# Learning Dense and Continuous Optical Flow from an Event Camera

Zhexiong Wan, Yuchao Dai, *Member, IEEE,* and Yuxin Mao

*Abstract*—**Event cameras such as DAVIS can simultaneously output high temporal resolution events and low frame-rate intensity images, which own great potential in capturing scene motion, such as optical flow estimation. Most of the existing optical flow estimation methods are based on two consecutive image frames and can only estimate *discrete flow* at a fixed time interval. Previous work has shown that *continuous flow* estimation can be achieved by changing the quantities or time intervals of events. However, they are difficult to estimate reliable *dense flow*, especially in the regions without any triggered events. In this paper, we propose a novel deep learning-based dense and continuous optical flow estimation framework from a single image with event streams, which facilitates the accurate perception of high-speed motion. Specifically, we first propose an event-image fusion and correlation module to effectively exploit the internal motion from two different modalities of data. Then we propose an iterative update network structure with bidirectional training for optical flow prediction. Therefore, our model can estimate reliable dense flow as two-frame-based methods, as well as estimate temporal continuous flow as event-based methods. Extensive experimental results on both synthetic and real captured datasets demonstrate that our model outperforms existing event-based state-of-the-art methods and our designed baselines for accurate dense and continuous optical flow estimation.**

*Index Terms*—**Event camera, event-based vision, optical flow estimation, multimodal learning.**

## I. INTRODUCTION

**E**VENT cameras are bio-inspired vision sensors that can trigger brightness change asynchronously and independently at each pixel with a microsecond time resolution [1]. Unlike the conventional frame-based shutter cameras that capture full resolution images at a fixed frame rate, event cameras such as DVS [2] and DAVIS [3] can output a discrete event stream at a very small and not fixed event rate. In particular, the DAVIS event camera [3] can simultaneously output image and event streams. The event data stream can be regarded as a frame sequence with up to millions of frames-per-second (fps) [4]–[6], which owns appealing advantages over the shutter frames, including high temporal resolution, high dynamic range, low latency, low redundancy, and low power consumption. These enable the broad applications of event cameras in feature tracking [7]–[9], depth estimation and
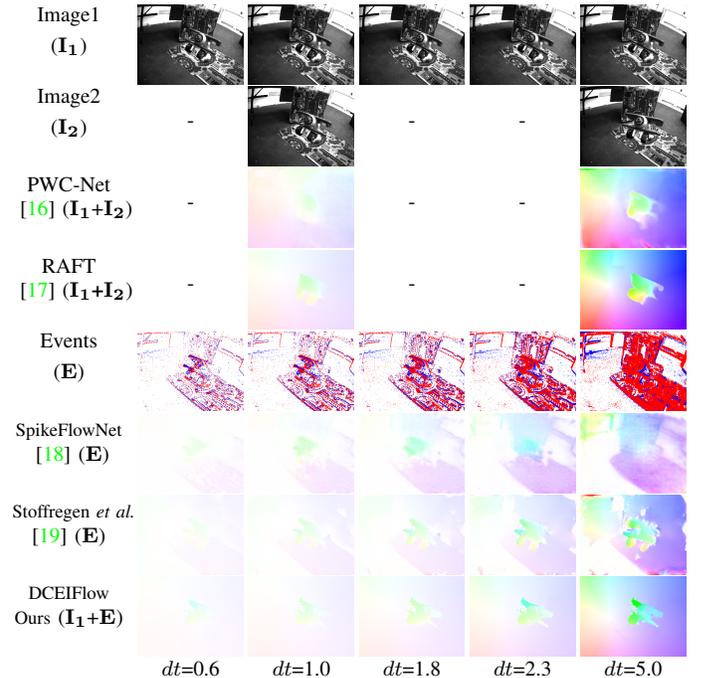
Fig. 1. **Visual comparisons of continuous flow prediction with different time intervals**. $dt$ denotes the frame interval (*e.g.*, $dt$=1.0 represents the time interval between two adjacent frames). The two-frame ($I_1$+$I_2$) approaches can only estimate dense optical flow between frames. The event-only ($E$) approaches can estimate continuous optical flow, but cannot predict accurate dense flow. Our model can estimate both dense and continuous optical flow by fusing events and the first image ($I_1$+$E$). Best viewed on screen.

3D reconstruction [10]–[13], frame synthesis [4], [14], [15], *etc*.

Optical flow estimation aims to predict the motion between two moments by exploiting the photometric consistency. Most of the existing event-based optical flow estimation approaches [18], [20]–[23] only use event streams. Although temporal continuous optical flow can be predicted, it is difficult to get reliable predictions in regions without any events, as shown in Fig. 1. Thus, we consider fusing a single image with events to improve the reliability of dense optical flow estimation. Due to the frame rate limitation of the first image, we cannot estimate the continuous flow from any start to end time like the event-only methods. However, we can still estimate continuous and reliable dense flow at varying time intervals from a fixed frame. We have shown a schematic diagram in Fig. 2 to illustrate the differences between these three types of input settings. The estimated continuous flow from a single image with events has notable significance for many event-based downstream applications, especially those associated with images, including image deblurring [5], video
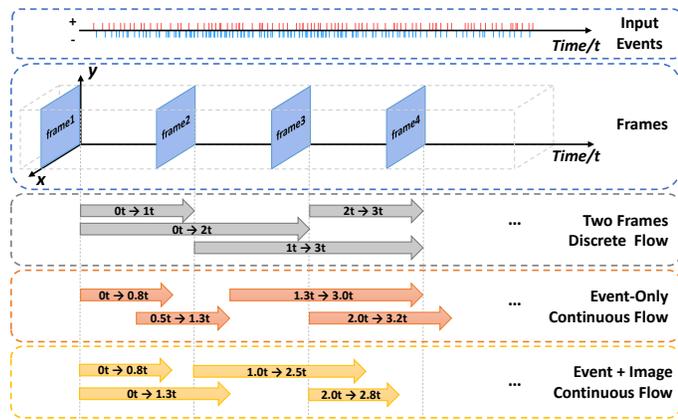
Fig. 2. Differences between **Discrete Flow** and **Continuous Flow**. Given a set of frames and event streams captured by a DAVIS camera. The two-frame optical flow estimation methods can only work with discrete integer frame intervals. The event-only approach can estimate continuous flow from any start to end time. Bringing the first image with events, we can still estimate continuous flow, but the start time is limited to the time of each frame.

synthesis [24], feature tracking [8], *etc.*

In this paper, we propose a novel deep learning-based optical flow estimation model from a single image with event streams, named DCEIFlow, which can effectively exploit the internal relation from two different modalities of data through our proposed event-image fusion and correlation module. Thus, our proposed DCEIFlow model can estimate reliable dense flow as two-frame-based methods, as well as estimate continuous flow as event-only methods. In particular, we first propose an event-image fusion module to effectively fuse the features of the first image and events by a multi-layer convolution fusion network. The fused feature is regarded as a pseudo second image feature that is constrained by the real second image feature in training. Then we use it to construct the feature correlation. On this basis, we propose a network with an iterative update structure to learn the optical flow from the constructed correlations. In addition, we propose a bidirectional flow training mechanism based on the reverse event, which can use the same network to estimate the backward flow by inputting the second image with the reversed event streams during the training stage.

Due to the lack of an event dataset with dense flow annotation, we first pre-train our model on the two-frame dataset FlyingChairs2 [25] with simulated events, and then evaluate the pre-trained model on the real captured dataset MVSEC [26]. Evaluation results and visual comparisons on both synthetic and real captured datasets show a significant improvement over the existing event-only or fused single image state-of-the-arts. Specifically, our pre-trained model achieves better results on MVSEC than existing methods with different time intervals ($dt$=1 and $dt$=4 frames). We also compare our model with two baseline networks, one with only input events and the other directly concatenating the image feature with events. The results show that our model achieves better results with a smaller model size than the baseline methods. In addition, we perform visual comparisons on a highly dynamic real captured dataset EV-IMO [27]. The results show that our model is superior to existing event-based methods in dense and accurate flow estimation and has advantages over existing two-frame-

based methods for detailed optical flow estimation.

Our main contributions are summarized as follows:

(1) We propose a novel deep learning-based dense and continuous optical flow estimation model from a single image with event streams, which can estimate reliable dense flow as the two-frame-based approaches, as well as estimate temporal continuous flow as the event-only approaches.

(2) We propose to build an event-image correlation to effectively exploit internal motion from two different modalities of data. We also propose bidirectional flow training based on reverse events to leverage the order of motion information in events.

(3) Extensive experiments on the MVSEC [26] and EV-IMO [27] datasets demonstrate that our proposed DCEIFlow model improves significantly compared to baselines and the existing event-based state-of-the-arts. We also verify the superiority of our method in dense and continuous optical flow estimation through further experiments and analysis.

## II. RELATED WORK

Optical flow estimation is a very active research area in computer vision, where various approaches have been proposed. In this section, we first review the development of learning-based two-frame optical flow estimation. Then we focus on event-based optical flow estimation, including dense and sparse flow from events or with a single image.

### A. Two-Frame-based Optical Flow Estimation

Recently, the success of deep Convolutional Neural Networks (CNNs) has been extended to various computer vision tasks such as optical flow estimation [28]. The first end-to-end CNN regression approach for estimating optical flow is FlowNet [29], which directly estimates flow from a pair of input images based on an encoder-decoder architecture. It has achieved a faster inference speed than optimization-based methods with higher accuracy. PWC-Net [16] exploits three well-known design principles from existing optimization approaches to deep learning scenarios, including pyramid structure, feature warping, and correlation construction. These key design principles have been widely used or improved in a series of recent works such as IRR-PWC [30], SelFlow [31], and VCN [32]. Recently, RAFT [17] introduced an all-pairs correlation iterative network structure and achieved significant improvements over existing methods.

### B. Event-based Optical Flow Estimation

Since event streams only encode the pixel-level brightness changes discretely, they cannot directly represent the absolute brightness. It is difficult to find the spatial photometric consistency between sparse pixels and estimate dense optical flow. According to the input and output, the existing event-based methods can be divided into the following three categories.

*1) Sparse Flow from Events:* Before deep networks were widely used, Benosman *et al.* [33] first proposed an event-based optical flow algorithm based on the Lucas-Kanade [34] brightness constancy assumption. However, it can only estimate the normal flow component perpendicular to the edge because the event data is usually triggered at the moving edge. After that, [22], [35] can estimate full flow, which introduce tangential flow and contain more motion information compared to normal flow [1]. Recently, [23] proposes to use SNNs [36] to estimate sparse flow efficiently, but it is not widely used because of the limited application of sparse flow.

*2) Dense Flow from Events:* Recent event-based methods tend to estimate dense optical flow, which can provide more spatial information than sparse flow. EV-FlowNet [20] is an end-to-end optical flow network learning from events in a self-supervised manner. It uses the grey image captured by DAVIS as the unsupervised supervision in training. After that, Zhu *et al.* [21] proposed an unsupervised training framework by using the predicted flow to remove the motion blur in the input events. EST [37], Matrix-LSTM [38] explore different event representations. SpikeFlowNet [18], LIF-EV-FlowNet [39], STE-FlowNet [40], E-RAFT [41] and Li *et al.* [42] explore the effects of introducing SNNs, recurrent structure and reducing network parameters. Because they only use events, the predicted dense flow in the regions without any triggered events, such as the constant brightness region, is relatively unreliable compared to the regions with events.

*3) Dense Flow from Single Image with Events:* In order to obtain more reliable dense estimates, researchers consider combining events with the image which contains per-pixel absolute intensity. Bardow *et al.* [43] jointly reconstructed the intensity image and estimated flow from events, but the accuracy of flow depends on the image reconstruction quality. Pan *et al.* [44] proposed to jointly use a set of events with a single image and introduced an event-based brightness constancy as the objective function for optimization. However, these non-convex optimization-based methods are not only time-consuming but also require complex post-processing and tuning. Very recently, Fusion-FlowNet [45] directly inputs events and an image to an end-to-end dual-branch fusion network. This concatenate fusion scheme is simple to explore the internal relationship between two modalities of data, and their results can not show the advantages of introducing the image.

## III. APPROACH

Given the event streams and first image, we build a learning-based framework for estimating dense and continuous optical flow. Our framework consists of five stages: (1) event volume representation, (2) event and image feature extraction, (3) event-image feature fusion, (4) event-image all-pairs correlation construction, and (5) iterative flow updater. Based on this, we propose a bidirectional optical flow training mechanism to constrain the training process.

### A. Preliminaries

*1) Event camera model:* The output data streams of the event camera can be regarded as a finite quaternion sequence,

which represents the per-pixel brightness changes during a period of time. Each event contains its space-time coordinate and a binary polarity representing its brightness change:

$$e_i = \{\mathbf{x}_i, t_i, p_i\}, \tag{1}$$

$$|\log(L(\mathbf{x}, t + \Delta t)) - \log(L(\mathbf{x}, t)| \geq c, \tag{2}$$

where $L(\mathbf{x}_i, t_i)$ is the brightness at camera coordinate $\mathbf{x}_i = (x_i, y_i)$ and microsecond timestamp $t_i$. When the logarithmic domain brightness changes reach the threshold $c$ after $\Delta t$ time, an event $e$ is triggered. The polarity $p = \pm 1$ indicates the direction of brightness change.

*2) **Discrete Flow** and **Continuous Flow**:* Within the standard setup, the optical flow is estimated from two image frames to represent the displacement of the corresponding pixels from the first frame to the second frame. Because the images from the shutter camera are usually at a fixed frame rate, two-frame-based methods can only estimate **discrete flow** with discrete integer frame intervals. Therefore, it has a limited ability to explain the motion with a high temporal resolution, such as the motion within frames. Here, we make a comparison for these different settings in Fig. 2. Note that the discrete flow in our paper is defined in the temporal domain, and it is different from the discrete flow defined in the previous two-frame-based approaches [46], [47], which represents estimating optical flow by discrete optimization.

Therefore, it is difficult to estimate the **continuous flow** with variable time intervals under the two-frame setting. In contrast, the event cameras capture the brightness change of each pixel at high time resolution. It can represent the reliable internal motion within frames with high temporal resolution. Thus, we can use events to estimate the continuous flow with the theoretical frame rate as high as the event camera's *eps* (events per second), which is very helpful for high-speed motion estimation or video analysis.

*3) Feature Correlation for Matching:* Most recent two-frame optical flow estimation networks use correlation to represent the pixel level matching similarity of two image features from a siamese encoder. The correlation is usually constructed by calculating the dot product similarity between the feature of the first image and the feature of the second image warped by the coarse flow.

Here, we review the general construction of local correlation in the two-frame setting [16], [48]. For the given two feature maps $P_{I_1}$ and $P_{I_2}$ (generated by two input images through a siamese encoder) and their corresponding optical flow $\boldsymbol{F}^{1 \to 2}$, the correlation volume are constructed as:

$$\begin{aligned} C_I(\mathbf{x}, \boldsymbol{\delta}_{uv}) &= P_{I_1}(\mathbf{x}) \cdot Warp\{P_{I_2}, \boldsymbol{F}^{1 \to 2}\}(\mathbf{x} + \boldsymbol{\delta}_{uv}) \\ &= P_{I_1}(\mathbf{x}) \cdot Warp\{P_{I_2}(\mathbf{x} + \boldsymbol{\delta}_{uv}), \boldsymbol{F}^{1 \to 2}(\mathbf{x} + \boldsymbol{\delta}_{uv})\} \\ &= P_{I_1}(\mathbf{x}) \cdot P_{I_2}(\mathbf{x} + \boldsymbol{F}^{1 \to 2}(\mathbf{x} + \boldsymbol{\delta}_{uv}) + \boldsymbol{\delta}_{uv}), \end{aligned} \tag{3}$$

where $\cdot$ is the vector dot product, $\boldsymbol{F}^{1 \to 2}$ is the forward optical flow. $\boldsymbol{\delta}_{uv} = (\delta_u, \delta_v) \in ([-d_u, d_u], [-d_v, d_v])$ represents the horizontal and vertical search range, the values of $d_u, d_v$ are two manually set hyper-parameters determines the 2D range of built correlation. The backward wrapping operation $Warp$ is implemented with bilinear interpolation and can compute the gradients for backpropagation [49].
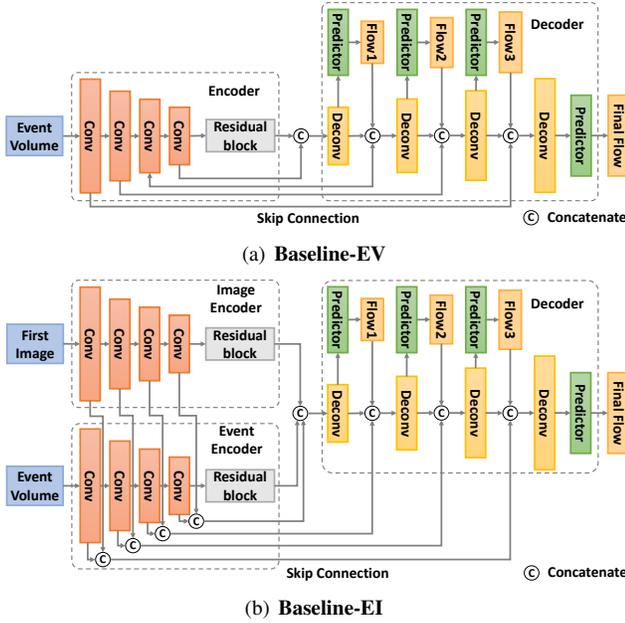
Fig. 3. **The network structure of two baselines.** *Baseline-EV* (a) is a U-Net-like network that only inputs the event volume. *Baseline-EI* (b) adds an image encoder to get the image features and concatenates it with event features as the flow decoder input.

*4) Baselines:* Most of the existing deep learning-based event-based flow methods employ a U-Net-like network, *e.g.*, EV-FlowNet [21]. So we adopt it as a baseline network (called *Baseline-EV* in Fig. 3(a)) using our event representation, training, and evaluation pipeline. Inspired by FusionFlowNet [45], we also introduce another simple network (called *Baseline-EI* in Fig. 3(b)) modified upon EV-FlowNet to estimate flow from a single image and events. We add another encoder to get the image feature and concatenate it with event features as the decoder input. The structures of these two baseline models are shown in Fig. 3.

### B. Event Representation

Since the original event streams are composed of a series of discrete events, following the setting in [4], [19], [21], [50], we aggregate it into a three-dimensional event volume as the input of the convolution network. This process preserves most of the spatial-temporal information in the original event streams.

For an event stream $(e_i)^N, i \in [0, N]$ with $N$ events, we divide it into $B$ temporal bins as the channel dimension of an event volume for each polarity, then sum the normalized timestamps at different pixel positions in each bin as below:

$$E(b, \mathbf{x}_i, p_i) = \sum_{i=0}^{N} \max\left(0, 1 - \left| b - \frac{t_i - t_{start}}{t_{end} - t_{start}}(B-1) \right|\right), \quad (4)$$

where $b \in [0, B)$ indicates the index of temporal bins, $t_{start}$ and $t_{end}$ are the start and end timestamps of the event streams, respectively. Finally, we concatenate these temporal bins by two polarities to an event volume with $(2B \times H \times W)$:

$$\underset{2B\,channels}{V(\mathbf{x})} = [\underset{B\,channels}{E(\mathbf{x}, p = 1)}, \underset{B\,channels}{E(\mathbf{x}, p = -1)}], \quad (5)$$

where $[\ ,\ ]$ is the concatenate operation and we perform it on the channel dimension.

According to [15], [19], [51], we divide the event stream into $B$=5 temporal bins for two polarities in our experiments, then the shape of represented event volume is $(10 \times H \times W)$.

### C. Dense Iterative Event-Image Flow Network

After representing the original discrete events as an event volume, we first extract the event and image feature using the feature extractor. Then we propose an event-image feature fusion module and construct the event-image correlation to effectively exploit the internal motion from two different modalities. After that, we adopt the event-image correlation to the iterative flow update structure for accurate dense optical flow estimation. Overall, our iterative network structure is shown in Fig. 4, which consists of four modules: event and image feature extractor, event-image feature fusion, event-image all-pairs correlation module, and residual flow updater.

**Event and Image Feature Extractor**. We use two convolutional encoders with the same structure but without sharing weights to extract the image and event features. Each encoder consists of 6 residual blocks with three times downsampling performed by convolution with a stride of 2. Then the encoders extract the input event volume size from $(H \times W \times 2B)$ to feature size $(H/8 \times W/8 \times C)$, and image size from $(H \times W \times 3)$ to $(H/8 \times W/8 \times C)$, where $C = 256$ is the channel size of feature maps. In addition, the image encoder also extracts features from the input second frame image during training.

**Event-Image Feature Fusion**. The correlation plays an important role in the two-frame optical flow estimation, thus we extend it to events. However, within our setting, the first image and event features are different data types, we cannot directly use the basic construction in Eq. 3. Therefore, we propose to fuse the single image feature with the event feature to build the correlation, as shown in Fig. 5. We use the motion contained in the event feature to establish the conversion relationship with the first image, and generate pseudo second frame image features. Then we construct the correlation of these two image features as the input of the flow updater.

We first consider the simple addition operation, the *Fusion by Add* structure. This method directly adds the event feature $P_E$ and the first image feature $P_{I_1}$ to obtain the pseudo second image feature $P_{pesudo}$.

$$P_{pesudo} = EIF_{add}(P_{I_1}, P_E) = P_{I_1} + P_E. \quad (6)$$

Since the modality between the two features is different, direct addition is not an intuitive choice. Therefore, we propose *Fusion by Convolutions* as shown in Fig. 5. There are three convolution layers in this module. The first two are used to encode the first image feature $P_{I_1}$ and event feature $P_E$, respectively. The last one is used for fusion by concatenating operation. The final pseudo second image feature $P_{pesudo}$ is obtained by residual addition.

$$\begin{aligned} P_{pesudo} &= EIF_{conv}(P_{I_1}, P_E) \\ &= Conv3([Conv1(P_{I_1}), Conv2(P_E)]) + P_{I_1}. \end{aligned} \quad (7)$$

We propose a feature similarity loss $L_{sim}$ to supervise the similarity between the fused pseudo second image feature with the real second image feature. We finally choose *Fusion by*
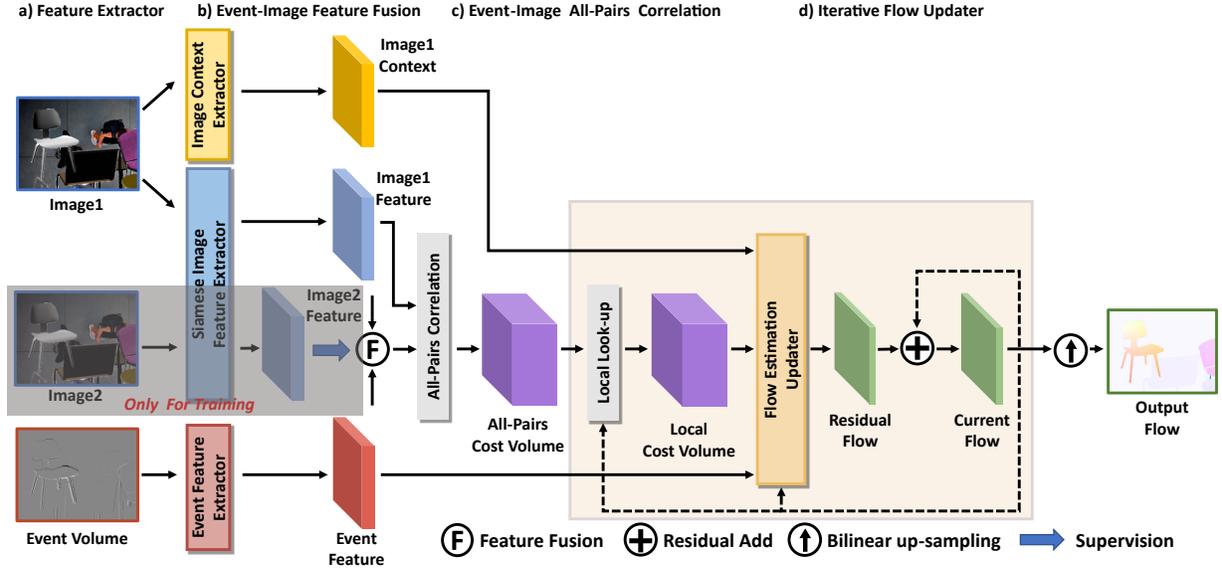
Fig. 4. **Our DCEIFlow model structure.** We use the feature extractor (left) to obtain event and image features and compute the matching correlation using our proposed event-image fusion and correlation construction module (middle). Then we feed them into the iterative flow updater (right) to update the flow iteratively. After the last iteration, we apply the up-sample operation to get the full-resolution output. The structures enclosed by the orange box need to be iteratively updated.
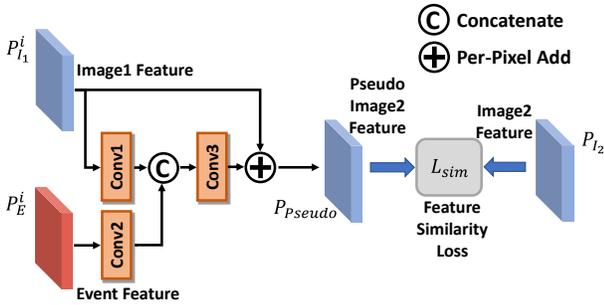


Fig. 5. The *Fusion by Convolutions* structure of our Event-Image Fusion **module**. The feature similarity loss $L_{sim}$ is only used in network training.

*Convolutions* for better performance. *Fusion by Add* is used for ablation studies IV-D3.

**Event-Image All-Pairs Correlation Module**. We construct the all-pairs correlation to enlarge the previous local path size described in III-A3 to full feature size. For the given first image feature $P_{I_1}$ and fused pseudo second image feature $P_{pseudo}$, the all-pair correlation $C^0_{EI}$ can be obtained by calculating the matrix multiplication:

$$C^0_{EI}(\mathbf{x}, \mathbf{y}) = P_{I_1}(\mathbf{x}, c)P^T_{pesudo}(\mathbf{y}, c), \quad (8)$$

where $c$ is the channel of the feature map, $\mathbf{x}, \mathbf{y}$ are the coordinate vectors of the two features.

For the input feature map with size $(H \times W \times C)$, the constructed correlation size is $(H \times W \times H \times W)$. We also introduce the pyramid correlation construction by three times average pooling to involve both large and small search ranges.

$$C^k_{EI}(\mathbf{x}, \mathbf{y}') = \frac{1}{2^{2k}} \sum_{\mathbf{q}}^{(2^k, 2^k)} C^0_{EI}\left(\mathbf{x}, 2^k \times \mathbf{y}' + \mathbf{q}\right), \quad (9)$$

where $\mathbf{y}'$ is the pooled coordinates of the last two dimensions, $k \in [1, 3]$ is the pyramid level, and $k=0$ means the original

correlation. Thus the size of each correlation is $(H \times W \times H/2^k \times W/2^k)$.

Finally, we perform the lookup operation by the coarse flow in the defined search range $(\delta_u, \delta_v)$ on each correlation.

$$LC^k_{EI}(\mathbf{x}, \boldsymbol{\delta}_{uv}) = C^k_{EI}\left(\mathbf{x}, \frac{\mathbf{x} + \boldsymbol{F}^{1 \to 2}(\mathbf{x})}{2^k} + \boldsymbol{\delta}_{uv}\right). \quad (10)$$

The size of $k$ level local correlation $LC^k_{EI}$ is $(H \times W \times (2 \times d_u + 1) \times (2 \times d_v + 1))$. We merge the last two dimensions and concatenate all of the pyramids $\{LC^0_{EI}, LC^1_{EI}, LC^2_{EI}, LC^3_{EI}\}$, then feed into the iterative flow updater.

**Iterative Flow Updater**. The iterative flow updater consists of a ConvGRU (Convolutional Gated Recurrent Unit [52]) and several convolution layers, which can estimate the residual flow $\Delta F$ from the concatenation of image and event features, as well as the pyramid correlation volume. At each iteration, the residual flow $\Delta F$ output by the flow updater is used to update the estimated flow $F$. The flow updater iterates $N$ times with the correlation lookup operation. The $i$-th updated optical flow $F^i$ is the sum of the previous and the current estimations: $F^i = F^{i-1} + \Delta F$.

Because the resolution of the extracted feature map is reduced to $1/8$ of the input image, the predicted optical flow needs to be upsampled to the input size. We use $8\times$ bilinear upsampling to the updated flow after the last iteration as the final predicted optical flow.

*D. Bidirectional Training*

Bidirectional flow training has been widely used in two-frame-based approaches [30], [31], which shows that the network with shared weight can estimate the forward and backward flow by exchanging the input order of two images. For event-based flow estimation, we can know the order of each event from the event timestamps. If we reverse the order of event timestamps and reverse the polarity of brightness
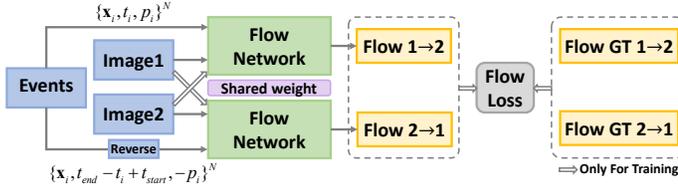
Fig. 6. **Our bidirectional Event-Image flow training framework.** It takes the first image and original events as input and outputs the forward flow. Meanwhile, the second image and reversed events are taken as input to output the backward flow. The ground truth forward and backward flow are used as supervision for training.

changes, the motion becomes reversed. Thus the estimated flow can be considered equivalent to the inverse motion.

Therefore, we propose our bidirectional flow training mechanism in Fig. 6. If we train our model on the datasets with bidirectional flow annotation, we can input the original and reversed events and the corresponding image to the network, respectively, and use the ground-truths to supervise the output forward and backward flow. However, it should be noted that the backward flow requires the second image as input, so we cannot obtain the backward flow using the first image during the inference stage. Our proposed bidirectional training mechanism is helpful to get more accurate flow results and improve the generalization ability without adding any network parameters, which has been verified in our ablation studies.

### E. Training Loss

*1) Flow Loss:* Assume the ground truth flow is $F^{gt}$, the predicted flow at the $i$-th iteration is $F^i$, where $i = 1, 2, ..., N$, $N$ is the total number of iterations. Then the predicted flows are supervised by using the following flow loss $L_f$:

$$L_f = \sum_{i=1}^{N} \phi^{N-i+1} \rho(\| F^i - F^{gt} \|_2), \qquad (11)$$

where the robust function $\rho(x) = (x^2 + \epsilon)^q$, $q \in (0, 1)$ is less sensitive to outliers, $\epsilon$ is a small number which is close to 0. $\phi$ is a hyper-parameter used to balance the loss weights of each prediction. In our experiments, we set $N = 6$ to balance the computation cost and performance and $\phi = 0.8$ to make the later predictions with bigger weights.

For bidirectional training, we define the bidirectional flow loss $L_{fb}$ as:

$$L_{fb} = \frac{1}{2} \cdot \sum_{i=1}^{N} \phi^{N-i+1} \big[ \rho(\| F_{1\rightarrow2}^i - F_{1\rightarrow2}^{gt} \|_2) \\ + \rho(\| F_{2\rightarrow1}^i - F_{2\rightarrow1}^{gt} \|_2) \big], \qquad (12)$$

where the ground truth forward and backward flow are $F_{1\rightarrow2}^{gt}$ and $F_{2\rightarrow1}^{gt}$, the predicted flow are $F_{1\rightarrow2}^i$ and $F_{2\rightarrow1}^i$.

*2) Feature similarity Loss:* In our event-image fusion module, the feature similarity loss $L_{sim}$ is used to supervise the pseudo second image feature $P_{pesudo}$ similar to the real second image feature $P_{I_2}$. Thus, we use $L_{sim}$ to compute $L_2$ distance between these two features:

$$L_{sim} = \| P_{I_2} - P_{pesudo} \|_2 . \qquad (13)$$

*3) Total Training Loss:* The total training loss is a weighted sum of those two losses. When training on the dataset with both forward and backward flow annotations, such as FlyingChairs2 [25], the bidirectional loss $L_{bi}$ is used. When training on the dataset with the only forward flow, such as MVSEC [26], the unidirectional loss $L_{un}$ is used.

$$L_{bi} = L_{fb} + \lambda \cdot L_{sim}, \\ L_{un} = L_f + \lambda \cdot L_{sim}. \qquad (14)$$

In our experiments, the feature similarity loss $L_{sim}$ can quickly converge to a small order of magnitude, so we set $\lambda = 100.0$ to balance the losses.

### IV. EXPERIMENTS

In this section, we first introduce our implementation details, including datasets, simulation, training details, and evaluation metrics. Then, we show the evaluation results of our model on both simulated and real datasets with comparisons to several baselines and existing methods. We further prove the effectiveness of each component in our network and the advantages of our network in dense and continuous optical flow estimation by model analysis. We conclude with discussions on failure cases and the limitations of our model.

### A. Implementation details

*1) Datasets:*
*- Selection:* The commonly used event camera optical flow dataset is MVSEC [26]. However, only using it to train our network is not a good practice because it only has sparse flow annotation with low spatial resolution. Most existing learning-based two-frame flow approaches usually pre-train on synthetic datasets and then finetune on other datasets to get benchmark results. Thus we use ESIM [56] to simulate the event data between two frames on the FlyingChairs2 dataset [25] to pre-train our model, because it provides full ground truth annotation of forward and backward flow.

The flow annotations of MVSEC are computed from the depth by LIDAR with the ego-motion by IMU, and there are only rigid scenes. Therefore, to verify the performance of our model in non-rigid dynamic scenes, we use an event-based highly dynamic moving object segmentation dataset, EV-IMO [27]. In addition, we also use the Sintel dataset [57] because it is commonly used in the two-frame methods.

*- Details:* Following the split of Stoffregen *et al.* [19], each indoor and outdoor sequence in MVSEC [26] contains 1,880∼ and 2,700∼ images with corresponding events, respectively. FlyingChairs2 [25] contains 22,232 training and 640 validation samples. Each sample includes two image pairs, forward & backward flow annotations, and our simulated event data. Sintel [57] provides naturalistic movie sequences with challenging long-range non-rigid motion, which includes Clean and Final passes with 1,041 pairs of training sets. Because it has no event data, and the flow annotations of the test set are not publicly available, we only simulate the event data on the training set to evaluate the generalization ability. The test set of EV-IMO [27] dataset includes 21 sequences, with a total of 8258 pairs of data captured by the DAVIS346C camera. Due to the lack of
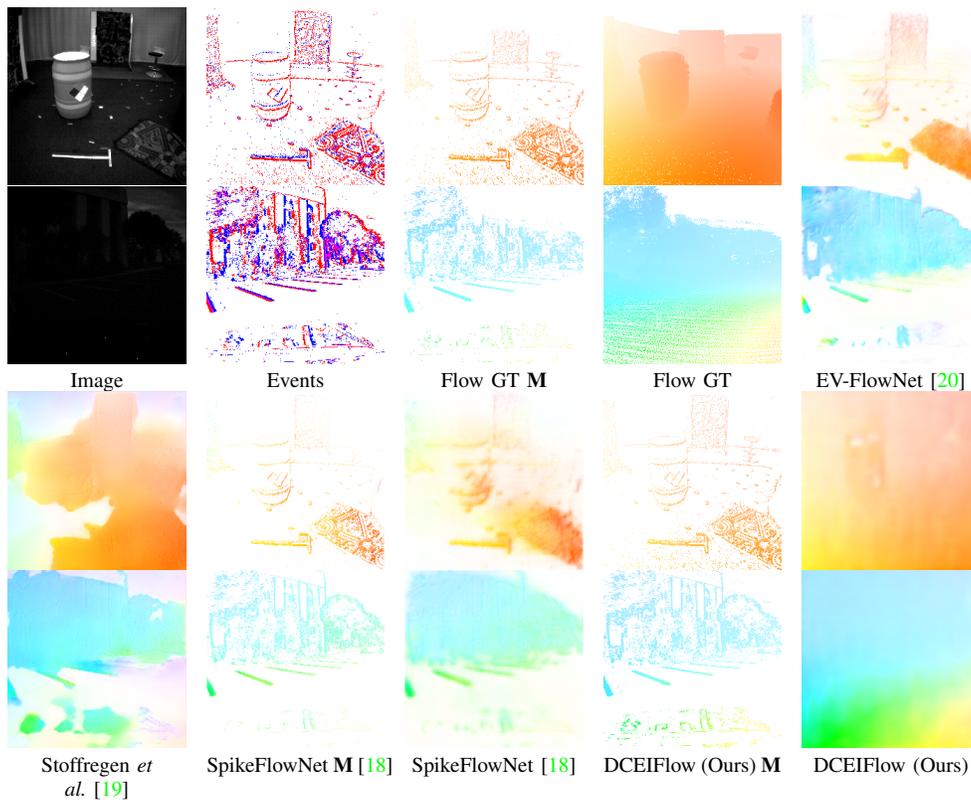
Fig. 7. **Visual comparisons on the MVSEC dataset [26]. M** is the masked flow at the pixels with events. Our model gets better visual results in both indoor (top) and outdoor scenes (bottom). Best viewed on screen.

optical flow annotation, we only use its image and event data to estimate optical flow for qualitative visual analysis.

*2) Event data simulation:* Due to the lack of an event-based dataset with high-quality optical flow annotations for training, we use the open-source ESIM simulator [56] to simulate events on FlyingChairs2. To simulate realistic events, ESIM requires a small displacement of the corresponding pixels between two frames. However, the pixel displacement of this dataset is not guaranteed to be always small, so we cannot directly input it into the simulator with the original two frames. Following Gehrig *et al.* [58], we first use Super-SloMo [59] to interpolate the two frames to more, and then use ESIM to simulate events on them. The amount of interpolating frames depends on the motion range between two frames.

*3) Model training details:* We train our model on the FlyingChairs2 training set by $L_{bi}$ in Eq. (14) for 100 epochs with a random cropped size $[368, 496]$ and a batch size of 8. Our model needs 25 minutes for one epoch, and it takes about 42 hours to complete the whole training process for 100 epochs. We train our model on two NVIDIA 2080Ti GPUs using PyTorch [60]. We use the AdamW optimizer [61] for training with a weight decay of $10^{-4}$ and default parameters $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$. We use several geometric and photometric augmentations, including random resize and crop, horizontal and vertical flips, contrast and brightness changes, etc. We use the OneCycle [62] policy and set the maximal learning rate to $4 \times 10^{-4}$. After model training, we directly use the same pre-trained model to evaluate on MVSEC under different time intervals ($dt = 1$ and $dt = 4$).

In the baseline comparison, we also conduct the experiments of training only on the *outdoor_day2* sequence of the MVSEC

dataset [26] (*i.e.*, M). Each baseline and our model follow the same pre-training settings on the MVSEC dataset (Train D.Set M). Both are trained for 200 epochs with a batch size of 16. Our model needs about 2 minutes for one epoch, and it takes about 10 hours to complete the training process for 300 epochs. Then we use the same pre-trained model to obtain the results of each method under two input interval settings ($dt = 1$ and $dt = 4$ frames).

*4) Evaluation metrics:* A commonly used metric for optical flow evaluation is the average End Point Error (**EPE**), which calculates the Euclidean distance between the predicted flow and the ground truth.

$$EPE = \frac{1}{m} \cdot \sum_m \sqrt{(\boldsymbol{F}_x^{pred} - \boldsymbol{F}_x^{gt})^2 + (\boldsymbol{F}_y^{pred} - \boldsymbol{F}_y^{gt})^2}. \quad (15)$$

For *dense* evaluation, $m$ is the pixels with valid flow annotation. For *sparse* evaluation, $m$ is the pixels with valid flow annotation and triggered at least one event. $\boldsymbol{F}^{pred}$ is the predicted flow vector and $\boldsymbol{F}^{gt}$ is ground-truth flow vector, the $x$ and $y$ subscripts indicate horizontal and vertical directions. Following KITTI [63] and EV-FlowNet [20], we also use the **outlier** metric to report the percentage of points with endpoint error greater than 3 pixels and 5% of the magnitude. Both of these two metrics are smaller the better.

In addition, we introduce a metric called **Dense Ratio** to measure the output density of an event-based optical flow estimation model, which is defined as follows,

$$\textbf{Dense Ratio} = \frac{\textbf{EPE}_{Dense} + \textbf{EPE}_{Event\ Masked}}{\textbf{EPE}_{Dense} + \textbf{EPE}_{Event\ Excluded}}, \quad (16)$$

where $\textbf{EPE}_{Dense}$ calculates the dense error of pixels annotated by the valid optical flow, $\textbf{EPE}_{Event\ Masked}$ calculates

TABLE I

PERFORMANCE EVALUATION ON THE MVSEC DATASET [26] COMPARED WITH EXISTING EVENT-BASED METHODS. THE RESULTS OF THE COMPARED METHODS ARE DIRECTLY EXTRACTED FROM THE ORIGINAL PAPERS. NOTE THAT EXISTING EVENT-BASED METHODS USUALLY TRAIN TWO SEPARATE MODELS FOR DIFFERENT TIME INTERVALS ($dt=1$ AND $dt=4$ FRAMES), BUT OUR RESULTS ARE OBTAINED ON THE SAME PRE-TRAINED MODEL.

| Input $dt=1$ | Method Reference | Train Mann. | Train D.Type | Train D.Set | Eval. Metric | indoor_flying1 EPE | %Out | indoor_flying2 EPE | %Out | indoor_flying3 EPE | %Out | outdoor_day1 EPE | %Out | outdoor_day2 EPE | %Out |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | EV-FlowNet [20] | USL | $I_1,I_2,E$ | M | sparse | (1.03) | (2.2) | (1.72) | (15.1) | (1.53) | (11.9) | [0.49] | [0.2] | × | × |
| | Zhu et al. [21] | USL | E | M | sparse | (0.58) | (0.0) | (1.02) | (4.0) | (0.87) | (3.0) | [0.32] | [0.0] | × | × |
| | EST [37] | SL | E | M | sparse | (0.97) | (0.91) | (1.38) | (8.20) | (1.43) | (6.47) | - | - | × | × |
| | Matrix-LSTM [38] | USL | $I_1,I_2,E$ | M | sparse | (0.82) | (0.53) | (1.19) | (5.59) | (1.08) | (4.81) | - | - | × | × |
| | Spike-FlowNet [18] | USL | $I_1,I_2,E$ | M | sparse | [0.84] | - | [1.28] | - | [1.11] | - | [0.49] | - | × | × |
| E | Stoffregen et al. [19] | SL | E | ESIM | dense | **0.56** | 1.00 | 0.66 | 1.00 | 0.59 | 1.00 | 0.68 | 0.99 | 0.82 | **0.96** |
| | Paredes et al. [51] | USL | E | M | sparse | (0.79) | (1.2) | (1.40) | (10.9) | (1.18) | (7.4) | [0.92] | [5.4] | × | × |
| | LIF-EV-FlowNet [39] | USL | E | FPV | sparse | 0.71 | 1.41 | 1.44 | 12.75 | 1.16 | 9.11 | **0.53** | 0.33 | - | - |
| | Deng et al. [53] | USL | $I_1,I_2,E$ | M | sparse | (0.89) | (0.66) | (1.31) | (6.44) | (1.13) | (3.53) | - | - | × | × |
| | Li et al. [42] | USL | $I_1,I_2,E$ | M | sparse | (0.59) | (0.83) | (**0.64**) | (2.26) | - | - | **[0.31]** | [0.03] | × | × |
| | STE-FlowNet [40] | USL | $I_1,I_2,E$ | M | sparse | [0.57] | [0.1] | [0.79] | [1.6] | [0.72] | [1.3] | [0.42] | **[0.0]** | × | × |
| $I_1+I_2$ +E | Fusion-FlowNet [45] | USL | $I_1,I_2,E$ | M | dense | (0.62) | - | (0.89) | - | (0.85) | - | [1.02] | - | × | × |
| | Fusion-FlowNet [45] | USL | $I_1,I_2,E$ | M | sparse | (0.56) | - | (0.95) | - | (0.76) | - | [0.59] | - | × | × |
| | Pan et al. [44]* | MB | - | - | M | 0.93 | 0.48 | 0.93 | 0.48 | 0.93 | 0.48 | 0.93 | 0.48 | - | - |
| $I_1+E$ | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | dense | **0.56** | **0.28** | **0.64** | **0.16** | **0.57** | **0.12** | 0.91 | 0.71 | **0.79** | 2.59 |
| | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | sparse | 0.57 | 0.30 | 0.70 | 0.30 | 0.58 | 0.15 | 0.74 | **0.29** | 0.82 | 2.34 |

\* Only the average EPE and outlier results of four sequences are given in Pan et al. [44].

| Input $dt=4$ | Method Reference | Train Mann. | Train D.Type | Train D.Set | Eval. Metric | indoor_flying1 EPE | %Out | indoor_flying2 EPE | %Out | indoor_flying3 EPE | %Out | outdoor_day1 EPE | %Out | outdoor_day2 EPE | %Out |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | EV-FlowNet [20] | USL | $I_1,I_2,E$ | M | sparse | (2.25) | (24.7) | (4.05) | (45.3) | (3.45) | (39.7) | [1.23] | [7.3] | × | × |
| | Zhu et al. [21] | USL | E | M | sparse | (2.18) | (24.2) | (3.85) | (46.8) | (3.18) | (47.8) | [1.30] | [9.7] | × | × |
| E | Spike-FlowNet [18] | USL | $I_1,I_2,E$ | M | sparse | [2.24] | - | [3.83] | - | [3.18] | - | [1.09] | - | × | × |
| | LIF-EV-FlowNet [39] | USL | E | FPV | sparse | 2.63 | 29.55 | 4.93 | 51.10 | 3.88 | 41.49 | 2.02 | 18.91 | - | - |
| | Li et al. [42] | USL | $I_1,I_2,E$ | M | sparse | (2.08) | (26.4) | (3.76) | (43.2) | - | - | [1.24] | [8.16] | × | × |
| | STE-FlowNet [40] | USL | $I_1,I_2,E$ | M | sparse | [1.77] | [14.7] | [2.52] | [26.1] | [2.23] | [22.1] | [0.99] | [3.9] | × | × |
| $I_1+I_2$ +E | Fusion-FlowNet [45] | USL | $I_1,I_2,E$ | M | dense | (1.81) | - | (2.90) | - | (2.46) | - | [3.06] | - | × | × |
| | Fusion-FlowNet [45] | USL | $I_1,I_2,E$ | M | sparse | (1.68) | - | (3.24) | - | (2.43) | - | [1.17] | - | × | × |
| $I_1+E$ | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | dense | **1.49** | 8.14 | **1.97** | 17.37 | **1.69** | 12.34 | 1.87 | 19.13 | 1.62 | 14.73 |
| | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | sparse | 1.52 | 8.79 | 2.21 | 22.13 | 1.74 | 13.33 | **1.37** | **8.54** | **1.61** | **14.38** |

- indicates these methods do not provide the corresponding results.

× indicates these methods are trained on the *ourdoor_day2* sequence.

The results with ( ) are obtained by evaluating the model trained on both *ourdoor_day1* and *ourdoor_day2* sequences.

The results with [ ] are obtained by evaluating the model trained on the *ourdoor_day2* sequence.

The results not enclosed by any brackets indicate that the model is not trained on any sequence of MVSEC.

the sparse error of pixels with the valid flow and triggered at least one event, and $\mathbf{EPE}_{Event\ Excluded}$ calculates the sparse error of pixels which do not trigger any event. Because events are usually triggered at moving objects or texture edges, we measure the dense flow prediction ability by calculating the ratio of the masked and excluded EPE. When the ratio is less than 1, the model has a smaller error at the edges than at other locations. The underlined results in Table V are obtained using unsupervised pre-trained models. Since the unsupervised objective function usually has higher energy at the edges, this result illustrates that unsupervised training tends to make the model fit to the edges.

### B. Results on the MVSEC dataset

Most event-based optical flow estimation methods report results on the MVSEC [26] dataset. Therefore, we compare our model pre-trained on the FlyingChairs2 dataset [25] (*i.e.*, FC2) with them in Table I. Existing event-based methods usually train two different models to separately evaluate on different time intervals ($dt$=1 or $dt$=4 frames) respectively, but we only use one model for these two settings. Note that the results in Shedligeri et al. [64] and Mostafavi et al. [65] are calculated with a fixed number of events ($dt = 15000$ or $dt = 30000$ events). So we do not compare with them because we follow the commonly used fixed frame intervals setting. We also show the results from four two-frame-based methods only as a reference comparison in Table II, including two supervised methods (PWC-Net [16] and RAFT [17]) and two unsupervised methods (ARFlow [54] and SMURF [55]).

We annotate the training manners for each method in Table I. In the column of training manners (**Train Mann.**), **SL** represents that the method is trained in a supervised manner, while **USL** represents an unsupervised manner. In addition, we annotate the data types (**Train D.Type**) used in the training process for each method. For the USL methods, the data type with $I_1,I_2,E$ indicates that the loss function of this method is

TABLE II
EXTENDED EVALUATION RESULTS ON THE MVSEC DATASET [26] COMPARED WITH EXISTING STATE-OF-THE-ART TWO FRAME-BASED METHODS
AND EVENT-BASED METHOD E-RAFT [41]. WE EVALUATE THE RESULTS OF THE COMPARED METHODS ON THEIR OPEN-SOURCE PRE-TRAINED
MODELS BY USING THE SAME DATASET SPLITTING (FOLLOWING [19]). WE USE THE SAME PRE-TRAINED MODEL TO GET THE RESULTS OF EACH
METHOD UNDER TWO INPUT INTERVAL SETTINGS ($dt=1$ AND $dt=4$ FRAMES).

| Input $dt=1$ | Method Reference | Train Mann. | Train D.Type | Train D.Set | Eval. Metric | indoor_flying1 EPE | %Out | indoor_flying2 EPE | %Out | indoor_flying3 EPE | %Out | outdoor_day1 EPE | %Out | outdoor_day2 EPE | %Out |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | PWC-Net [16] | SL | $I_1,I_2$ | C+T | dense | 1.57 | 3.11 | 1.62 | 3.29 | 1.55 | 2.70 | 1.83 | 11.50 | 1.67 | 7.88 |
| | PWC-Net [16] | SL | $I_1,I_2$ | C+T | sparse | 1.59 | 3.27 | 1.69 | 4.80 | 1.58 | 3.10 | 1.75 | 7.68 | 1.64 | 8.39 |
| $I_1+I_2$ | RAFT [17] | SL | $I_1,I_2$ | C | dense | 0.44 | 0.13 | 0.54 | 0.05 | 0.50 | 0.00 | 0.86 | 0.21 | 0.62 | 2.81 |
| | RAFT [17] | SL | $I_1,I_2$ | C | sparse | 0.48 | 0.12 | 0.62 | 0.07 | 0.54 | 0.00 | 0.77 | 0.21 | 0.56 | 2.37 |
| As reference | ARFlow [54] | USL | $I_1,I_2$ | SR+S | dense | 0.39 | 0.13 | 0.46 | 0.07 | 0.43 | 0.02 | 1.44 | 12.29 | 0.86 | 7.19 |
| | ARFlow [54] | USL | $I_1,I_2$ | SR+S | sparse | 0.38 | 0.11 | 0.48 | 0.05 | 0.41 | 0.00 | 0.89 | 5.44 | 0.70 | 4.49 |
| | SMURF [55] | USL | $I_1,I_2$ | C | dense | 0.42 | 0.14 | 0.50 | 0.27 | 0.46 | 0.15 | 1.50 | 12.90 | 0.95 | 8.79 |
| | SMURF [55] | USL | $I_1,I_2$ | C | sparse | 0.39 | 0.11 | 0.50 | 0.09 | 0.43 | 0.04 | 0.99 | 6.37 | 0.72 | 5.09 |
| E | E-RAFT [41] | SL | E | DSEC | dense | 0.70 | **0.16** | 0.94 | 2.97 | 0.82 | 1.48 | 0.95 | 4.55 | 1.04 | 6.47 |
| | E-RAFT [41] | SL | E | DSEC | sparse | 0.78 | 0.33 | 1.20 | 5.70 | 0.93 | 2.25 | **0.65** | 2.19 | 0.92 | 4.73 |
| $I_1+E$ | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | dense | **0.56** | 0.28 | **0.64** | **0.16** | **0.57** | **0.12** | 0.91 | 0.71 | **0.79** | 2.59 |
| | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | sparse | 0.57 | 0.30 | 0.70 | 0.30 | 0.58 | 0.15 | 0.74 | **0.29** | 0.82 | **2.34** |
| Input $dt=4$ | Method Reference | Train Mann. | Train D.Type | Train D.Set | Eval. Metric | indoor_flying1 EPE | %Out | indoor_flying2 EPE | %Out | indoor_flying3 EPE | %Out | outdoor_day1 EPE | %Out | outdoor_day2 EPE | %Out |
| | PWC-Net [16] | SL | $I_1,I_2$ | C+T | dense | 1.94 | 14.35 | 2.19 | 21.01 | 2.03 | 17.06 | 3.03 | 37.99 | 2.33 | 19.52 |
| | PWC-Net [16] | SL | $I_1,I_2$ | C+T | sparse | 1.96 | 14.95 | 2.31 | 24.60 | 2.05 | 17.46 | 2.48 | 26.62 | 2.28 | 19.44 |
| $I_1+I_2$ | RAFT [17] | SL | $I_1,I_2$ | C | dense | 1.45 | 7.85 | 1.80 | 13.89 | 1.65 | 11.02 | 3.10 | 38.77 | 1.43 | 12.50 |
| | RAFT [17] | SL | $I_1,I_2$ | C | sparse | 1.48 | 7.82 | 1.91 | 15.94 | 1.67 | 11.29 | 2.47 | 27.15 | 1.39 | 11.65 |
| As reference | ARFlow [54] | USL | $I_1,I_2$ | SR+S | dense | 1.31 | 6.21 | 1.58 | 9.51 | 1.44 | 8.05 | 3.43 | 39.55 | 1.53 | 13.06 |
| | ARFlow [54] | USL | $I_1,I_2$ | SR+S | sparse | 1.31 | 6.59 | 1.72 | 12.05 | 1.47 | 8.81 | 1.89 | 17.36 | 1.42 | 11.50 |
| | SMURF [55] | USL | $I_1,I_2$ | C | dense | 1.34 | 6.80 | 1.63 | 10.38 | 1.49 | 8.76 | 3.98 | 46.49 | 1.73 | 15.49 |
| | SMURF [55] | USL | $I_1,I_2$ | C | sparse | 1.32 | 6.76 | 1.73 | 12.27 | 1.50 | 9.17 | 2.53 | 25.65 | 1.41 | 11.57 |
| E | E-RAFT [41] | SL | E | DSEC | dense | 1.82 | 15.58 | 2.64 | 25.47 | 2.12 | 17.60 | 1.93 | 19.55 | 1.66 | 14.05 |
| | E-RAFT [41] | SL | E | DSEC | sparse | 1.89 | 16.41 | 3.22 | 33.23 | 2.27 | 19.81 | 1.43 | 9.17 | **1.59** | **11.83** |
| $I_1+E$ | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | dense | **1.49** | **8.14** | **1.97** | **17.37** | **1.69** | **12.34** | 1.87 | 19.13 | 1.62 | 14.73 |
| | DCEIFlow (Ours) | SL | $I_1,I_2,E$ | C2 | sparse | 1.52 | 8.79 | 2.21 | 22.13 | 1.74 | 13.33 | **1.37** | **8.54** | 1.61 | 14.38 |

based on the warping of the APS images, while the data type with E indicates that it is based on the warping of events. In particular, for the model-based optimization method Pan *et al.* [44], we annotate it as **MB**. The model-based methods usually do not need pre-training but need to manually adjust the hyper-parameters for different input data.

As shown in Table I, our model achieves state-of-the-art performance on indoor sequences for both EPE and outlier metrics. Although we do not get the best EPE on the *ourdoor_day*1 sequence, the outlier metric is significantly lower than others, and the performance on *indoor*1-3 sequences is good enough to verify our advantage. Especially for the longer time and larger motion evaluation setting with $dt$=4, our results have been greatly improved compared with others. Moreover, as a dense optical flow estimation method using a single image with events, our performance is comparable to the existing two-frame-based SOTAs. This superior performance shows the effectiveness of our framework in fusing the first image and events for accurate dense flow prediction. As shown in Table II, we use the pre-trained model of E-RAFT [41] on DSEC [66] to compare with our model pre-trained on FlyingChairs2 [25]. The DSEC dataset used for pre-training E-RAFT is real captured on outdoor vehicles, while the FlyingChairs2 dataset we used is simulated with multiple chairs superimposed on a random image. Therefore, the performance

of E-RAFT is comparable to our DCEIFlow model when evaluated on the outdoor sequences of the MVSEC dataset, while our model performs significantly better in the indoor sequences due to the introduction of image data. Moreover, using a single 2080ti GPU, our model only takes 28ms to process data with MVSEC size and get flow prediction, while E-RAFT takes 62ms.

In addition, we make visual comparisons with several event-based methods, which have open-sourced models in Fig. 7. EV-FlowNet [20] and SpikeFlowNet [18] only use events, and their claim is to predict sparse flow. Their visualizations are also sparse and include many incorrect predictions (such as the upper left corner of the first sample in Fig. 7). For another event-only method Stoffregen *et al.* [19], although its claim is dense prediction, it is difficult to predict a complete dense flow in the area without events. Most of the motion in the MVSEC dataset is caused by the camera, and the ground-truth optical flow labels are calculated from sparse depth and camera motion. Thus there are spatial mismatches between images and optical flow in some scenes, which also increases the difficulty of visual comparison. Despite the slight mismatches, we can still conclude that our proposed DCEIFlow model produces not only fewer errors but also more dense estimations. This is consistent with the conclusion of the above quantitative comparison.

TABLE III
ABLATION STUDY RESULTS. BOTH MODELS ARE PRE-TRAINED ON THE FLYINGCHAIRS2 [25] TRAINING SET WITH THE SAME TRAINING SETTING, AND DIRECTLY EVALUATED ON THE FLYINGCHAIRS2 VALIDATION SET, SINTEL [57] TRAINING SET AND MVSEC [26] *indoor_flying1-3* SEQUENCES WITH $dt = 1$.

| Model ID | Event Pol. | Corr. Module | E-I Fusion | Sim Loss | Bi. Train. | Network Structure | Param. Num. (M) | FlyingChairs2 EPE | %Out | Sintel EPE | %Out | MVSEC EPE | %Out |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| (a) | × | × | × | × | × | Pyramid | 10.88 | 2.01 | 12.00 | 8.96 | 45.33 | 1.16 | 2.58 |
| (b) | ✓ | × | × | × | × | Pyramid | 10.88 | 1.97 | 11.71 | 8.52 | 42.41 | 1.05 | 2.64 |
| (c) | ✓ | ✓ | Add | ✓ | × | Pyramid | 12.28 | 1.90 | 11.36 | 8.12 | 38.41 | 0.88 | 1.49 |
| (d) | ✓ | ✓ | Add | ✓ | × | Iterative | 6.08 | 1.85 | 10.51 | 7.69 | 36.42 | 0.76 | 0.72 |
| (e) | ✓ | ✓ | Conv | ✓ | × | Pyramid | 13.27 | 1.83 | 10.79 | 7.55 | 36.01 | 0.74 | 0.93 |
| (f) | ✓ | ✓ | Conv | ✓ | × | Iterative | 7.07 | 1.66 | 8.92 | 7.01 | 34.51 | 0.62 | 0.40 |
| (g) | ✓ | ✓ | Conv | ✓ | ✓ | Pyramid | 13.27 | 1.80 | 10.23 | 7.21 | 35.29 | 0.68 | 0.81 |
| (h) | ✓ | ✓ | Conv | × | ✓ | Iterative | 7.07 | 1.74 | 9.20 | 7.50 | 35.46 | 0.73 | 0.89 |
| (i) | ✓ | ✓ | Conv | ✓ | ✓ | Iterative | 7.07 | **1.58** | **7.88** | **6.47** | **32.23** | **0.59** | **0.18** |



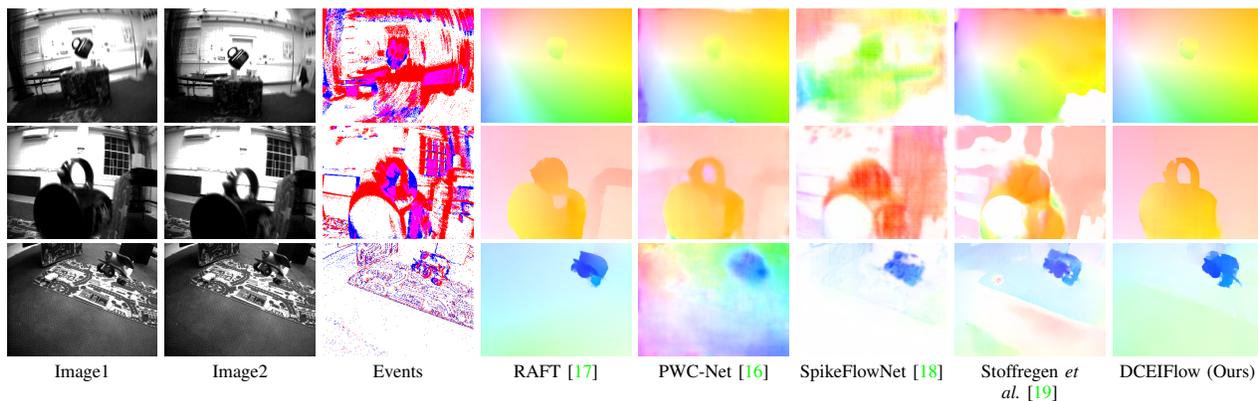| Image1 | Image2 | Events | RAFT [17] | PWC-Net [16] | SpikeFlowNet [18] | Stoffregen *et al.* [19] | DCEIFlow (Ours) |

Fig. 8. **Visual comparisons on the EV-IMO [27] dataset.** The models used for inference are all the pre-trained models that have not been trained on the EV-IMO dataset. Since the ground-truth flow annotations are not publicly available, we use the results of the two-frame methods as a reference. Best viewed on screen.

## C. Results on the EV-IMO dataset

To verify the performance of our model on a more challenging dataset EV-IMO [27] with fast-moving objects, we also run SpikeFlowNet [18] and Stoffregen *et al.* [19] with their open sourced models, and the visual comparisons are shown in Fig. 8. Because there is no ground-truth optical flow, we take the results of the two-frame methods (PWC-Net [16], and RAFT [17]) as a reference for comparative analysis. We found that RAFT achieved perfect flow visualization using two frames, but we also found its shortcomings in some detailed areas compared with ours, especially the edge or hole position. We believe this is the advantage of introducing the events which contain detailed motion. For event-based methods, it is obvious that their predictions are not dense and accurate enough, especially SpikeFlowNet makes some unique estimates for foreground objects. Our results are denser than them, and most consistent with RAFT. This further verifies the superior generalization performance of our model compared with the existing event-based methods.

## D. Model analysis

*1) Baselines for comparison:* To fully illustrate the superiority of our model, we experiment with the same training setting for both baseline models, i.e., 100 epochs on FlyingChairs2 or 300 epochs on the MVSEC dataset using the same hyper-parameters. See Sec. IV-A3 for more training

details. Table IV show the results of two baselines on the MVSEC dataset. In addition, we also evaluate their dense flow prediction ability in Table V. By comparing the two baselines, our proposed DCEIFlow model improves significantly when the model size is less than theirs. This shows that our model is more suitable and powerful than the baselines for event-based dense flow estimation. When trained with only one frame interval setting, our model also achieves better performance than the compared baselines and the two event-based methods SpikeFlownet [18] and Stoffregen *et al.* [19]. In addition, we found that the model trained with two interval settings can achieve better performance than with only one. We think this is because larger data sizes help supervised learning to achieve a better model.

Furthermore, Table IV shows that our proposed DCEIFlow model outperforms the competing event-based baseline models even when directly trained on the MVSEC dataset without pre-training. Compared with the existing event-based unsupervised methods in Table I, our model achieves comparable, if not better, performance. These results show that the methods with unsupervised training usually have better generalization ability. The methods with supervised training may not achieve superior results when trained on limited data. Thus our model achieves a significant performance improvement when pre-trained with a large-scale dataset FlyingChairs2 (*i.e.*, FC2).

*2) Supervised and unsupervised training:* As shown in Table I, some existing methods use unsupervised loss functions

TABLE IV

**PERFORMANCE COMPARISON BETWEEN DIFFERENT BASELINES AND OUR PROPOSED DCEIFLOW MODEL ON THE MVSEC [26] DATASET. WE USE THE SAME PRE-TRAINED MODEL TO GET THE RESULTS OF EACH METHOD UNDER TWO INPUT INTERVAL SETTINGS ($dt=1$ AND $dt=4$ FRAMES).**

| Input $dt=1$ | Method | Train Mann. | Train D.Type | Train D.Set | Eval. Metric | *indoor_flying1* EPE | %Out | *indoor_flying2* EPE | %Out | *indoor_flying3* EPE | %Out | *outdoor_day1* EPE | %Out | *outdoor_day2* EPE | %Out |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **E** | Baseline-EV | SL | E | C2 | dense | 0.91 | 1.06 | 0.98 | 1.04 | 1.02 | 1.60 | 1.09 | 1.56 | 0.99 | 4.07 |
| | Baseline-EV | SL | E | C2 | sparse | 0.93 | 1.11 | 1.04 | 2.47 | 1.00 | 1.22 | 0.95 | 0.82 | 0.98 | 2.52 |
| | Baseline-EV | SL | E | M | dense | 0.80 | 0.65 | 0.95 | 2.04 | 0.89 | 1.56 | 0.44 | 0.03 | × | × |
| | Baseline-EV | SL | E | M | sparse | 0.89 | 1.17 | 1.14 | 4.31 | 0.97 | 2.50 | 0.47 | 0.10 | × | × |
| **I$_1$+E** | Baseline-EI | SL | I$_1$,E | C2 | dense | 0.78 | 0.49 | 0.81 | 0.58 | 0.80 | 0.25 | 0.92 | 0.75 | 0.84 | 3.33 |
| | Baseline-EI | SL | I$_1$,E | C2 | sparse | 0.82 | 0.60 | 0.89 | 1.37 | 0.82 | 0.40 | 0.80 | 0.54 | 0.83 | 2.69 |
| | Baseline-EI | SL | I$_1$,E | M | dense | 0.75 | 0.54 | 0.80 | 0.77 | 0.80 | 0.93 | 0.36 | 0.00 | × | × |
| | Baseline-EI | SL | I$_1$,E | M | sparse | 0.78 | 0.68 | 0.91 | 1.45 | 0.82 | 1.03 | 0.36 | 0.00 | × | × |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | C2 | dense | **0.56** | **0.28** | **0.64** | **0.16** | **0.57** | **0.12** | 0.91 | 0.71 | **0.79** | 2.59 |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | C2 | sparse | 0.57 | 0.30 | 0.70 | 0.30 | 0.58 | 0.15 | **0.74** | **0.29** | 0.82 | **2.34** |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | M | dense | 0.64 | 0.87 | 0.74 | 1.16 | 0.70 | 1.08 | 0.20 | 0.00 | × | × |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | M | sparse | 0.75 | 1.55 | 0.90 | 2.10 | 0.80 | 1.77 | 0.22 | 0.00 | × | × |
| Input $dt=4$ | Method | Train Mann. | Train D.Type | Train D.Set | Eval. Metric | *indoor_flying1* EPE | %Out | *indoor_flying2* EPE | %Out | *indoor_flying3* EPE | %Out | *outdoor_day1* EPE | %Out | *outdoor_day2* EPE | %Out |
| **E** | Baseline-EV | SL | E | C2 | dense | 1.76 | 13.21 | 2.05 | 19.09 | 1.99 | 18.73 | 2.63 | 30.84 | 2.10 | 20.83 |
| | Baseline-EV | SL | E | C2 | sparse | 1.72 | 12.37 | 2.24 | 22.98 | 1.91 | 16.57 | 1.96 | 17.91 | 2.01 | 18.86 |
| | Baseline-EV | SL | E | M | dense | 2.59 | 29.67 | 3.53 | 39.82 | 2.87 | 31.21 | 1.55 | 12.24 | × | × |
| | Baseline-EV | SL | E | M | sparse | 3.04 | 36.87 | 4.55 | 52.58 | 3.38 | 36.72 | 1.64 | 13.64 | × | × |
| **I$_1$+E** | Baseline-EI | SL | I$_1$,E | C2 | dense | 1.66 | 11.17 | 2.13 | 20.17 | 1.78 | 14.56 | 2.25 | 25.19 | 1.92 | 19.39 |
| | Baseline-EI | SL | I$_1$,E | C2 | sparse | 1.65 | 11.21 | 2.37 | 25.70 | 1.79 | 14.57 | 1.66 | 13.20 | 1.83 | 17.55 |
| | Baseline-EI | SL | I$_1$,E | M | dense | 2.08 | 21.13 | 2.77 | 29.20 | 2.32 | 24.59 | 1.16 | 5.98 | × | × |
| | Baseline-EI | SL | I$_1$,E | M | sparse | 2.21 | 23.61 | 3.37 | 36.99 | 2.47 | 26.50 | 1.11 | 5.01 | × | × |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | C2 | dense | **1.49** | **8.14** | **1.97** | **17.37** | **1.69** | **12.34** | 1.87 | 19.13 | 1.62 | 14.73 |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | C2 | sparse | 1.52 | 8.79 | 2.21 | 22.13 | 1.74 | 13.33 | **1.37** | **8.54** | **1.61** | **14.38** |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | M | dense | 1.90 | 17.43 | 2.97 | 34.38 | 2.32 | 26.07 | 0.87 | 3.12 | × | × |
| | DCEIFlow (Ours) | SL | I$_1$,I$_2$,E | M | sparse | 2.08 | 21.47 | 3.48 | 42.05 | 2.51 | 29.73 | 0.89 | 3.19 | × | × |

for training, while our model is obtained by supervised training using ground-truth supervision. For the four methods with input data is $I_1 + I_2$, we only evaluate their pre-trained model for reference comparisons. Because SMURF and ARFlow have similar structures to RAFT and PWC-Net, respectively, the experiments on the MVSEC dataset (*cf.* Table II) indicate that the unsupervised methods usually have better generalization performance. However, in the dense prediction analysis (*cf.* Table V), the unsupervised methods usually have worse results in the weakly textured regions (i.e., without events) than in the richly textured regions (i.e., with events), while the supervised methods are usually better in the weakly textured regions. This conclusion has also been verified by the results of the event-based unsupervised methods. We conclude that the supervised methods can generally achieve better dense flow estimation, but their generalization ability is worse than unsupervised methods under the same training protocols.

*3) Ablation study:* We conducted ablations to confirm the effectiveness of each module in our framework, including 1) event polarity representation, 2) event-image correlation module, 3) event-image feature fusion, and 4) bidirectional flow training. In addition, we also experimented with the commonly used pyramid structure [16] to verify the effectiveness of our iterative structure. More details of this pyramid structure are described in the supplementary material.

The results are shown in Table III. We use the same training setting to pre-train on the FlyingChairs2 dataset for the nine models composed of the above five parts. In models (a)&(b), separating events by their polarity during represent events into event volume can reduce the information loss caused by positive and negative coexistence, and improve the accuracy. In models (b)&(c), introducing the correlation construction can greatly improve the accuracy of flow prediction, which is also proved in the existing two-frame methods [16]. At the same time, because the input dimension of the decoder is increased, the number of network parameters is also increased. In models (c)&(e) and (d)&(f), our proposed fusion by convolution structure is more suitable compared to simple addition. In models (e)&(g) and (f)&(i), the proposed bidirectional training mechanism can further improve the performance without increasing the number of network parameters, and the results of the iterative structure are better than pyramid structure with smaller parameters. In models (h)&(i), since there is no constraint on the output of the fusion module, the correlation module is challenging to realize the function to search for matching points in the neighborhood. The comparison results also illustrate that the fusion module and the similar loss need to be used together to perform better. In summary, the model with both parts obtains the best results, which demonstrates the effectiveness of each component in our framework. In addition, we also compare the dense flow estimation performance of each model on the *indoor_flying1-3* sequences of the MVSEC

TABLE V
**DENSE FLOW PREDICTION ANALYSIS.** THE RESULTS ARE AVERAGED ON THE MVSEC *indoor_flying1-3* SEQUENCES.

| Input | Method | Claim Setting | Train. set ($dt$) | Eval. $dt$ | Full valid pixels EPE | %Out | Event Masked EPE | %Out | Event Excluded EPE | %Out | Dense Ratio |
|---|---|---|---|---|---|---|---|---|---|---|---|
| $I_1+I_2$ | PWC-Net [16] | dense | C+T | 1 | 1.58 | 3.03 | 1.62 | 3.73 | 1.58 | 2.96 | 1.013 |
| | RAFT [17] | dense | C | 1 | 0.49 | 0.06 | 0.55 | 0.07 | 0.49 | 0.06 | 1.060 |
| | ARFlow [54] | dense | SR+S | 1 | 0.43 | 0.07 | 0.42 | 0.05 | 0.43 | 0.08 | 0.996 |
| | SMURF [55] | dense | C | 1 | 0.46 | 0.19 | 0.44 | 0.08 | 0.46 | 0.20 | 0.981 |
| | PWC-Net [16] | dense | C+T | 4 | 2.05 | 17.47 | 2.11 | 18.99 | 2.03 | 16.83 | 1.019 |
| | RAFT [17] | dense | C | 4 | 1.63 | 10.92 | 1.69 | 11.68 | 1.61 | 10.49 | 1.024 |
| | ARFlow [54] | dense | SR+S | 4 | 1.44 | 7.93 | 1.50 | 9.15 | 1.41 | 7.20 | 1.034 |
| | SMURF [55] | dense | C | 4 | 1.53 | 9.29 | 1.54 | 9.72 | 1.52 | 8.86 | 1.007 |
| E | SpikeFlowNet [18] | sparse | M (1) | 1 | 1.14 | 4.67 | 1.08 | 4.03 | 1.15 | 5.16 | 0.969 |
| | Stoffregen *et al.* [19] | dense | ESIM | 1 | 0.75 | 0.57 | 0.65 | 0.41 | 0.76 | 0.60 | 0.927 |
| | E-RAFT [41] | dense | DSEC | 1 | 0.82 | 1.54 | 0.97 | 2.76 | 0.80 | 1.33 | 1.105 |
| | Baseline-EV | dense | M (1) | 1 | 0.83 | 2.38 | 0.99 | 4.60 | 0.80 | 1.96 | 1.114 |
| | Baseline-EV | dense | M (1&4) | 1 | 0.88 | 1.42 | 1.00 | 2.66 | 0.86 | 1.15 | 1.080 |
| | Baseline-EV | dense | C2 | 1 | 0.97 | 1.24 | 0.99 | 1.60 | 0.97 | 1.18 | 1.010 |
| | SpikeFlowNet [18] | sparse | M (4) | 4 | 3.65 | 45.42 | 3.08 | 33.45 | 3.78 | 49.01 | 0.906 |
| | Stoffregen *et al.* [19] | dense | ESIM | 4 | 3.08 | 35.91 | 2.29 | 21.03 | 3.35 | 40.79 | 0.835 |
| | E-RAFT [41] | dense | DSEC | 4 | 2.19 | 19.55 | 2.46 | 23.15 | 2.05 | 17.97 | 1.098 |
| | Baseline-EV | dense | M (4) | 4 | 3.12 | 35.51 | 3.60 | 39.99 | 2.84 | 33.36 | 1.128 |
| | Baseline-EV | dense | M (1&4) | 4 | 3.00 | 33.57 | 3.66 | 42.06 | 2.66 | 30.01 | 1.177 |
| | Baseline-EV | dense | C2 | 4 | 1.93 | 17.01 | 1.96 | 17.31 | 1.90 | 16.46 | 1.015 |
| $I_1+E$ | Baseline-EI | dense | M (1) | 1 | 0.77 | 1.91 | 0.91 | 3.61 | 0.75 | 1.59 | 1.106 |
| | Baseline-EI | dense | M (1&4) | 1 | 0.78 | 0.75 | 0.84 | 1.05 | 0.77 | 0.68 | 1.040 |
| | Baseline-EI | dense | C2 | 1 | 0.80 | 0.44 | 0.84 | 0.79 | 0.79 | 0.40 | 1.033 |
| | DCEIFlow (Ours) | dense | M (1) | 1 | 0.78 | 1.60 | 0.92 | 2.55 | 0.76 | 1.45 | 1.102 |
| | DCEIFlow (Ours) | dense | M (1&4) | 1 | 0.67 | 0.48 | 0.77 | 0.89 | 0.66 | 0.42 | 1.080 |
| | DCEIFlow (Ours) | dense | C2 | 1 | **0.59** | **0.18** | **0.62** | **0.25** | **0.58** | **0.18** | 1.027 |
| | Baseline-EI | dense | M (4) | 4 | 2.73 | 32.52 | 3.00 | 36.79 | 2.57 | 30.34 | 1.081 |
| | Baseline-EI | dense | M (1&4) | 4 | 2.39 | 24.97 | 2.69 | 29.03 | 2.22 | 23.02 | 1.100 |
| | Baseline-EI | dense | C2 | 4 | 1.85 | 15.30 | 1.94 | 17.16 | 1.80 | 14.18 | 1.036 |
| | DCEIFlow (Ours) | dense | M (4) | 4 | 2.58 | 22.00 | 2.93 | 25.80 | 2.41 | 20.41 | 1.106 |
| | DCEIFlow (Ours) | dense | M (1&4) | 4 | 2.24 | 22.42 | 2.52 | 27.33 | 2.10 | 20.27 | 1.098 |
| | DCEIFlow (Ours) | dense | C2 | 4 | **1.72** | **12.62** | **1.82** | **14.75** | **1.67** | **11.69** | 1.045 |

dataset [26] with $dt = 1$. The comparison shows that models with better results on other datasets are usually superior on MVSEC, where the iteration-based models usually perform fewer outliers (%Out) than the pyramid-based models.

*4) Dense and Continuous:* In Table V, we evaluate the performance of dense flow prediction on MVSEC [26]. Besides the dense and sparse (Event Masked) metrics, we also report the mean accuracy of pixels that do not trigger any events and have flow ground-truth (Event Excluded). The event-only methods have consistently better masked results than the excluded results, even Stoffregen *et al.* [19] and Baseline-EV use dense supervision. On the contrary, the excluded results are better than those that use events with an image. Our two models are much better than the others for each metric and frame interval setting. This shows the superiority of our framework in predicting dense flow using events.

For continuous flow estimation, we not only evaluate the results with $dt$=1 and $dt$=4 of each model on MVSEC in Table I, but also compare the results with the non-integer frame number time window on the EV-IMO dataset in Fig. 1 and Fig. 9. Since there is no corresponding second frame image,

the two-frame-based method cannot deal with non-integer intervals. Our model obviously achieves higher accuracy and denser visualization results than other event-based methods. This illustrates the advantages of our framework in predicting continuous and dense optical flow.

*E. Failure cases and limitation*

From the previous experiments, our predictions are better than the two-frame-based start-of-the-art approach RAFT [17] in some cases, especially in the detailed areas, but sometimes worse. This aroused our interest in further analysis. We found two examples on the EV-IMO [27] dataset, as shown in Fig. 10. For the first example, the tail part lacks texture, and events are rarely triggered in such weakly textured regions. Thus for our setting with a single image and events, our model outputs incorrect predictions in the tail part of the toy airplane. We think this shows that the two-frame-based methods such as RAFT can still match the structural association in these challenging regions. For the second example, a part of the object suddenly reflects light. However, the brightness change
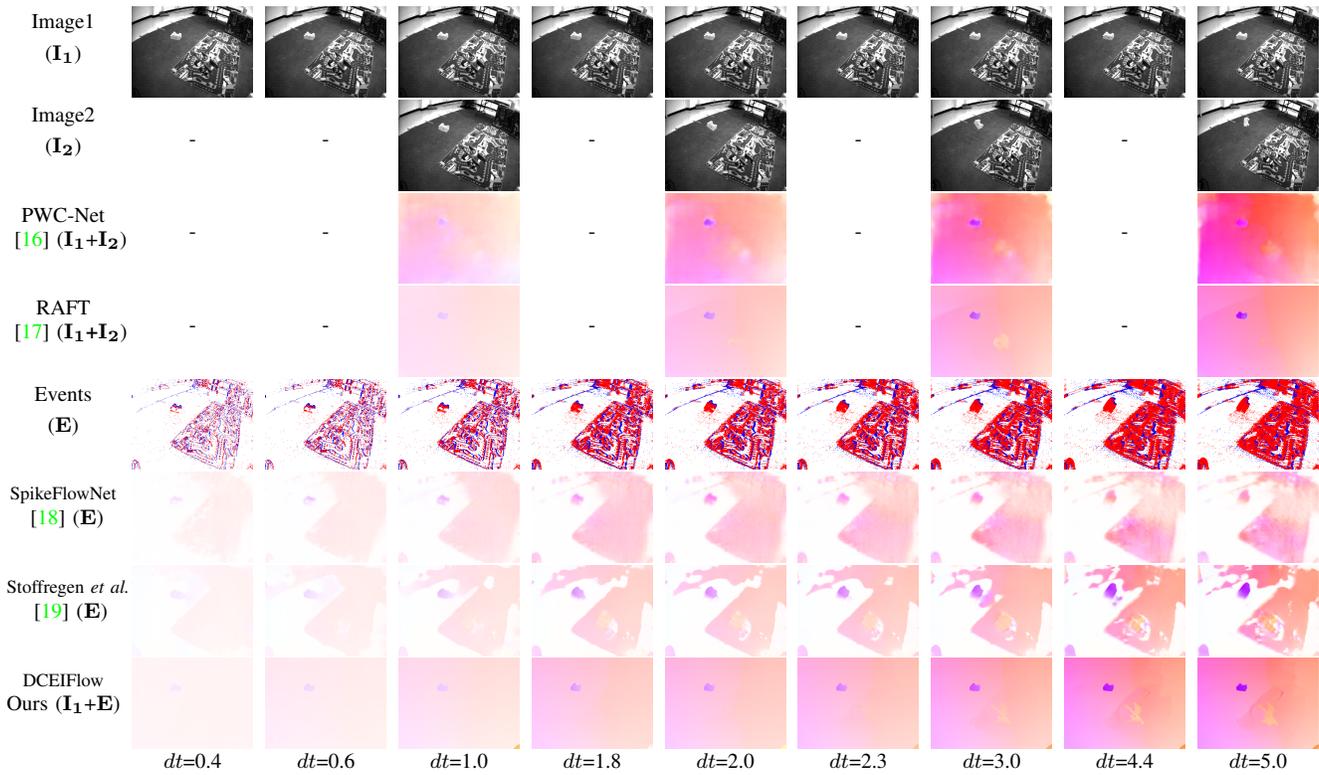
Fig. 9.  **Another visual comparisons of continuous flow prediction with different time intervals on the EV-IMO dataset**. Best viewed on screen.
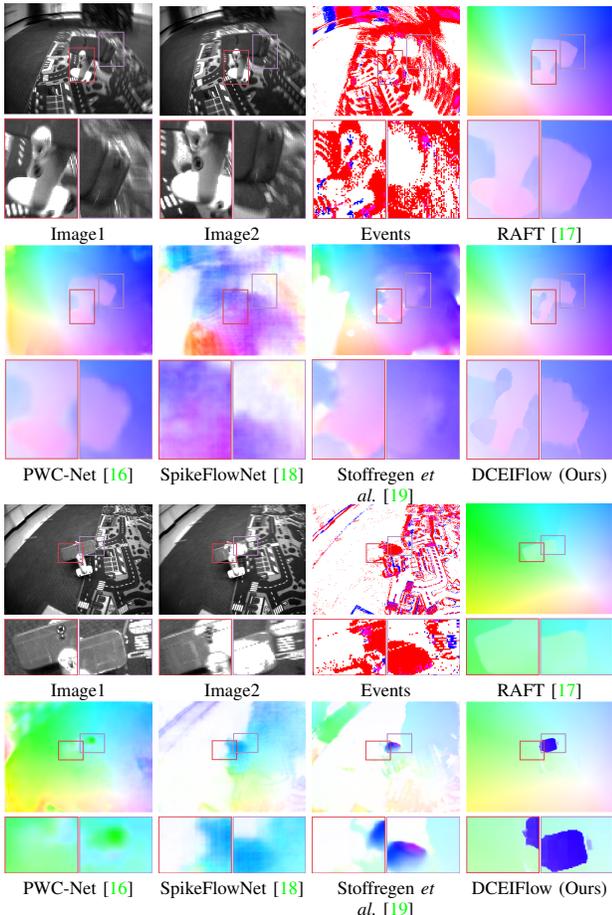


Fig. 10.  **Failure cases on the EV-IMO dataset.** Best viewed on screen.

is not caused by motion, and the object did not move so much. The reflection produces a lot of events, which seriously affects the predictions of event-based methods, including ours.

For event-based applications, the efficiency of the algorithm is worth considering. Using a single 2080ti GPU, our model takes 28ms to process data with MVSEC size. This result meets the *real-time* standard (30 fps) and is much better than the model-based optimization method Pan *et al.* [5], which requires an uncertain running time of more than 1 second. However, compared to EV-FlowNet [20] and Spike-FlowNet [18] that use only events and require only 5ms and 15ms, we have achieved better results with increased computation cost. Higher time consumption will limit the application, which is what we need to improve next.

## V. CONCLUSION

In this paper, we have proposed a deep learning-based dense and continuous optical flow estimation approach from a single image with event streams. Our network can effectively exploit the internal relation of two different modalities of data through an event-image fusion and correlation module, and predict the dense optical flow by the iterative flow update network structure, combined with our bidirectional training strategy. Thus our framework can reliably estimate dense flow as two-frame-based methods, as well as estimate continuous flow as event-based methods. Extensive experimental evaluation on multiple datasets demonstrates the superiority of our proposed framework in estimating dense and continuous optical flow compared with existing state-of-the-art event-only or fused single-image methods.

## REFERENCES

[1] G. Gallego, T. Delbruck, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. Davison, J. Conradt, K. Daniilidis *et al.*, "Event-based vision: A survey," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2020. 1, 3

[2] P. Lichtsteiner, C. Posch, and T. Delbruck, "A 128 × 128 120db 30mw asynchronous vision sensor that responds to relative intensity change," in *IEEE International Solid State Circuits Conference*, 2006, pp. 2060–2069. 1

[3] C. Brandli, R. Berner, M. Yang, S.-C. Liu, and T. Delbruck, "A 240 × 180 130db 3μs latency global shutter spatiotemporal vision sensor," *IEEE Journal of Solid-State Circuits*, vol. 49, no. 10, pp. 2333–2341, 2014. 1

[4] H. Rebecq, R. Ranftl, V. Koltun, and D. Scaramuzza, "High speed and high dynamic range video with an event camera," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, 2019. 1, 4

[5] L. Pan, C. Scheerlinck, X. Yu, R. Hartley, M. Liu, and Y. Dai, "Bringing a blurry frame alive at high frame-rate with an event camera," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 6820–6829. 1, 13

[6] Q. Sabatier, S.-H. Ieng, and R. Benosman, "Asynchronous event-based fourier analysis," *IEEE Trans. on Image Processing (TIP)*, vol. 26, no. 5, pp. 2192–2202, 2017. 1

[7] D. Tedaldi, G. Gallego, E. Mueggler, and D. Scaramuzza, "Feature detection and tracking with the dynamic and active-pixel vision sensor (davis)," in *Second International Conference on Event-based Control, Communication, and Signal Processing (EBCCSP)*, 2016, pp. 1–7. 1

[8] D. Gehrig, H. Rebecq, G. Gallego, and D. Scaramuzza, "Asynchronous, photometric feature tracking using events and frames," in *European Conf. on Computer Vision (ECCV)*, 2018, pp. 750–765. 1, 2

[9] B. Kueng, E. Mueggler, G. Gallego, and D. Scaramuzza, "Low-latency visual odometry using event-based feature tracks," in *IEEE Int. Conf. on Intelligent Robots and Systems (IROS)*, 2016, pp. 16–23. 1

[10] S. Tulyakov, F. Fleuret, M. Kiefel, P. Gehler, and M. Hirsch, "Learning an event sequence embedding for dense event-based deep stereo," in *IEEE Int. Conf. on Computer Vision (ICCV)*, 2019, pp. 1527–1537. 1

[11] H. Kim, S. Leutenegger, and A. J. Davison, "Real-time 3d reconstruction and 6-dof tracking with an event camera," in *European Conf. on Computer Vision (ECCV)*, 2016, pp. 349–364. 1

[12] Y. Zhou, G. Gallego, H. Rebecq, L. Kneip, H. Li, and D. Scaramuzza, "Semi-dense 3d reconstruction with a stereo event camera," in *European Conf. on Computer Vision (ECCV)*, 2018, pp. 235–251. 1

[13] G. Gallego, J. E. Lund, E. Mueggler, H. Rebecq, T. Delbruck, and D. Scaramuzza, "Event-based, 6-dof camera tracking from photometric depth maps," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 40, no. 10, pp. 2402–2412, 2017. 1

[14] Z. Jiang, Y. Zhang, D. Zou, J. Ren, J. Lv, and Y. Liu, "Learning event-based motion deblurring," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 3320–3329. 1

[15] P. R. G. Cadena, Y. Qian, C. Wang, and M. Yang, "Spade-e2vid: Spatially-adaptive denormalization for event-based video reconstruction," *IEEE Trans. on Image Processing (TIP)*, vol. 30, pp. 2488–2500, 2021. 1, 4

[16] D. Sun, X. Yang, M. Liu, and J. Kautz, "Models matter, so does training: An empirical study of cnns for optical flow estimation," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 42, no. 6, pp. 1408–1423, 2019. 1, 2, 3, 8, 9, 10, 11, 12, 13

[17] Z. Teed and J. Deng, "Raft: Recurrent all-pairs field transforms for optical flow," in *European Conf. on Computer Vision (ECCV)*, 2020, pp. 402–419. 1, 2, 8, 9, 10, 12, 13

[18] C. Lee, A. K. Kosta, A. Z. Zhu, K. Chaney, K. Daniilidis, and K. Roy, "Spike-flownet: event-based optical flow estimation with energy-efficient hybrid neural networks," in *European Conf. on Computer Vision (ECCV)*, 2020, pp. 366–382. 1, 3, 7, 8, 9, 10, 12, 13

[19] T. Stoffregen, C. Scheerlinck, D. Scaramuzza, T. Drummond, N. Barnes, L. Kleeman, and R. Mahony, "Reducing the sim-to-real gap for event cameras," in *European Conf. on Computer Vision (ECCV)*, 2020, pp. 534–549. 1, 4, 6, 7, 8, 9, 10, 12, 13

[20] A. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "Ev-flownet: Self-supervised optical flow estimation for event-based cameras," in *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018. 1, 3, 7, 8, 9, 13

[21] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "Unsupervised event-based learning of optical flow, depth, and egomotion," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 989–997. 1, 3, 4, 8

[22] G. Gallego, H. Rebecq, and D. Scaramuzza, "A unifying contrast maximization framework for event cameras, with applications to motion, depth, and optical flow estimation," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 3867–3876. 1, 3

[23] F. Paredes-Vallés, K. Y. Scheper, and G. C. de Croon, "Unsupervised learning of a hierarchical spiking neural network for optical flow estimation: From events to global motion perception," *IEEE Trans. on Pattern Analysis and Machine Intelligence (TPAMI)*, vol. 42, no. 8, pp. 2051–2064, 2019. 1, 3

[24] S. Tulyakov, D. Gehrig, S. Georgoulis, J. Erbach, M. Gehrig, Y. Li, and D. Scaramuzza, "Time lens: Event-based video frame interpolation," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 16 155–16 164. 2

[25] E. Ilg, T. Saikia, M. Keuper, and T. Brox, "Occlusions, motion and depth boundaries with a generic network for disparity, optical flow or scene flow estimation," in *European Conf. on Computer Vision (ECCV)*, 2018, pp. 614–630. 2, 6, 8, 9, 10

[26] A. Z. Zhu, D. Thakur, T. Özaslan, B. Pfrommer, V. Kumar, and K. Daniilidis, "The multivehicle stereo event camera dataset: An event camera dataset for 3d perception," *IEEE Robotics and Automation Letters*, vol. 3, no. 3, pp. 2032–2039, 2018. 2, 6, 7, 8, 9, 10, 11, 12

[27] A. Mitrokhin, C. Ye, C. Fermüller, Y. Aloimonos, and T. Delbruck, "Ev-imo: Motion segmentation dataset and learning pipeline for event cameras," in *IEEE Int. Conf. on Intelligent Robots and Systems (IROS)*, 2019, pp. 6105–6112. 2, 6, 10, 12

[28] J. Hur and S. Roth, *Optical Flow Estimation in the Deep Learning Age*. Cham: Springer International Publishing, 2020, pp. 119–140. 2

[29] A. Dosovitskiy, P. Fischer, E. Ilg, P. Hausser, C. Hazirbas, V. Golkov, P. Van Der Smagt, D. Cremers, and T. Brox, "Flownet: Learning optical flow with convolutional networks," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2015, pp. 2758–2766. 2

[30] J. Hur and S. Roth, "Iterative residual refinement for joint optical flow and occlusion estimation," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 5754–5763. 2, 5

[31] P. Liu, M. Lyu, I. King, and J. Xu, "Selflow: Self-supervised learning of optical flow," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4571–4580. 2, 5

[32] G. Yang and D. Ramanan, "Volumetric correspondence networks for optical flow," *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 32, pp. 794–805, 2019. 2

[33] R. Benosman, S.-H. Ieng, C. Clercq, C. Bartolozzi, and M. Srinivasan, "Asynchronous frameless event-based optical flow," *Neural Networks*, vol. 27, pp. 32–37, 2012. 3

[34] B. D. Lucas, T. Kanade *et al.*, "An iterative image registration technique with an application to stereo vision," in *Int. Joint Conf.s on Artificial Intelligence (IJCAI)*, vol. 81, 1981, p. 674–679. 3

[35] M. Liu and T. Delbrück, "Adaptive time-slice block-matching optical flow algorithm for dynamic vision sensors," in *British Machine Vision Conf. (BMVC)*, 2018, p. 280. 3

[36] W. Maass, "Networks of spiking neurons: the third generation of neural network models," *Neural networks*, vol. 10, no. 9, pp. 1659–1671, 1997. 3

[37] D. Gehrig, A. Loquercio, K. G. Derpanis, and D. Scaramuzza, "End-to-end learning of representations for asynchronous event-based data," in *IEEE Int. Conf. on Computer Vision (ICCV)*, 2019, pp. 5633–5643. 3, 8

[38] M. Cannici, M. Ciccone, A. Romanoni, and M. Matteucci, "A differentiable recurrent surface for asynchronous event-based data," in *European Conf. on Computer Vision (ECCV)*, 2020, pp. 136–152. 3, 8

[39] J. Hagenaars, F. Paredes-Vallés, and G. De Croon, "Self-supervised learning of event-based optical flow with spiking neural networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021, pp. 7167–7179. 3, 8

[40] Z. Ding, R. Zhao, J. Zhang, T. Gao, R. Xiong, Z. Yu, and T. Huang, "Spatio-temporal recurrent networks for event-based optical flow estimation," in *AAAI Conf. on Artificial Intelligence (AAAI)*, vol. 36, no. 1, 2022, pp. 525–533. 3, 8

[41] M. Gehrig, M. Millhäusler, D. Gehrig, and D. Scaramuzza, "E-raft: Dense optical flow from event cameras," in *Int. Conf. on 3D Vision (3DV)*, 2021, pp. 197–206. 3, 9, 12

[42] Z. Li, J. Shen, and R. Liu, "A lightweight network to learn optical flow from event data," in *Int. Conf. on Pattern Recognition (ICPR)*, 2021, pp. 1–7. 3, 8

[43] P. Bardow, A. J. Davison, and S. Leutenegger, "Simultaneous optical flow and intensity estimation from an event camera," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 884–892. 3

[44] L. Pan, M. Liu, and R. Hartley, "Single image optical flow estimation with an event camera," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 1669–1678. 3, 8, 9

[45] C. Lee, A. K. Kosta, and K. Roy, "Fusion-flownet: Energy-efficient optical flow estimation using sensor fusion and deep fused spiking-analog network architectures," in *Int. Conf. on Robotics and Automation (ICRA)*, 2022, pp. 6504–6510. 3, 4, 8

[46] C. Lei and Y.-H. Yang, "Optical flow estimation on coarse-to-fine region-trees using discrete optimization," in *IEEE Int. Conf. on Computer Vision (ICCV)*, 2009, pp. 1562–1569. 3

[47] F. Güney and A. Geiger, "Deep discrete flow," in *Asian Conf. on Computer Vision (ACCV)*, 2016, pp. 207–224. 3

[48] M. Hofinger, S. R. Bulò, L. Porzi, A. Knapitsch, T. Pock, and P. Kontschieder, "Improving optical flow on a pyramid level," in *European Conf. on Computer Vision (ECCV)*, 2020, pp. 770–786. 3

[49] E. Ilg, N. Mayer, T. Saikia, M. Keuper, A. Dosovitskiy, and T. Brox, "Flownet 2.0: Evolution of optical flow estimation with deep networks," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017, pp. 2462–2470. 3

[50] H. Rebecq, R. Ranftl, V. Koltun, and D. Scaramuzza, "Events-to-video: Bringing modern computer vision to event cameras," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 3857–3866. 4

[51] F. Paredes-Vallés and G. C. de Croon, "Back to event basics: Self-supervised learning of image reconstruction for event cameras via photometric constancy," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3446–3455. 4, 8

[52] K. Cho, B. van Merrienboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2014, pp. 1724–1734. 5

[53] Y. Deng, H. Chen, H. Chen, and Y. Li, "Learning from images: A distillation learning framework for event cameras," *IEEE Trans. on Image Processing (TIP)*, vol. 30, pp. 4919–4931, 2021. 8

[54] L. Liu, J. Zhang, R. He, Y. Liu, Y. Wang, Y. Tai, D. Luo, C. Wang, J. Li, and F. Huang, "Learning by analogy: Reliable supervision from transformations for unsupervised optical flow estimation," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 6489–6498. 8, 9, 12

[55] A. Stone, D. Maurer, A. Ayvaci, A. Angelova, and R. Jonschkowski, "Smurf: Self-teaching multi-frame unsupervised raft with full-image warping," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021, pp. 3887–3896. 8, 9, 12

[56] H. Rebecq, D. Gehrig, and D. Scaramuzza, "Esim: an open event camera simulator," in *Conference on Robot Learning*, 2018, pp. 969–982. 6, 7

[57] D. J. Butler, J. Wulff, G. B. Stanley, and M. J. Black, "A naturalistic open source movie for optical flow evaluation," in *European Conf. on Computer Vision (ECCV)*, 2012, pp. 611–625. 6, 10

[58] D. Gehrig, M. Gehrig, J. Hidalgo-Carrió, and D. Scaramuzza, "Video to events: Recycling video datasets for event cameras," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 3586–3595. 7

[59] H. Jiang, D. Sun, V. Jampani, M.-H. Yang, E. Learned-Miller, and J. Kautz, "Super slomo: High quality estimation of multiple intermediate frames for video interpolation," in *IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 9000–9008. 7

[60] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 8024–8035. 7

[61] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *Int. Conf. on Learning Representations (ICLR)*, 2015. 7

[62] L. N. Smith and N. Topin, "Super-convergence: Very fast training of neural networks using large learning rates," in *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, vol. 11006, 2019, p. 1100612. 7

[63] M. Menze, C. Heipke, and A. Geiger, "Joint 3d estimation of vehicles and scene flow." *ISPRS Annals of Photogrammetry, Remote Sensing*, vol. 2, 2015. 7

[64] P. Shedligeri and K. Mitra, "High frame rate optical flow estimation from event sensors via intensity estimation," *Computer Vision and Image Understanding (CVIU)*, vol. 208, p. 103208, 2021. 8

[65] M. Mostafavi, L. Wang, and K.-J. Yoon, "Learning to reconstruct hdr images from events, with applications to depth and flow prediction," *Int.*

*Journal of Computer Vision (IJCV)*, vol. 129, no. 4, pp. 900–920, 2021. 8

[66] M. Gehrig, W. Aarents, D. Gehrig, and D. Scaramuzza, "Dsec: A stereo event camera dataset for driving scenarios," *IEEE Robotics and Automation Letters (RA-L)*, 2021. 9

**Zhexiong Wan** is currently a PhD student with School of Electronics and Information, Northwestern Polytechnic University, Xi'an, China. He received his Bachelor of Engineering degree from Northwestern Polytechnical University in 2019.

**Yuchao Dai** is currently a Professor with School of Electronics and Information at the Northwestern Polytechnical University. He received the B.E. degree, M.E degree and Ph.D. degree all in signal and information processing from Northwestern Polytechnical University, Xi'an, China, in 2005, 2008 and 2012, respectively. He was an ARC DECRA Fellow with the Research School of Engineering at the Australian National University, Canberra, Australia from 2014 to 2017 and a Research Fellow with the Research School of Computer Science at the Australian National University, Canberra, Australia from 2012 to 2014. His research interests include 3D vision, multi-view geometry, low-level computer vision, deep learning, and optimization. He won the Best Paper Award in IEEE CVPR 2012, Best Paper Nominee in IEEE CVPR 2020, the DSTO Best Fundamental Contribution to Image Processing Paper Prize at DICTA 2014, the Best Algorithm Prize in NRSFM Challenge at CVPR 2017, the Best Student Paper Prize at DICTA 2017 and the Best Deep/Machine Learning Paper Prize at APSIPA ASC 2017. He served/serves as Area Chair at CVPR, ICCV, ACM MM, ACCV, etc. He serves as the Publicity Chair at ACCV 2022 and the Distinguished Lecturer of APSIPA.

**Yuxin Mao** is currently a PhD student with School of Electronics and Information, Northwestern Polytechnical University, Xi'an, China. He received his Bachelor of Engineering degree from Southwest Jiaotong University in 2020. He won the best Paper Award Nominee at ICIUS 2019.