odometry. Using the estimated 6-DOF pose and depth, they directly warp the intensity image to the intermediate location. Liu et al. [18] assume a scene should

have static background. Thus, their method needs an extra sharp static foreground image as input and the event data are used to align the foreground with the background.

Image deblurring. Traditional deblurring methods usually make assumptions on the scenes (such as a static scene) or exploit multiple images (such as stereo, or video) to solve the deblurring problem. Significant progress has been made in the field of single image deblurring. Methods using gradient based regularizers, such as Gaussian scale mixture [7], $l_1 \setminus l_2$ norm [14], edge-based patch priors [33] and l_0 -norm regularizer [36], have been proposed. Non-gradient-based priors such as the color line based prior [16], and the extreme channel (dark/bright channel) prior [22, 37] have also been explored. Another family of image deblurring methods tends to use multiple images [5, 10, 30, 23, 24].

Driven by the success of deep neural networks, Sun *et al.* [32] propose a convolutional neural network (CNN) to estimate locally linear blur kernels. Gong *et al.* [9] learn optical flow from a single blurry image through a fully-convolutional deep neural network. The blur kernel is then obtained from the estimated optical flow to restore the sharp image. Nah *et al.* [21] propose a multi-scale CNN that restores latent images in an end-to-end learning manner without assuming any restricted blur kernel model. Tao *et al.* [34] propose a light and compact network, SRN-DeblurNet, to deblur the image. However, deep deblurring methods generally need a large dataset to train the model and usually require sharp images provided as supervision. In practice, blurry images do not always have corresponding ground-truth sharp images.

Blurry image to sharp video. Recently, two deep learning based methods [11, 25] propose to restore a video from a single blurry image with a fixed sequence length. However, their reconstructed videos do not obey the 3D geometry of the scene and camera motion. Although deep-learning based methods achieve impressive performance in various scenarios, their success heavily depend on the consistency between the training datasets and the testing datasets, thus hinder the generalization ability for real-world applications.

3. Formulation

In this section, we develop an **EDI** model of the relationships between the events, the latent image and the blurry image. Our goal is to reconstruct a high frame-rate, sharp video from a single image and its corresponding events. This model can tackle various blur types and work stably in highly dynamic contexts and low lighting conditions.

3.1. Event Camera Model

Event cameras are bio-inspired sensors that asynchronously report logarithmic intensity changes [3, 17]. Unlike conventional cameras that produce the full image

at a fixed frame-rate, event cameras trigger events whenever the change in intensity at a given pixel exceeds a preset threshold. Event cameras do not suffer from the limited dynamic ranges typical of sensors with synchronous exposure time, and are able to capture high-speed motion with microsecond accuracy.

Inherent in the theory of event cameras is the concept of the latent image $L_{xy}(t)$, denoting the instantaneous intensity at pixel (x, y) at time t , related to the rate of photon arrival at that pixel. The latent image $L_{xy}(t)$ is not directly output by the camera. Instead, the camera outputs a sequence of *events*, denoted by (x, y, t, σ) , which record changes in the intensity of the latent image. Here, (x, y) are image coordinates, t is the time the event takes place, and polarity $\sigma = \pm 1$ denotes the direction (increase or decrease) of the intensity change at that pixel and time. Polarity is given by,

$$\sigma = \mathcal{T} \left(\log \left(\frac{\mathbf{L}_{xy}(t)}{\mathbf{L}_{xy}(t_{\text{ref}})} \right), c \right), \quad (1)$$

where $\mathcal{T}(\cdot, \cdot)$ is a truncation function,

$$\mathcal{T}(d, c) = \begin{cases} +1, & d \geq c, \\ 0, & d \in (-c, c), \\ -1, & d \leq -c. \end{cases}$$

Here, c is a threshold parameter determining whether an event should be recorded or not, $\mathbf{L}_{xy}(t)$ is the intensity at pixel (x, y) in the image $\mathbf{L}(t)$ and t_{ref} denotes the timestamp of the previous event. When an event is triggered, $\mathbf{L}_{xy}(t_{\text{ref}})$ at that pixel is updated to a new intensity level.

3.2. Intensity Image Formation

In addition to the event sequence, event cameras can provide a full-frame grey-scale intensity image, at a much slower rate than the event sequence. The grey-scale images may suffer from motion blur due to their long exposure time. A general model of image formation is given by,

$$\mathbf{B} = \frac{1}{T} \int_{f-T/2}^{f+T/2} \mathbf{L}(t) dt, \quad (2)$$

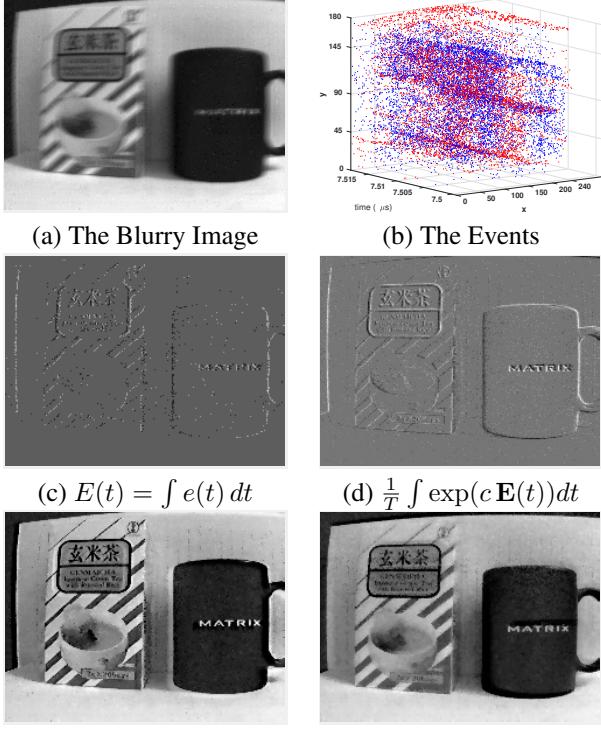
where \mathbf{B} is a blurry image, equal to the average value of the latent image during the exposure time $[f - T/2, f + T/2]$. This equation applies to each pixel (x, y) independently, and subscripts x, y denoting pixel location are often omitted henceforth.

3.3. Event-based Double Integral Model

We aim to recover a sequence of latent intensity images by exploiting both the blur model and the event model. We define $e_{xy}(t)$ as a function of continuous time t such that

$$e_{xy}(t) = \sigma \delta_{t_0}(t).$$

whenever there is an event (x, y, t_0, σ) . Here, $\delta_{t_0}(t)$ is an impulse function, with unit integral, at time t_0 , and the sequence of events is turned into a continuous time signal,



(e) Sample Frames of Our Reconstructed Video

Figure 2. The event data and our reconstructed result, where (a) and (b) are the input of our method. (a) The intensity image from the event camera. (b) Events from the event camera plotted in 3D space-time (x, y, t) (blue: positive event; red: negative event). (c) The first integral of several events during a small time interval. (d) The second integral of events during the exposure time. (e) Samples from our reconstructed video from $\mathbf{L}(0)$ to $\mathbf{L}(200)$.

consisting of a sequence of impulses. There is such a function $e_{xy}(t)$ for every point (x, y) in the image. Since each pixel can be treated separately, we omit the subscripts x, y .

During an exposure period $[f - T/2, f + T/2]$, we define $\mathbf{E}(t)$ as the sum of events between time f and t at a given pixel,

$$\mathbf{E}(t) = \int_f^t e(s) ds,$$

which represents the proportional change in intensity between time f and t . Except under extreme conditions, such as glare and no-light conditions, the latent image sequence $\mathbf{L}(t)$ is expressed as,

$$\mathbf{L}(t) = \mathbf{L}(f) \exp(c \mathbf{E}(t)) = \mathbf{L}(f) \exp(c)^{\mathbf{E}(t)}. \quad (3)$$

In particular, an event (x, y, t, σ) is triggered when the intensity of a pixel (x, y) increases or decreases by an amount c at time t . We put a tilde on top of things to denote logarithm, e.g. $\tilde{\mathbf{L}}(t) = \log(\mathbf{L}(t))$.

$$\tilde{\mathbf{L}}(t) = \tilde{\mathbf{L}}(f) + c \mathbf{E}(t). \quad (4)$$

Given a sharp frame, we can reconstruct a high frame-rate video from the sharp starting point $\mathbf{L}(f)$ by using Eq. (4). When the input image is blurry, a trivial solution would be to first deblur the image with an existing deblurring method and then to reconstruct the video using Eq. (4) (see Fig. 6 for details). However, in this way, the event data between intensity images is not fully exploited, thus resulting in inferior performance. Instead, we propose to reconstruct the video by exploiting the inherent connection between event and blur, and present the following model.

As for the blurred image,

$$\begin{aligned} \mathbf{B} &= \frac{1}{T} \int_{f-T/2}^{f+T/2} \mathbf{L}(t) dt \\ &= \frac{\mathbf{L}(f)}{T} \int_{f-T/2}^{f+T/2} \exp\left(c \int_f^t e(s) ds\right) dt. \end{aligned} \quad (5)$$

In this manner, we construct the relation between the captured blurry image \mathbf{B} and the latent image $\mathbf{L}(f)$ through the double integral of the event. We name Eq.(5) the **Event-based Double Integral (EDI)** model.

Taking the logarithm on both sides of Eq. (5) and rearranging, yields

$$\tilde{\mathbf{L}}(f) = \tilde{\mathbf{B}} - \log\left(\frac{1}{T} \int_{f-T/2}^{f+T/2} \exp(c \mathbf{E}(t)) dt\right), \quad (6)$$

which shows a linear relation between the blurry image, the latent image and the integral of the events in the log space.

3.4. High Frame-Rate Video Generation

The right-hand side of Eq. (6) is known, apart from perhaps the value of the contrast threshold c , the first term from the grey-scale image, the second term from the event sequence, it is possible to compute $\tilde{\mathbf{L}}(f)$, and hence $\mathbf{L}(f)$ by exponentiation. Subsequently, from Eq. (4) the latent image $\mathbf{L}(t)$ at any time may be computed.

To avoid accumulated errors of constructing a video from many frames of a blurred video, it is more suitable to construct each frame $\mathbf{L}(t)$ using the closest blurred frame.

Theoretically, we could generate a video with frame-rate as high as the DVS's eps (events per second). However, as each event carries little information and is subject to noise, several events must be processed together to yield a reasonable image. We generate a reconstructed frame every 50-100 events, so for our experiments, the frame-rate of the reconstructed video is usually 200 times greater than the input low frame-rate video. Furthermore, as indicated by Eq. (6), the challenging blind motion deblurring problem has been reduced to a single variable optimization problem of how to find the best value of the contrast threshold c . In the following section, we use $\mathbf{L}(c, t)$ to present the latent sharp image $\mathbf{L}(t)$ with different c .

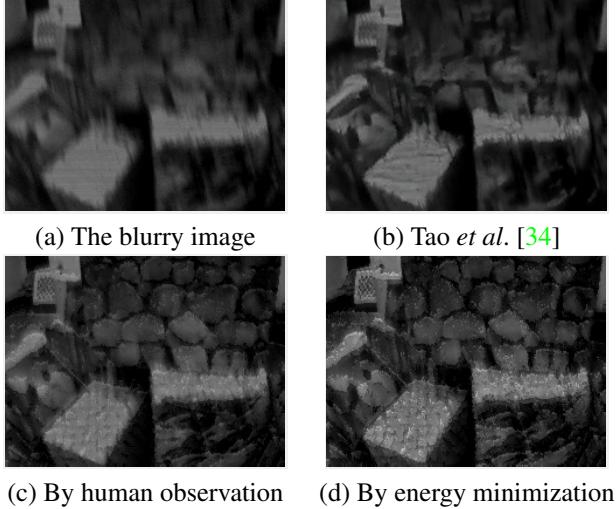


Figure 3. An example of our reconstruction result using different methods to estimate c , from the real dataset [20]. (a) The blurry image. (b) Deblurring result of [34] (c) Our result where c is chosen by manual inspection. (d) Our result where c is computed automatically by our proposed energy minimization (9).

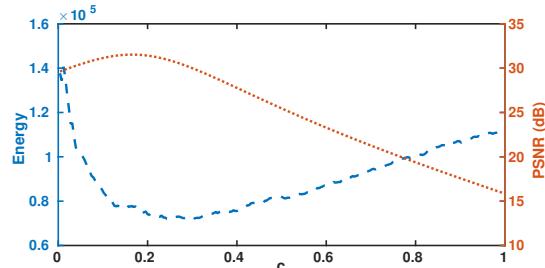


Figure 4. The figure plot deblurring performance against the value of c . The image is clearer with higher PSNR value.

4. Optimization

The unknown contrast threshold c represents the minimum change in log intensity required to trigger an event. By choosing an appropriate c in Eq. (5), we can generate a sequence of sharper images. To this end, we first need to evaluate the sharpness of the reconstructed images. Here, we propose two different methods to estimate the unknown variable c : manually chosen and automatically optimized.

4.1. Manually Chosen c

According to our EDI model in Eq. (5), given a value for c , we can obtain a sharp image. Therefore, we develop a method for deblurring by manually inspecting the visual effect of the deblurred image. In this way, we incorporate human perception into the reconstruction loop and the deblurred images should satisfy human observation. In Fig. 3, we give an example for manually chosen and automatically optimized results on dataset from [20].

4.2. Automatically Chosen c

To automatically find the best c , we need to build an evaluation metric (energy function) that can evaluate the quality of the deblurred image $L(c, t)$. Specifically, we propose to exploit different prior knowledge for sharp images and the event data.

4.2.1 Edge Constraint for Event Data

As mentioned before, when a proper c is given, our reconstructed image $L(c, t)$ will contain much sharper edges compared with the original input intensity image. Furthermore, event cameras inherently yield responses at moving intensity boundaries, so edges in the latent image may be located where (and when) events occur. This allows us to find latent image edges. An edge at time t corresponds to an event (at the pixel in question) during some time interval around t so we convolve the event sequence with an exponentially decaying window, to obtain an edge map,

$$M(t) = \int_0^T \exp(-(\alpha|t-s|)) e(t) ds,$$

where α is a weight parameter for time attenuation, which we set to 1. As $M(t)$ mainly reports the intensity changes but not the intensity value itself, we need to change $M(t)$ to an edge map $S(M(t))$ using the Sobel filter, which is also applied to $L(c, t)$. (See Fig. 5 for details).

Here, we use cross-correlation between $S(L(c, t))$ and $S(M(t))$ to evaluate the sharpness of $L(c, t)$.

$$\phi_{\text{edge}}(L(c, t)) = \sum_{x,y} S(L(c, t))(x, y) \cdot S(M(t))(x, y). \quad (7)$$

4.2.2 Regularizing the Intensity Image

In our model, total variation is used to suppress noise in the latent image while preserving edges, and penalize the spatial fluctuations[28]. Therefore, we use conventional total variation (TV) based regularization

$$\phi_{\text{TV}}(L(c, t)) = |\nabla L(c, t)|_1, \quad (8)$$

where ∇ represents the gradient operators.

4.2.3 Energy Minimization

The optimal c can be estimate by solving Eq. (9),

$$\min_c \phi_{\text{TV}}(L(c, t)) + \lambda \phi_{\text{edge}}(L(c, t)), \quad (9)$$

where λ is a trade-off parameter. The response of cross-correlation reflect the matching rate of $L(c, t)$ and $M(t)$ which makes $\lambda < 0$. This single-variable minimization problem can be solved by the nonlinear least-squares method [19], Scatter-search[35] or Fibonacci search [6].

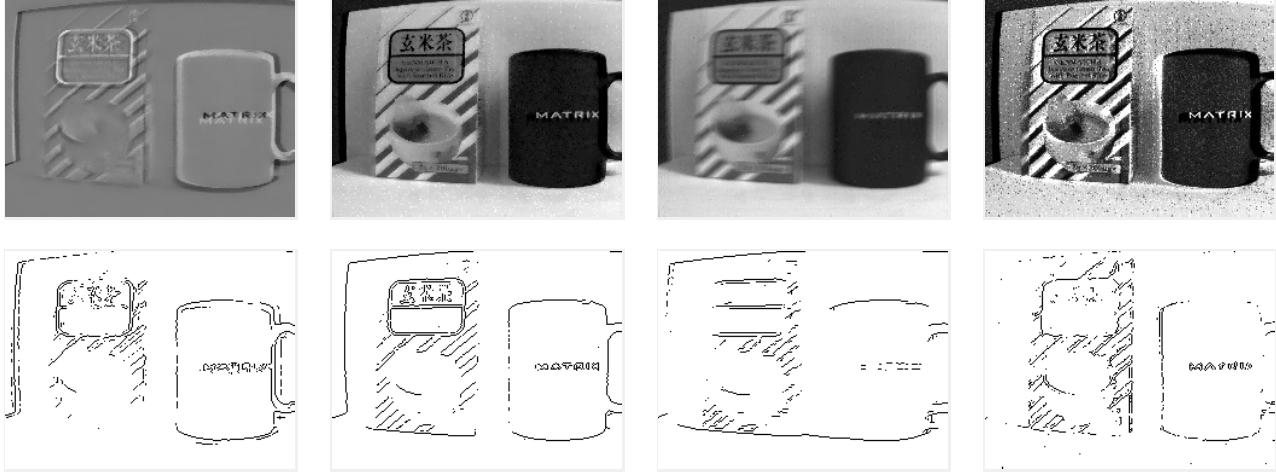


Figure 5. At left, the edge image $M(f)$ and below, its Sobel edge map. To the right are 3 reconstructed latent images using different values of c , low 0.03, middle 0.11 and high 0.55. Above, the reconstructed images, below, their Sobel edge maps. The optimal value of the threshold c is found by computing the cross-correlation of such images with the edge map at the left. (Best viewed on screen).

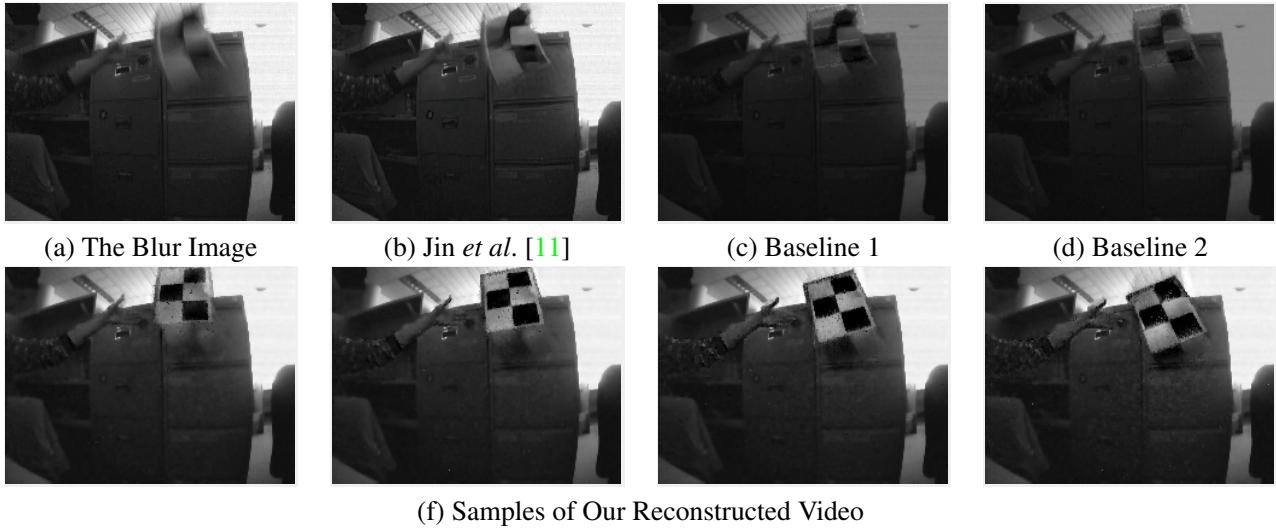


Figure 6. Deblurring and reconstruction results on our real blurry event dataset. (a) Input blurry images. (b) Deblurring result of [11]. (c) Baseline 1 for our method. We first use the state-of-the-art video-based deblurring method [11] to recover a sharp image. Then use the sharp image as input to a state-of-the-art reconstruction method [29] to get the intensity image. (d) Baseline 2 for our method. We first use method [29] to reconstruct an intensity image. Then use a deblurring method [11] to recover a sharp image. (e) Samples from our reconstructed video from $L(0)$ to $L(150)$. (Best viewed on screen).

In Fig. 4, we illustrate the clearness of the reconstructed image against the value of c . Meanwhile, we also provide the PSNR of the corresponding reconstructed image. As demonstrated in the figure, our proposed reconstruction metric could locate/identify the best deblurred image with peak PSNR properly.

5. Experiment

5.1. Experimental Setup

Synthetic dataset. In order to provide a quantitative comparison, we build a synthetic dataset based on the GoPro

blurry dataset [21]. It supplies ground truth videos which are used to generate the blurry images. Similarly, we employ the ground-truth images to generate event data based on the methodology of *event camera model*.

Real dataset. We evaluate our method on a public Event-Camera dataset [20], which provides a collection of sequences captured by the event camera for high-speed robotics. Furthermore, we present our real *blurry event dataset*², where each real sequence is captured with the DAVIS[3] under different conditions, such as indoor, out-

²To be released with codes

Table 1. Quantitative comparisons on the Synthetic dataset [21]. This dataset provides videos can be used to generate not only blurry images but also event data. All methods are tested under the same blurry condition, where methods [21, 11, 34, 38] use GoPro dataset [21] to train their models. Jin [11] achieves their best performance when the image is down-sampled to 45% mentioned in their paper.

Average result of the deblurred images on dataset[21]								
	Pan et al. [22]	Sun et al. [32]	Gong et al. [9]	Jin et al. [11]	Tao et al. [34]	Zhang et al. [38]	Nah et al. [21]	Ours
PSNR(dB)	23.50	25.30	26.05	26.98	30.26	29.18	29.08	29.06
SSIM	0.8336	0.8511	0.8632	0.8922	0.9342	0.9306	0.9135	0.9430
Average result of the reconstructed videos on dataset[21]								
	Baseline 1 [34] + [29]	Baseline 2 [29] + [34]	Scheerlinck et al. [29]	Jin et al. [11]	Ours			
PSNR(dB)	25.52	26.34	25.84	25.62	28.49			
SSIM	0.7685	0.8090	0.7904	0.8556	0.9199			



Figure 7. An example of the reconstructed result on our synthetic event dataset based on the GoPro dataset [21]. [21] provides videos to generate the blurry images and event data. (a) The blurry image. The red close-up frame is for (b)-(e), the yellow close-up frame is for (f)-(g). (b) The deblurring result of Jin et al. [11]. (c) Our deblurring result. (d) The crop of their reconstructed images and the frame number is fixed at 7. Jin et al. [11] uses the GoPro dataset added with 20 scenes as training data and their model is supervised by 7 consecutive sharp frames. (e) The crop of our reconstructed images. (f) The crop of Reinbacher [27] reconstructed images from only events. (g) The crop of Scheerlinck [29] reconstructed image, they use both events and the intensity image. For (e)-(g), the shown frames are the chosen examples, where the length of the reconstructed video is based on the number of events.

door scenery, low lighting conditions, and different motion patterns (e.g., camera shake, objects motion) that naturally introduce motion blur into the APS intensity images.

Implementation details. For all our real experiments, we use the DAVIS that shares photosensor array to simultaneously output events (DVS) and intensity images (APS). The framework is implemented by using MATLAB with C++ wrappers. It takes around 1.5 second to process one image on a single i7 core running at 3.6 GHz.

5.2. Experimental Results

We compare our proposed approach with state-of-the-art blind deblurring methods, including conventional deblurring methods [22, 37], deep based dynamic scene deblurring methods [21, 11, 34, 38, 32], and event-based image reconstruction methods [27, 29]. Moreover, Jin et al. [11] can restore a video from a single blurry image based on a deep network, where the middle frame in the restored odd-numbered sequence is the best.

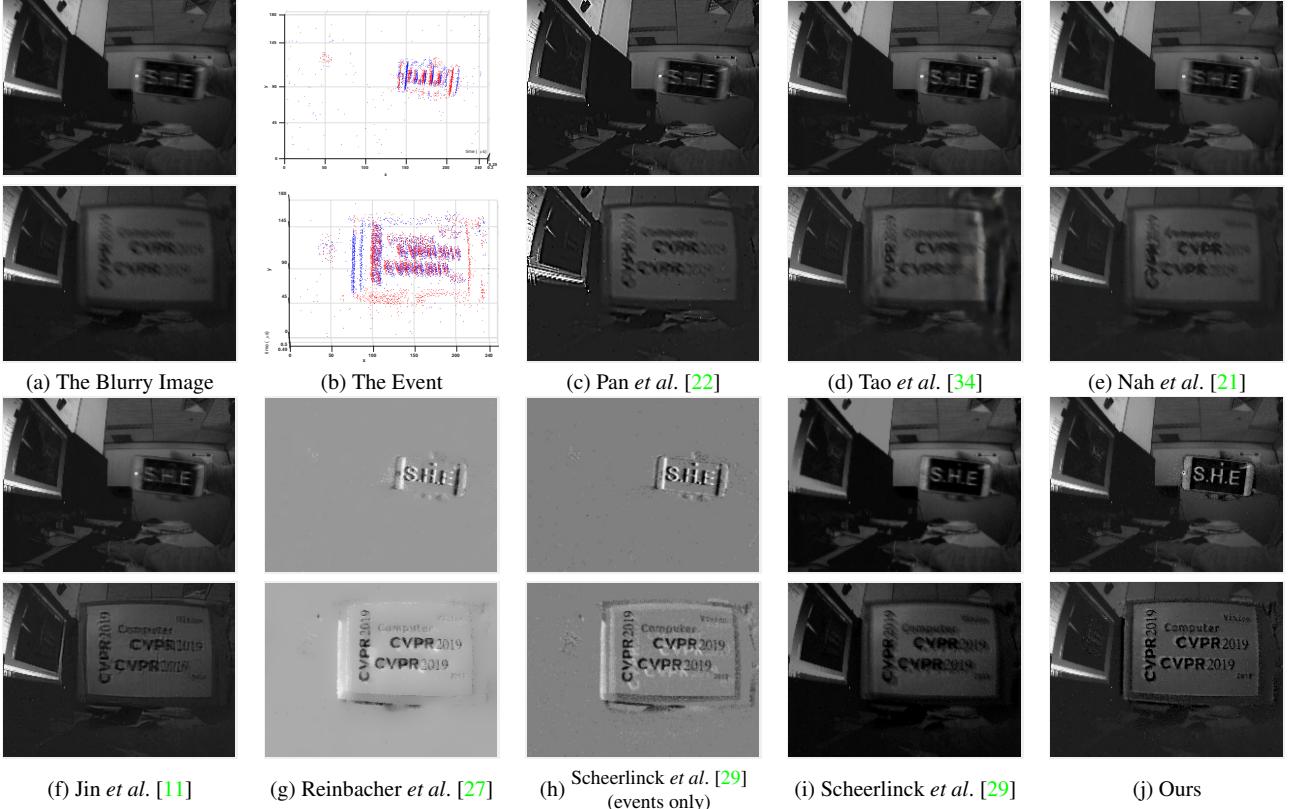


Figure 8. Examples of reconstruction result on our real blurry event dataset in low lighting and complex dynamic conditions (a) Input blurry images. (b) The event information. (c) Deblurring results of [22]. (d) Deblurring results of [34]. (e) Deblurring results of [21]. (f) Deblurring results of [11] and they use video as training data. (g) Reconstruction result of [27] from only events. (h)-(i) Reconstruction results of [29], (h) from only events, (i) from combining events and frames. (j) Our reconstruction result. Results in (c)-(f) show that real high dynamic settings and low light condition is still challenging in the deblurring area. Results in (g)-(h) show that while intensity information of a scene is still retained with an event camera recording, color, and delicate texture information cannot be recovered.

In order to prove the effective of our **EDI** model, we show some baseline comparisons in Fig. 6 and Table 1. For baseline 1, we first apply a state-of-the-art deblurring method [34] to recover a sharp image, and then the recovered image as an input is then fed to a reconstruction method [29]. For baseline 2, we first use the video reconstruction method to construct a sequence of intensity images, and then apply the deblurring method to each frame. As seen in Table 1, our approach obtains higher PSNR and SSIM in comparison to both baseline 1 and baseline 2. This also implies that our approach better exploits the event data to not only recover sharp images but also reconstruct high frame-rate videos.

In Table 1, we show the quantitative comparisons with the state-of-the-art image deblurring approaches [32, 22, 9, 11, 34, 38, 21], and the video reconstruction method [29] on our synthetic dataset, respectively. As indicated in Table 1, our approach achieves the best performance on SSIM and competitive result on PSNR compared to the state-of-the-art methods, and attains significant performance improvements

on high-frame video reconstruction.

In Fig. 7, we compare our generated video frames with the state-of-the-art deblurring methods [22, 11, 34, 21] qualitatively. Furthermore, image reconstruction methods [27, 29] are also included for comparisons. Fig. 7 shows that our method can generate more frames from a single blurry image and the recovered frames are much sharper.

We also report our reconstruction results on the real dataset, including text images and low-lighting images, in Fig. 1, Fig. 2, Fig. 3 and Fig. 8. Compared with state-of-the-art deblurring methods, our method achieves superior results. In comparison to existing event-based image reconstructed methods [27, 29], our reconstructed images are not only more realistic but also contain richer details. More deblurring results and **high-temporal resolution videos** are shown in the supplementary material.

6. Conclusion

In this paper, we propose a **Event-based Double Integral (EDI)** model to naturally connect intensity images and event data captured by the event camera, which also takes

the blur generation process into account. In this way, our model can be used to not only recover latent sharp images but also reconstruct intermediate frames at high frame-rate. We also propose a simple yet effective method to solve our EDI model. Due to the simplicity of our optimization process, our method is efficient as well. Extensive experiments show that our method can generate high-quality high frame-rate videos efficiently under different conditions, such as low lighting and complex dynamic scenes.

References

- [1] P. Bardow, A. J. Davison, and S. Leutenegger. Simultaneous optical flow and intensity estimation from an event camera. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 884–892, 2016. [1](#), [2](#)
- [2] S. Barua, Y. Miyatani, and A. Veeraraghavan. Direct face detection and video reconstruction from event cameras. In *IEEE Winter Conf. Appl. Comput. Vis. (WACV)*, pages 1–9, 2016. [1](#), [2](#)
- [3] C. Brandli, R. Berner, M. Yang, S.-C. Liu, and T. Delbrück. A 240×180 130 db 3 μ s latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014. [1](#), [2](#), [3](#), [6](#), [11](#)
- [4] C. Brandli, L. Muller, and T. Delbrück. Real-time, high-speed video decompression using a frame- and event-based DAVIS sensor. In *IEEE Int. Symp. Circuits Syst. (ISCAS)*, pages 686–689, June 2014. [1](#), [2](#)
- [5] S. Cho, J. Wang, and S. Lee. Video deblurring for hand-held cameras using patch-based synthesis. *ACM Transactions on Graphics (TOG)*, 31(4):64, 2012. [3](#)
- [6] R. A. Dunlap. *The golden ratio and Fibonacci numbers*. World Scientific, 1997. [5](#)
- [7] R. Fergus, B. Singh, A. Hertzmann, S. T. Roweis, and W. T. Freeman. Removing camera shake from a single photograph. *ACM Trans. Graph.*, 25:787–794, 2006. [3](#)
- [8] D. Gehrig, H. Rebecq, G. Gallego, and D. Scaramuzza. Asynchronous, photometric feature tracking using events and frames. In *Eur. Conf. Comput. Vis. (ECCV)*, 2018. [2](#)
- [9] D. Gong, J. Yang, L. Liu, Y. Zhang, I. 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 *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 2319–2328, 2017. [3](#), [7](#), [8](#)
- [10] T. Hyun Kim and K. Mu Lee. Generalized video deblurring for dynamic scenes. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 5426–5434, 2015. [3](#)
- [11] M. Jin, G. Meishvili, and P. Favaro. Learning to extract a video sequence from a single motion-blurred image. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, June 2018. [1](#), [2](#), [3](#), [6](#), [7](#), [8](#), [11](#), [12](#), [14](#)
- [12] H. Kim, A. Handa, R. Benosman, S.-H. Ieng, and A. J. Davison. Simultaneous mosaicing and tracking with an event camera. In *British Machine Vis. Conf. (BMVC)*, 2014. [2](#)
- [13] H. Kim, S. Leutenegger, and A. J. Davison. Real-time 3D reconstruction and 6-DoF tracking with an event camera. In *Eur. Conf. Comput. Vis. (ECCV)*, pages 349–364, 2016. [2](#)
- [14] D. Krishnan, T. Tay, and R. Fergus. Blind deconvolution using a normalized sparsity measure. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 233–240, 2011. [3](#)
- [15] B. Kueng, E. Mueggler, G. Gallego, and D. Scaramuzza. Low-latency visual odometry using event-based feature tracks. In *IEEE/RSJ Int. Conf. Intell. Robot. Syst. (IROS)*, pages 16–23, Daejeon, Korea, Oct. 2016. [2](#)
- [16] W.-S. Lai, J.-J. Ding, Y.-Y. Lin, and Y.-Y. Chuang. Blur kernel estimation using normalized color-line prior. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 64–72, 2015. [3](#)
- [17] P. Lichtsteiner, C. Posch, and T. Delbrück. A 128×128 120 db 15 μ s latency asynchronous temporal contrast vision sensor. *IEEE journal of solid-state circuits*, 43(2):566–576, 2008. [1](#), [2](#), [3](#)
- [18] H.-C. Liu, F.-L. Zhang, D. Marshall, L. Shi, and S.-M. Hu. High-speed video generation with an event camera. *The Visual Computer*, 33(6-8):749–759, 2017. [2](#)
- [19] J. J. Moré. The levenberg-marquardt algorithm: implementation and theory. In *Numerical analysis*, pages 105–116. Springer, 1978. [5](#)
- [20] E. Mueggler, H. Rebecq, G. Gallego, T. Delbrück, and D. Scaramuzza. The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and slam. *The International Journal of Robotics Research*, 36(2):142–149, 2017. [5](#), [6](#)
- [21] S. Nah, T. H. Kim, and K. M. Lee. Deep multi-scale convolutional neural network for dynamic scene deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, July 2017. [3](#), [6](#), [7](#), [8](#), [11](#), [12](#), [14](#)
- [22] J. Pan, D. Sun, H. Pfister, and M.-H. Yang. Deblurring images via dark channel prior. *IEEE Trans. Pattern Anal. Mach. Intell.*, 2017. [2](#), [3](#), [7](#), [8](#), [11](#), [14](#)
- [23] L. Pan, Y. Dai, M. Liu, and F. Porikli. Simultaneous stereo video deblurring and scene flow estimation. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, July 2017. [3](#)

- [24] L. Pan, Y. Dai, M. Liu, and F. Porikli. Depth map completion by jointly exploiting blurry color images and sparse depth maps. In *Applications of Computer Vision (WACV), 2018 IEEE Winter Conference on*, pages 1377–1386. IEEE, 2018. 3
- [25] K. Purohit, A. Shah, and A. Rajagopalan. Bringing alive blurred moments! *arXiv preprint arXiv:1804.02913*, 2018. 3
- [26] H. Rebecq, T. Horstschäfer, G. Gallego, and D. Scaramuzza. EVO: A geometric approach to event-based 6-DOF parallel tracking and mapping in real-time. *IEEE Robot. Autom. Lett.*, 2, 2017. 2
- [27] C. Reinbacher, G. Graber, and T. Pock. Real-time intensity-image reconstruction for event cameras using manifold regularisation. In *British Machine Vis. Conf. (BMVC)*, 2016. 1, 2, 7, 8
- [28] L. I. Rudin, S. Osher, and E. Fatemi. Nonlinear total variation based noise removal algorithms. *Physica D: nonlinear phenomena*, 60(1-4):259–268, 1992. 5
- [29] C. Scheerlinck, N. Barnes, and R. Mahony. Continuous-time intensity estimation using event cameras. *arXiv e-prints*, Nov. 2018. 1, 2, 6, 7, 8, 12
- [30] A. Sellent, C. Rother, and S. Roth. Stereo video deblurring. In *Eur. Conf. Comput. Vis. (ECCV)*, pages 558–575. Springer, 2016. 3
- [31] P. A. Sheldinger, K. Shah, D. Kumar, and K. Mitra. Photorealistic image reconstruction from hybrid intensity and event based sensor. *arXiv preprint arXiv:1805.06140*, 2018. 1, 2
- [32] J. Sun, W. Cao, Z. Xu, and J. Ponce. Learning a convolutional neural network for non-uniform motion blur removal. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 769–777, 2015. 3, 7, 8
- [33] L. Sun, S. Cho, J. Wang, and J. Hays. Edge-based blur kernel estimation using patch priors. In *IEEE Int. Conf. Comput. Photography (ICCP)*, 2013. 3
- [34] X. Tao, H. Gao, X. Shen, J. Wang, and J. Jia. Scale-recurrent network for deep image deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, June 2018. 2, 3, 5, 7, 8, 11, 12, 14
- [35] Z. Ugray, L. Lasdon, J. Plummer, F. Glover, J. Kelly, and R. Martí. Scatter search and local nlp solvers: A multistart framework for global optimization. *INFORMS Journal on Computing*, 19(3):328–340, 2007. 5
- [36] L. Xu, S. Zheng, and J. Jia. Unnatural l0 sparse representation for natural image deblurring. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 1107–1114, 2013. 3
- [37] Y. Yan, W. Ren, Y. Guo, R. Wang, and X. Cao. Image deblurring via extreme channels prior. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, July 2017. 3, 7, 11, 14
- [38] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 7, 8
- [39] Y. Zhou, G. Gallego, H. Rebecq, L. Kneip, H. Li, and D. Scaramuzza. Semi-dense 3d reconstruction with a stereo event camera. *arXiv preprint arXiv:1807.07429*, 2018. 2
- [40] A. Zhu, L. Yuan, K. Chaney, and K. Daniilidis. Event-flownet: Self-supervised optical flow estimation for event-based cameras. In *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018. 2
- [41] A. Z. Zhu, N. Atanasov, and K. Daniilidis. Event-based visual inertial odometry. In *IEEE Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 5816–5824, 2017. 2

Appendix

In this supplementary material, we provide more details about our datasets (Sec. A). Section B show details in high frame-rate video generation. We also give an example of how to run our testing code (Sec. C).

A. Datasets

A.1. Synthetic Dataset

In order to qualitatively comparing our method with state-of-the-art deblurring methods [22, 37, 11, 34, 21], we build a synthetic dataset based on the GoPro blurry dataset [21]. We use the provided videos from [21] to generate event data. The comparing results are presented in Fig. 11. It clearly shows that our method can generate more frames from a single blurry image and the recovered frames are much sharper.

A.2. Real Dataset

A Dynamic and Active-pixel Vision Sensor (DAVIS) [3] asynchronously measures the intensity changes at each pixel independently with microsecond temporal resolution. It uses a shared photo-sensor array to simultaneously output events from its Dynamic Vision Sensor (DVS) and intensity images from its Active Pixel Sensor (APS).

In all our real experiments, we use a DAVIS to build our real *Blurry event dataset*. The sequences are taken under different conditions, such as indoor, outdoor, low-lighting, and different motion patterns (camera shake, moving object), *etc.*

Our intensity images are recorded at a frame rate from 5 fps (frames per second) to 20 fps for different sequences. The resolution of the intensity images is 240×180 . The maximum frame-rate of our recorded events is 200 keps (thousands of events per second) since the events depend on the intensity changes of the scene.

B. High Frame-Rate Videos

The input of our method is a single image and its event data during the exposure time. Given a sharp frame $\mathbf{L}(f)$, setting it as a starting point, we can reconstruct a high frame-rate video by using Eq. (10).

$$\mathbf{L}(t) = \mathbf{L}(f) \exp \left(c \int_f^t e(s) ds \right). \quad (10)$$

When the input image is blurry, a trivial solution would be to first deblur the image with an existing deblurring method and then using Eq. (10) to reconstruct the video (see Fig. 9 for details). However, in this way, the events between consecutive intensity images are not fully exploited, resulting in inferior performance. We propose to exploit the inherent connection between event and blur to reconstruct the

video, and give the following model:

$$\log(\mathbf{L}(f)) = \log(\mathbf{B}) - \log \left(\frac{1}{T} \int_{f-T/2}^{f+T/2} \exp(c \mathbf{E}(t)) dt \right). \quad (11)$$

The right-hand side of Eq. (11) is known, apart from perhaps the value of the contrast threshold c , the first term from the grey-scale image, the second term from the event sequence, it is possible to compute $\log(\mathbf{L}(f))$, and hence $\mathbf{L}(f)$ by exponentiation. Subsequently, from Eq. (10) the latent image $\mathbf{L}(t)$ at any time may be computed.

When tackling a video, we process each frame in the video separately to generate our reconstructed video. To avoid accumulating errors, it is more suitable to construct each frame $\mathbf{L}(t)$ using the closest blurred frame.

Theoretically, we could generate a video with frame-rate as high as the DVS’s eps (events per second). However, as each event carries little information and is subject to noise, several events must be processed together to yield a reasonable image. We generate a reconstructed frame every $50 - 100$ events, so for our experiments, the frame-rate of the reconstructed video is usually 200 times greater than the input low frame-rate video. Our model (Eq. (11)) can get a better result compared to the state-of-the-art deblurring methods. (see Fig. 9 for an example.)

We put several results of our reconstructed video in the supplementary material and details are given in Table 2. For all our videos, the left part is the input image and the right part is our reconstructed video.

C. Code

All our codes will be released after publication. Here, we only provide a testing code with data ‘snowman’ due to the size limitation of supplementary material (100MB). Please run ‘./snowman_codefortest/maindalta.m’, and manually change the two slide bars to visualize the reconstructed image. ‘Threshold c’ is the threshold to trigger an event. ‘Frame t’ is the reconstructed image at timestamp t . (see Fig. 10).

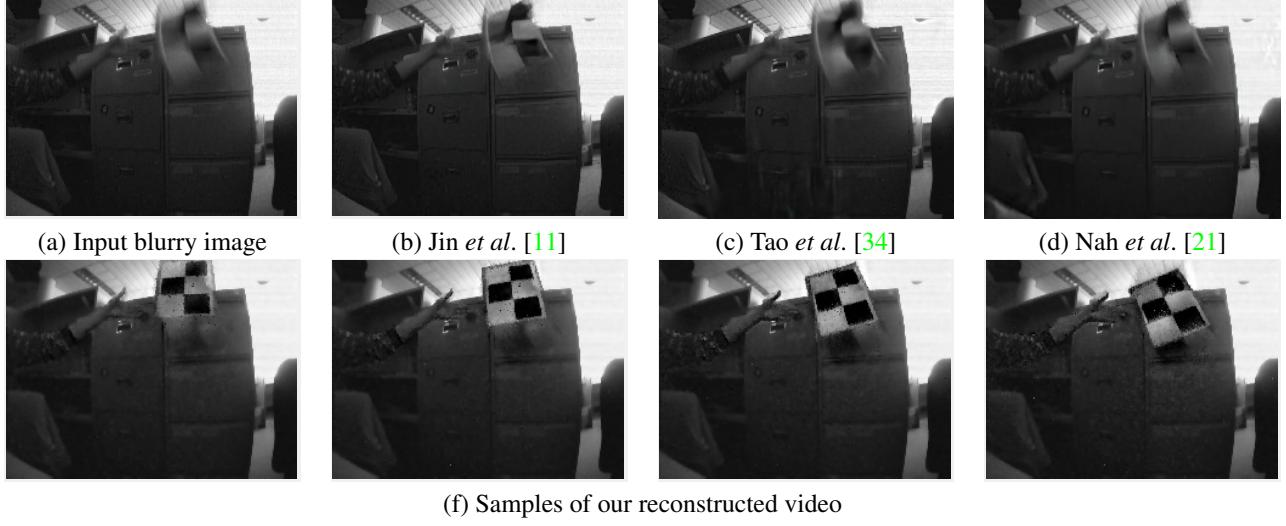


Figure 9. *Deblurring and reconstruction results on our real blurry event dataset.* (a) Input blurry image. (b) Deblurring result of [11]. (c) Deblurring result of [34]. (d) Deblurring result of [21]. (e) Samples of our reconstructed video from $\mathbf{L}(0)$ to $\mathbf{L}(150)$. For baseline 1, we first apply a deblurring method to recover a sharp image, and the recovered image is then fed to a reconstruction method [29]. For baseline 2, we first use a video reconstruction method [29] to construct a sequence of intensity images, and then apply a deblurring method to each intensity image. On synthetic dataset, we use [34] as the deblurring method for its state-of-the-art performance. On our real blurry event dataset, Jin [11] achieves better results than [34] (see Fig. 1 and 8 in our paper). Thus, we use Jin [11] as the deblurring method.

Table 2. *Details of our real blurry event dataset*

Num.	Video name	Condition	Motion type
1	Tea	indoor	camera shake
2	text_SHE	indoor & low light	moving object (The text shown in the phone is also moving)
3	chessboard1	indoor	moving object
4	chessboard2	indoor	moving object
5	jumping	indoor	moving object
6	snowman	indoor	camera shake
7	pillow	indoor	moving object
8	outdoorrun	outdoor & dusk	moving object and camera shake
9	runningman	indoor	moving object
10	lego	indoor	moving objects
11*	nightrun	outdoor & low light	moving object

*:Data from [29].

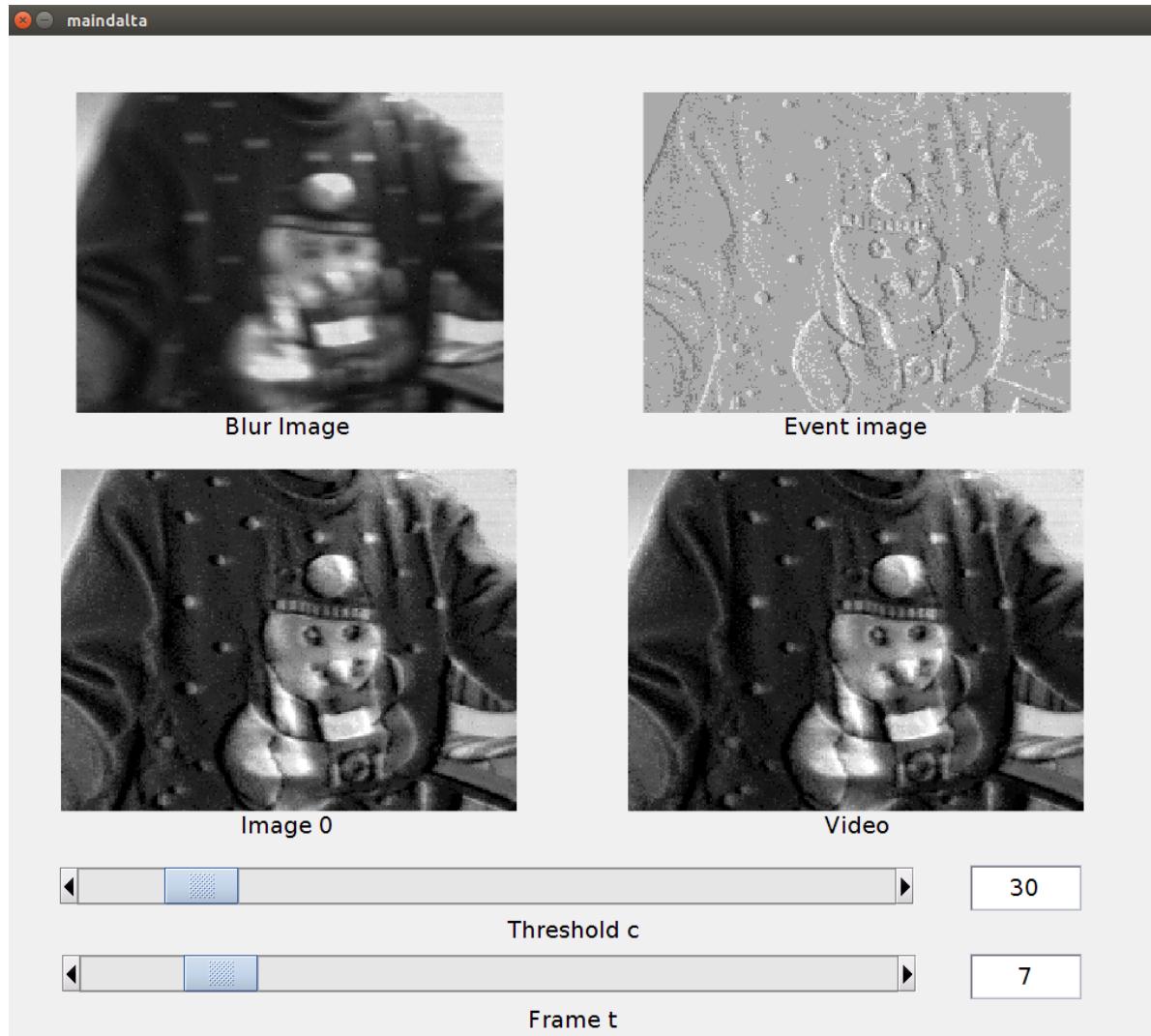


Figure 10. Our program interface.

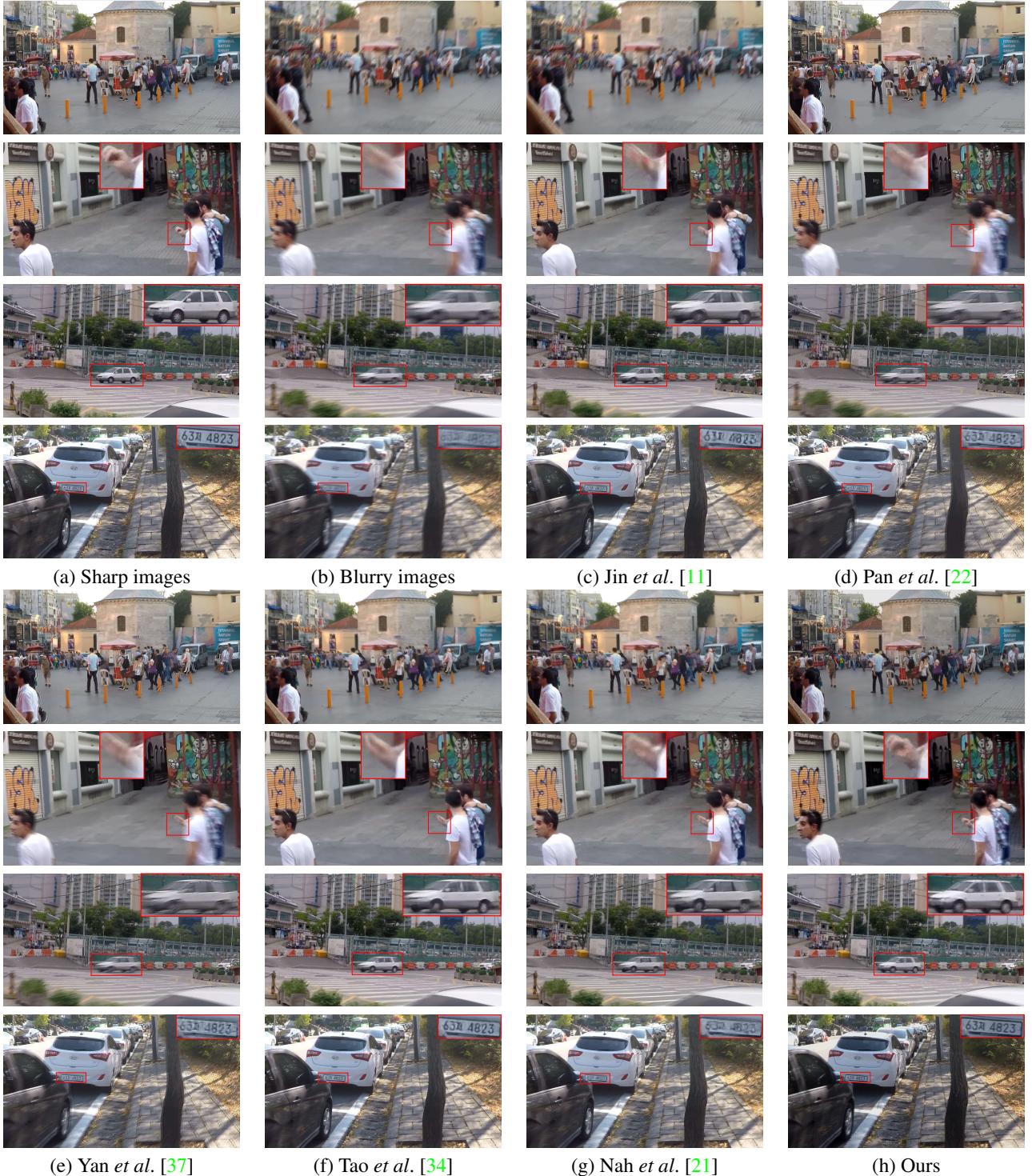


Figure 11. Examples of deblurring results on our synthetic event dataset. (a) Sharp images. (b) Generated blurry images. (c) Deblurring results of [11]. (d) Deblurring results of [22]. (e) Deblurring results of [37]. (f) Deblurring results of [34]. (g) Deblurring results of [21]. (h) Our deblurring results. (Best view in color on screen).