



# High Frame Rate Optical Flow Estimation from Event Sensors via Intensity Estimation

Prasan **Shedligeri**<sup>a,\*\*</sup>, Kaushik **Mitra**<sup>a</sup>

<sup>a</sup>Department of Electrical Engineering, Indian Institute of Technology Madras, Chennai, 600036, India

## ABSTRACT

Optical flow estimation forms the core of several computer vision tasks and its estimation requires accurate spatial and temporal gradient information. However, if there are fast-moving objects in the scene or if the camera moves rapidly, then the acquired images will suffer from motion blur, which will lead to poor optical flow estimation. Such challenging cases can be handled by event sensors which are a novel generation of sensors that acquire pixel-level brightness changes as binary events at a very high temporal resolution. Brightness constancy constraint, which is the basis of several optical flow algorithms cannot be directly used on event sensors making it challenging to estimate optical flow. We overcome this challenge by imposing brightness constancy constraint on intensity images predicted from event sensor data. For this task, we design a recurrent neural network that jointly predicts a sparse optical flow and intensity images from the event data. While intensity estimation is supervised using ground truth frames, optical flow estimation is self-supervised using the predicted intensity frames. However, in our case the temporal resolution of the ground truth intensity frames is far lower than the temporal resolution of the predicted intensity frames, making it challenging to supervise. As we use recurrent neural network, such a challenge can be overcome by sharing the weights for each of the predicted intensity frames. Quantitatively our predicted optical flow is better than previously proposed algorithms for optical flow estimation from event sensors. We also show our algorithm's robustness against challenging cases of fast motion and high dynamic range scenes.

© 2021 Elsevier Ltd. All rights reserved.

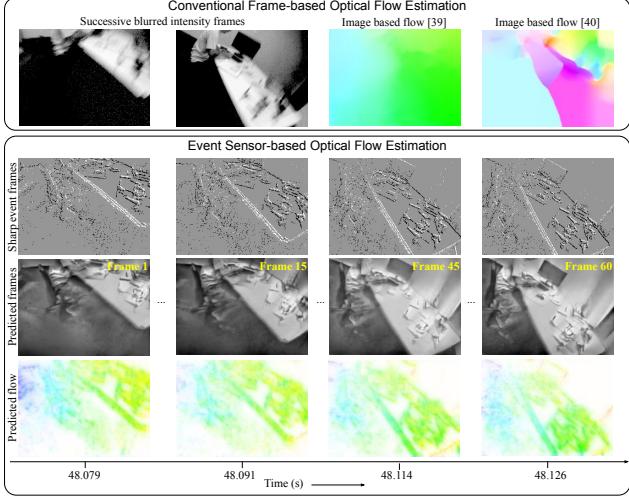
## 1. Introduction

Many of the modern computer vision applications rely on acquiring data from conventional image sensors. Optical flow forms basis for many of the computer vision tasks such as object-tracking, moving object segmentation, autonomous navigation, etc (Fortun et al., 2015). The dense texture rich information acquired from conventional image sensors, enable dense optical flow prediction. The brightness constancy based energy functional introduced by Horn and Schunck (Horn and Schunck, 1981; Fortun et al., 2015) is the basis of many modern optical flow estimation algorithms. This energy functional relies on accurate sensing of image intensities between successive frames. This brightness constancy constraint fails to hold when the acquired images are degraded from motion blur due to fast-moving objects or due to the rapid camera motion as

shown in Fig. 1. Again, due to low frame rate of image sensors, it becomes challenging to estimate optical flow for cases of large scene motion even without significant blur. This challenge can be overcome if we use an image sensor with a very high temporal resolution. Conventional image sensors, that acquire high temporal resolution video are significantly expensive and require large data bandwidth and hence event-based sensors can provide a viable alternative. Event-based sensors are a novel generation of neuromorphic sensors which asynchronously sense only the pixel-level brightness changes with a temporal resolution of the order of microseconds (Delbrück et al., 2010). At each pixel, the event sensor outputs a positive/negative event when it senses an increase/decrease in brightness over a specified threshold. Its extremely high temporal resolution has been demonstrated by reconstructing intensity frames at a frame rate of several thousand frames per second. These sensors also have a much higher dynamic range compared to conventional image sensors making them attractive for several computer vision applications(Gallego et al., 2019).

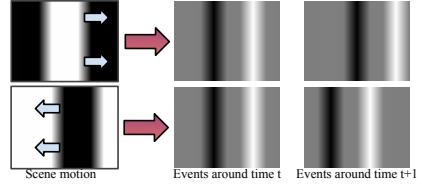
\*\*Corresponding author:

e-mail: ee16d409@ee.iitm.ac.in (Prasan Shedligeri)



**Fig. 1.** Conventional frame-based optical flow algorithms suffer when the input images are degraded with motion blur as shown in the top row. Event sensors on the other hand operate at much higher temporal resolution and can sense much higher dynamic range than the frame-based sensors. We accumulate the events triggered between the two successive intensity images as event frames and show some of them in the second row. Our proposed algorithm takes these intermediate event frames as input and predicts corresponding intensity images and optical flow. In this example, optical flow and intensity images are predicted at 60 intermediate temporal locations corresponding to a 60x temporal super-resolution.

Optical flow estimation directly from event sensors is attractive but a challenging task as the brightness constancy based energy functional cannot be used directly. Despite this challenge, several algorithms have been proposed in the literature for event based optical flow estimation (Liu and Delbrück, 2018; Nagata et al., 2019; Paredes-Vallés et al., 2019; Khoei et al., 2019; Bardow et al., 2016; Zhu et al., 2018c; Haessig et al., 2018; Gallego et al., 2018; Almatrafi and Hirakawa, 2020). While learning based methods have shown significant improvement in optical flow prediction accuracy, they fail to exploit the advantages provided by the event sensor. EV-FlowNet(Zhu et al., 2018a), is one of the first learning-based algorithms proposed to predict optical flow from event sensor data. It used the low dynamic range and low frame rate intensity frames and use the brightness constancy as a supervisory signal, thus ignoring the high dynamic range and high temporal resolution offered by event sensors. In (Zhu et al., 2019), the authors propose to use a contrast maximization framework to estimate optical flow. This is an unsupervised algorithm, where the event sensor data alone is used to supervise optical flow prediction, thus fully utilizing the high dynamic range nature of event sensors. However, this algorithm requires an event volume of 30,000 events as input and hence cannot predict optical flow at very high frame rates. This algorithm also makes a limiting assumption of linear object motion thus affecting the optical flow prediction accuracy. To make full use of the event sensor advantages, algorithms for optical flow estimation from event sensors should have the following desirable properties: (a) the optical flow should be predicted at high temporal resolution, (b) predicted optical flow should be reliable even for challenging high dynamic range scenes, (c) should not require difficult to acquire ground truth optical flow,



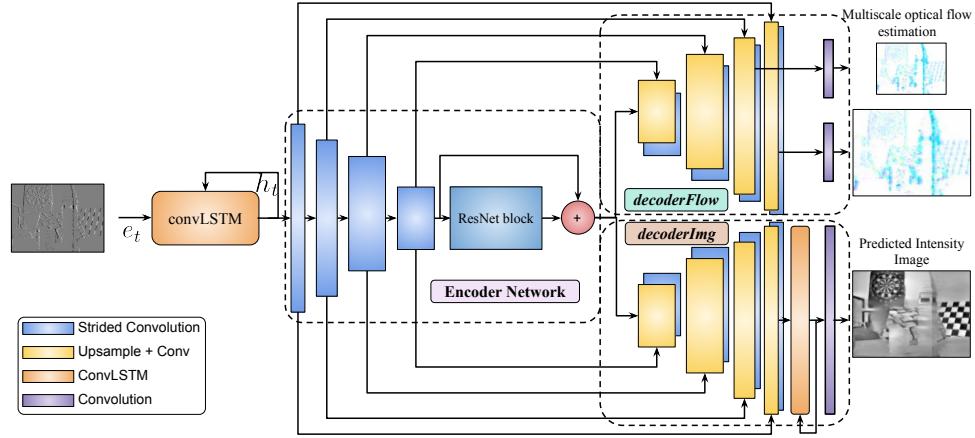
**Fig. 2.** Ambiguity in intensity image prediction from a single event frame. The first column shows two different scenes which have opposite motion with respect to the camera. These two scenes produce the same event frame at time  $t$  making it ambiguous to predict the corresponding scene intensity from the single event frame. However, when we consider the next event frame at time  $t + 1$ , we clearly see the motion in the scene. Modeling this temporal information using recurrent neural network helps in predicting the intensity frames unambiguously from event data alone.

(d) should not make non-generalizable assumptions such as linear motion of the objects.

In our proposed method, intensity frames and a sparse optical flow are simultaneously predicted from the input event sensor data. The event sensor data is first converted to a series of event frames by stacking a fixed number of events per frame following the stacking by number (SBN) principle of (Wang et al., 2019). A sequence of event frames are given as input one-by-one to the neural network which predicts the corresponding intensity frame and optical flow. The intensity frame prediction is supervised using the temporally sparse ground truth intensity frames. While our proposed algorithm predicts intensity frame at a very high temporal resolution (at the rate of incoming events) the frames acquired from hybrid intensity and event based sensors (Brandli et al., 2014) are at a much lower temporal resolution. Thus, it is not possible for us to have a supervised loss for every predicted intensity frame. We overcome this challenge by using recurrent neural network architecture that makes it possible to use supervision only at a few time-steps by sharing weights across all the time-steps. Recurrent neural networks have already been used in (Rebecq et al., 2019b) to predict high frame rate intensity frames. We adapt this network to simultaneously predict intensity frames and optical flow. As demonstrated for optical flow prediction from conventional image sensors (Jason et al., 2016; Ren et al., 2017; Meister et al., 2018), we use the brightness constancy constraint as a supervisory signal for optical flow prediction from event sensors.

In summary, we make the following contributions:

- We propose a semi-supervised learning algorithm to predict high frame rate, sparse optical flow for high dynamic range scenes.
- Optical flow prediction is self-supervised using the high frame rate and high dynamic range intensity frames predicted directly from the event sensor data. Thus, ground truth optical flow is not necessary for training our proposed algorithm.
- We also demonstrate the generalizability of our proposed algorithm on a wide variety of open source event datasets captured with different sensors and in different environments.



**Fig. 3. Overall flow of our proposed method:** Our proposed methods takes in a single event frame at each time-step, which is then input to a ConvLSTM (Convolutional Long-Short Term Memory) network. The updated hidden state from the convLSTM network is input to an encoder network consisting of four strided convolutional layers followed by a ResNet block. The hidden representation from the encoder network is then fed as input to two decoder networks, *decoderImg* and *decoderFlow*, which predict the intensity image and the optical flow, respectively.

## 2. Related Work

**Motion estimation from event sensors:** Although it’s a challenging task to estimate optical flow from event sensors, several algorithms have been proposed (Liu and Delbrück, 2018; Nagata et al., 2019; Paredes-Vallés et al., 2019; Khoei et al., 2019; Bardow et al., 2016; Zhu et al., 2018a, 2019, 2018c; Haessig et al., 2018; Gallego et al., 2018). Works such as (Gallego et al., 2018; Zhu et al., 2018c, 2019) use motion compensation on the space-time volume of events to estimate optical flow. In (Haessig et al., 2018), the authors design a spiking neural network to estimate optical flow and demonstrate their proposed algorithm on IBM’s neuromorphic chip. A few learning based methods have also been proposed for estimating optical flow from event sensors (Zhu et al., 2019, 2018a).

**Intensity image reconstruction:** Previously researchers have attempted to estimate intensity frames from event sensor (Reinbacher et al., 2016; Scheerlinck et al., 2018; Bardow et al., 2016; Shedligeri and Mitra, 2019; Rebecq et al., 2019a; Wang et al., 2019), so that the intensity frames could be used as an input to off-the-shelf frame based computer vision algorithms. Recent learning based algorithms (Rebecq et al., 2019a; Wang et al., 2019) have shown a great improvement in intensity image quality compared to traditional methods. The closest work to ours is (Bardow et al., 2016), where the authors propose a framework to simultaneously estimate intensity and optical flow directly from the event sensor data.

## 3. Optical Flow Estimation from Event Sensors

### 3.1. Modeling events as sequential data

The output of an event sensor is a 4-tuple  $(x, y, t, p)$  where  $x$  and  $y$  represent the spatial location,  $t$  represents the time instant and  $p$  denotes the polarity (+1 or -1) of the triggered event. Following (Wang et al., 2019), we stack these events into a sequence of event frames to form the input to our algorithm. The temporal information is obviously lost due to this projection of

spatio-temporal data as a spatial frame. In Fig. 2, we show a toy example where two different video sequences are used to generate an event frame at time  $t$ . Both the event frames look identical as they lack any temporal information about the events, leading to ambiguity in prediction of intensity frames.

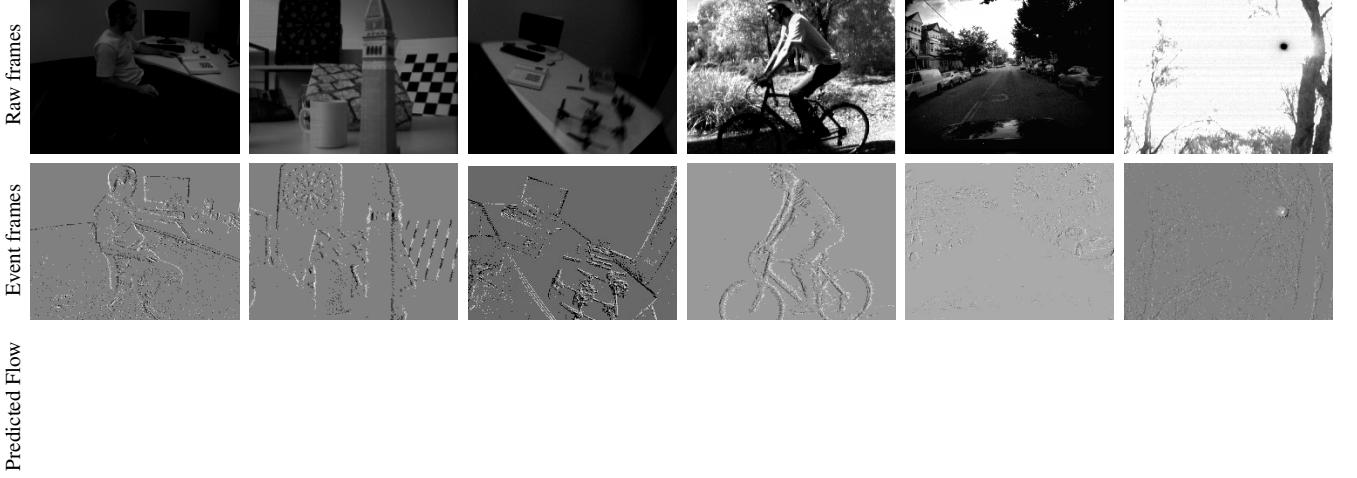
To tackle this loss of temporal information we use a sequence of event frames akin to a sequence of image frames forming a temporal video. The effectiveness of this simple representation can be seen from Fig. 2 where a clear distinction emerges between the two cases of scene motion when considering a video sequence instead of looking at each frame independently. It’s imperative for us to design a neural network that can effectively incorporate this temporal information so as to unambiguously predict the intensity images. LSTM (Long-Short Term Memory) (Gers et al., 1999) networks have been shown to be effective for such tasks and we use them to model the long-term temporal dependency in the sequence of event frames. Although the input to the algorithm at each timestep is a single event frame, the intensity frame is still unambiguously predicted, demonstrating the effectiveness of the proposed LSTM network to model sequential information.

### 3.2. Joint estimation of intensity image and optical flow

Fig. 3 shows our overall model to predict the intensity frames and optical flow from input event sensor data. The intensity frame prediction is supervised using temporally sparse raw intensity images acquired from the conventional image sensor present in DAVIS (Brandli et al., 2014). DAVIS is a hybrid sensor consisting of co-located intensity and event based sensors. The input frames are formed by accumulating events occurring in  $N$  non-overlapping sub-intervals between successive intensity frames. Each of these sub-intervals contain a fixed, pre-determined number of events. These  $N$  event frames are given as input and at the output we obtain the  $N$  intensity frames and corresponding  $N - 1$  optical flow estimates. In the following sections we elaborate on the training algorithm for intensity and the optical flow estimation.







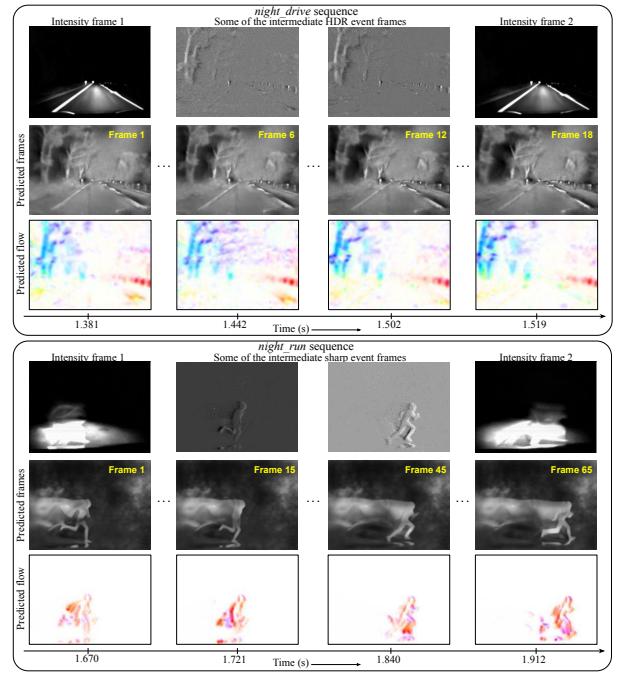
**Fig. 5.** We test our proposed optical flow model for its generalizability on various test sequences obtained from (Mueggler et al., 2017; Scheerlinck et al., 2018; Zhu et al., 2018b). (If document is opened in Adobe Reader, videos can be viewed by clicking on the images). We provide further results in the supplementary video ([link to the video](#)).

Method	indoor flying 1		indoor flying 2		indoor flying 3	
	AEE	% outliers	AEE	% outliers	AEE	% outliers
(Zhu et al., 2018a)	0.83	0.84	1.19	6.75	1.07	4.97
(Zhu et al., 2019)	0.58	0	1.02	4	0.87	3
Ours	<b>0.49</b>	<b>0.02</b>	<b>0.55</b>	<b>0.05</b>	<b>0.53</b>	<b>0.03</b>

**Table 1.** Quantitative comparison of the predicted optical flow on event sequences from (Zhu et al., 2018b).

quences from MVSEC dataset. Following (Zhu et al., 2018a), we choose the metrics (a) Average End-point Error (AEE) which measures the mean absolute error and (b) percentage outliers for quantitative comparison. Percentage outlier (% outlier) measures the percentange of pixels with end-point error above 3 pixels and 5% of the magnitude of the flow vector. For fair comparison, we select two state of the art *unsupervised* learning-based optical flow algorithms (Zhu et al., 2019, 2018a) to benchmark our proposed algorithm. In (Zhu et al., 2018a), all the events between two successive intensity frames are accumulated into a frame based representation and fed to the trained network. In (Zhu et al., 2019), a volume consisting of 30,000 events divided over 10 event frames is fed into the optical flow network. Effectively, each of event frames in (Zhu et al., 2019) is formed by accumulating 3000 events from the event data. For a fair comparison, we too accumulate successive 3000 events into a single event frame which is then sequentially fed to our trained model.

In Table 1 we provide the quantitative metrics to compare our optical flow algorithms with the state of the art methods. We qualitatively compare the optical flow predicted from our model to that of the (Zhu et al., 2018a) in Fig. 4. We show optical flow predicted from various test sequences from datasets proposed by (Scheerlinck et al., 2018; Mueggler et al., 2017) in Fig. 5. Note that these test sequences do not have ground truth optical flow to be compared against. We also provide the video of the predicted optical flow for most of the sequences in the accompanying supplementary video.



**Fig. 6.** The top figure shows the *night\_drive* sequence shot in low-light conditions, demonstrating the ability of event sensors to sense objects at a high dynamic range, allowing the prediction of optical flow in extreme challenging cases. The *night\_run* sequence combines two challenging scenarios, low-light and motion blur. With the help of event sensors we are able to predict the optical flow and intensity images at an effective rate of 1300 frames per second.

#### 4.4. Advantages of event-based optical flow prediction

In this section, we demonstrate the advantages event sensors can provide over conventional image sensors for challenging scenes with fast motion and high dynamic range. In Fig. 1, we show an indoor scene with significant motion blur in the acquired image frames. A significant temporal information has also been lost between the two intensity frames. However, due to the high temporal resolution of the event sensors we are able









used for training and the rest for testing. The training sequences were *boxes\_6dof*, *boxes\_translation*, *boxes\_rotation*, *office\_spiral*, *office\_zigzag*, *outdoors\_running*, *outdoors\_walking*, *poster\_6dof*, *poster\_rotation*, *poster\_translation*, *shapes\_6dof*, *shapes\_translation*, *shapes\_rotation*. The validation and testing sequences were *dynamic\_6dof*, *dynamic\_translation*, *dynamic\_rotation*, *slider\_depth*, *slider\_close*. Video sequences that are captured in similar environments were put in either the training set or the test set, but not both. E.g. sequences such as *boxes* or *dynamic* appear either in the training data or in the test data, but not in both.

We quantitatively evaluate our proposed optical flow algorithm with the ground truth optical flow available in the dataset proposed by Zhu *et al.* (Zhu et al., 2018b). This dataset, also known as *MVSEC* dataset (Zhu et al., 2018b), consists of event sequences captured using a available hybrid sensor named DAVIS346. This sensor has a spatial resolution of  $260 \times 346$ . The dataset also contains ground truth depth maps captured using a LiDAR and the relative 6-DoF pose captured using a motion-capture system. For quantitative evaluation, the ground truth optical flow is also provided in the dataset. The ground truth optical flow is computed using the ground truth depth and the relative 6-DoF camera pose (Zhu et al., 2018a) under the assumption of a static scene. As the authors assume static scene to compute ground truth optical flow, we exclude the *outdoor\_driving* sequence from optical flow evaluation and only use the *indoor\_flying* sequences. This is because the *outdoor\_driving* sequence has many moving objects such as cars, pedestrians, etc. where the static scene assumption does not hold and hence affects the optical flow evaluation. For a fair evaluation, we follow the procedure defined in (Zhu et al., 2018a) to compute metrics for optical flow evaluation by computing the error only at pixels where an event has fired.

### C. More qualitative comparison on intensity image prediction

Intensity image prediction from event sensor data has been investigated by many researchers in the past few years. In this work we predict intensity images in order to facilitate the learning of optical flow directly from event sensors. However, in order to predict a good estimate of the optical flow, the predicted intensity frames should be temporally consistent, have high dynamic range and be free of any noise or other artifacts.

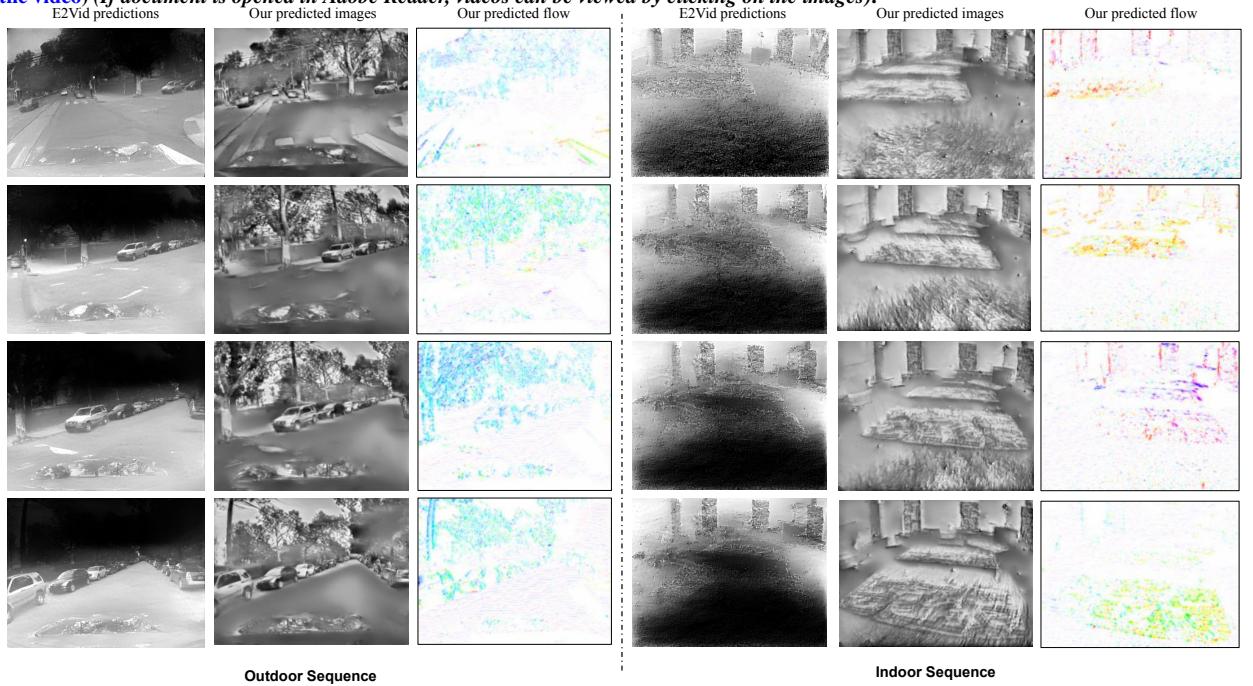
In Fig. S.1 of supplementary material we provide qualitative comparison of the intensity images from our proposed method

with that of the other state-of-the-art intensity image-only estimation algorithms such as (Scheerlinck et al., 2018; Rebecq et al., 2019a; Reinbacher et al., 2016). Images predicted from (Reinbacher et al., 2016) contains various artifacts such as trailing edges as it relies on the hand-crafted image prior based on event manifolds. The algorithm proposed in (Scheerlinck et al., 2018) is sensitive to noisy events and can be seen producing noisy intensity estimates when a large number of events are being triggered in the scene. This sensitivity arises from the lack of any spatial regularization in the complementary filter model proposed in (Scheerlinck et al., 2018). In (Rebecq et al., 2019a), the authors propose a learning based method to predict intensity images from event sensor data. The method proposed in (Rebecq et al., 2019a) demonstrate high quality intensity image prediction from event sensor data. In contrast to MR (Reinbacher et al., 2016) and CF (Scheerlinck et al., 2018) which produce images with trailing edge artifacts, our reconstructions are smooth and are free of most of the artifacts. As seen in Fig. S.1, the predicted intensity frames from our model are comparable to the state of the art, learning-based, intensity only reconstruction method E2Vid (Rebecq et al., 2019a). The comparisons for all the methods in Fig. S.1 can be better visualized in the accompanying supplementary video ([link to the video](#)).

The learning based algorithm proposed in (Rebecq et al., 2019a) eliminates most of the artifacts dominant in the event based intensity reconstruction such as bleeding edges. However, in some cases the reconstructed intensity images start to show a dark region as shown in Fig. S.2. We also observe that the dark region on the reconstructed images grow and occupy larger area as more frames are reconstructed. The images from (Rebecq et al., 2019a) were reconstructed by accumulating  $N_e$  events into an event voxel-grid consisting of 5 frames. We consider  $N_e = 0.35 \times H \times W$ , where  $H$  and  $W$  represent the sensor resolution and this has shown to produce impressive results in most cases. The scenes shown in Fig. S.2 are from the *indoor\_flying* and the *outdoor\_driving* sequences from (Zhu et al., 2018b). These sequences are acquired with a hybrid sensor with a sensor resolution of  $H = 260$ ,  $W = 346$ . As our algorithm requires only a single event frame as input, we accumulate  $N_e/5$  events into each event frame to predict the intensity frames. In Fig. S.2, we provide the qualitative comparison for the predicted images from the two sequences *indoor\_flying* and the *outdoor\_driving*.



**Fig. S.1.** Qualitative comparison of intensity frame reconstruction on various event sensor sequences. We see that the reconstructed intensity images from our method do not have trailing edges and noisy regions as compared to (Reinbacher et al., 2016; Scheerlinck et al., 2018). Ours and (Rebecq et al., 2019a), both use learning based intensity reconstruction method and produce comparable results which can be better seen from the supplementary video material. ([link to the video](#)) (*If document is opened in Adobe Reader, videos can be viewed by clicking on the images*).



**Fig. S.2.** We show two sequences, one outdoor and another indoor, never seen by our method or the one proposed in (Rebecq et al., 2019a). We see that the images predicted by (Rebecq et al., 2019a) degrade with a growing dark region in the predicted intensity images in this particular case. Our proposed method generalizes enough to provide a reliable estimate of the intensity image and the optical flow.

