

# Fusion-FlowNet: Energy-Efficient Optical Flow Estimation using Sensor Fusion and Deep Fused Spiking-Analog Network Architectures

Chankyu Lee, Adarsh Kumar Kosta and Kaushik Roy  
 Purdue University, West Lafayette, IN 47907, USA  
 {lee2216, akosta, kaushik}@purdue.edu

## Abstract

*Standard frame-based cameras that sample light intensity frames are heavily impacted by motion blur for high-speed motion and fail to perceive scene accurately when the dynamic range is high. Event-based cameras, on the other hand, overcome these limitations by asynchronously detecting the variation in individual pixel intensities. However, event cameras only provide information about pixels in motion, leading to sparse data. Hence, estimating the overall dense behavior of pixels is difficult. To address such issues associated with the sensors, we present Fusion-FlowNet, a sensor fusion framework for energy-efficient optical flow estimation using both frame- and event-based sensors, leveraging their complementary characteristics. Our proposed network architecture is also a fusion of Spiking Neural Networks (SNNs) and Analog Neural Networks (ANNs) where each network is designed to simultaneously process asynchronous event streams and regular frame-based images, respectively. Our network is end-to-end trained using unsupervised learning to avoid expensive video annotations. The method generalizes well across distinct environments (rapid motion and challenging lighting conditions) and demonstrates state-of-the-art optical flow prediction on the Multi-Vehicle Stereo Event Camera (MVSEC) dataset. Furthermore, our network offers substantial savings in terms of the number of network parameters and computational energy cost.*

## 1. Introduction

Optical flow estimation is a fundamental computer vision problem, allowing us to visualize the motion field in scenes. It involves estimating the spatio-temporal motion patterns of pixels and forms the groundwork for more complex tasks such as motion segmentation [22] and action recognition [29]. Over the past years, the optical flow estimation has been largely dominated by conventional computer vision algorithms such as differential [18], phase cor-

relation and block-based methods [2]. Recently, deep Analog Neural Networks (ANNs<sup>1</sup>) based approaches for optical flow estimation have gained immense popularity [7, 24]. In general, these methods rely on the standard frame-based cameras as input sensors that capture pixel intensities over the entire frame at a regular sampling rate. However, the frame-based images suffer from a variety of issues such as motion blur and temporal aliasing when capturing high speed motion due to fixed low temporal resolution. They are also unable to perceive information accurately in high dynamic range scenes due to uneven exposure [8].

Event-based cameras, such as Dynamic Vision Sensors (DVS) [16], address these problems by asynchronously sampling intensity changes on each pixel element, generating a stream of asynchronous events. This grants promising advantages, namely high temporal resolution (10 $\mu$ s vs 3ms), high dynamic range (140dB vs 60dB) and low power consumption (10mW vs 3W) compared to standard frame-based cameras [8]. Note, event cameras only capture the varying components of visual signals, generating sparse event streams. Hence, the output prediction becomes limited only at pixels-points where events exist, adding difficulty towards encoding the scene context.

As is evident from the above discussions, none of the above sensors by themselves is able to effectively capture all relevant information of a scene. The limited applicability of each individual camera gives rise to the need for an optimal sensor-fusion technique, enabling the sensors to complement the limitations of each other. Such a technique would provide a practical solution towards accurately estimating dense pixel-wise motion in challenging scenarios such as rapid motion and high dynamic range environments.

Conventional computer vision and ANN-based methods are incompatible at handling the discrete and asynchronous event streams from event-based camera in their native form. This is due to the fact that these methods are generally de-

<sup>1</sup>We refer to the conventional deep learning networks as ANNs owing to their analog nature of inputs and computations. This nomenclature helps to distinguish them from Spiking Neural Networks (SNNs) which perform event-based computations.

signed for frame-based images, assuming regular frame rate and brightness consistency over the entire frames. In this regard, Spiking Neural Networks (SNNs), inspired from biological neuronal mechanisms, show a great promise for directly handling event-camera outputs. Moreover, SNNs perform event-based operations by carrying out the computations only at the arrival of the input events, exploiting the inherent sparsity of spatio-temporal event streams and thus enabling energy-efficient computations on specialized neuromorphic hardware such as Loihi from Intel Labs [6] and TrueNorth from IBM [21].

In this work, we propose a method for combining the advantages of regular frame-based images and a stream of asynchronous events. For this purpose, we present Fusion-FlowNet, a deep fused spiking-analog architecture for estimating optical flow that uses sensors of different modalities (standard frame-based images and asynchronous event streams). Our main contributions are as follows:

- We propose Fusion-FlowNet architecture composed of a fusion of SNNs and ANNs for simultaneously processing event streams and frame-based images, leveraging their complementary sensing capabilities.
- We present a Signed Integrate-and-Fire (SIF) neuron model for SNNs which can generate spike outputs with polarity (either positive or negative). In addition, we show that the SIF model coupled with a surrogate gradient method enables end-to-end learning in SNNs.
- We show that Fusion-FlowNet outperforms the corresponding previous works in terms of optical flow estimation on the Multi-Vehicle Stereo Event Camera (MVSEC) dataset. Furthermore, we analyze that Fusion-FlowNet provides substantial savings in terms of network parameters and computational energy cost.

## 2. Related Works

Over the past few years, there have been major advancements towards optical flow estimation using event-cameras. Conventional computer vision algorithms have been adapted to encompass the asynchronous event stream from these sensors in [1, 3, 9]. In ANN-based approaches, the event streams are essentially accumulated for fixed time intervals to generate synchronous frames. In EV-FlowNet [32], the recent event counts as well as pixel-wise last timestamp information are encoded in a frame-based representation. However, this approach heavily suffers during rapid motion and in scenarios with dense localized events, resulting in loss of rich spatio-temporal information. Researchers in [34] proposed a 3D input representation of events interpolated in a 3D volume with time dimension comprising the input channels to retain the temporal fidelity. Nevertheless, this approach struggled to es-

timate the dense predictions in image regions with fewer events.

In general, SNNs provide advantages towards directly handling the asynchronous events and exploiting the inherent temporal information. Recently, authors in Spike-FlowNet [15] aimed to overcome a noticeable drawback of SNNs – namely the “spike vanishing” phenomenon where the number of spikes drastically reduce in the deeper layers, hindering learning. They proposed to effectively integrate SNNs and ANNs into a single network with the SNN layers enabling efficient event stream handling and the ANN layers addressing the spike vanishing problem. Note, since they used only the event streams as input, the predictions were limited to only the pixel locations containing non-zero number of events. Hence, estimating dense motion behavior was greatly limited.

In contrast, researchers in [23] presented a two-step approach to estimate optical flow by jointly using a set of events and a single frame-based image. They employed an optimization-based method to restore a sharp intensity image from the inputs, followed by ANN-based flow estimation methods [17, 28] on the restored frame-based image to generate an optical flow prediction. Contrary to [23], we explore an end-to-end learning approach that can directly process event streams and frame-based images for predicting final outputs while skipping the image restoration step. Moreover, our proposed method utilizes all available event streams as well as frame-based images within a time window, enabling accurate optical flow estimations over longer time windows.

## 3. Method

### 3.1. Sensors and Input representation

#### 3.1.1 Frame-based Camera

Frame-based cameras have been widely popular for computer vision applications. They provide dense and highly accurate pixel intensity information as frames over regular time intervals. In general, the frame intensity information is pivotal in various computer vision applications that require the high degree of accuracy such as face and object recognition [20]. Optical flow estimation using ANNs requires consecutive frame-based images to pass through separate input channels to the network. In our work, this input representation is utilized for the ANN part of Fusion-FlowNet (Sec. 3.3).

#### 3.1.2 Event-based Camera

Event-based cameras are novel vision sensors, emulating the functionality of biological retina cells [19]. Event cameras transmit a stream of asynchronous events as the outcome of tracking intensity changes ( $I$ ) at each pixel ele-

ment, thereby capturing the relative motion of objects in the scene. Whenever the logarithmic intensity change at a pixel element surpasses a specified threshold ( $\theta$ ), a discrete event is asynchronously generated as follows:

$$\|\log(I_{t+1}) - \log(I_t)\| \geq \theta \quad (1)$$

Event cameras provide the data in Address Event Representation (AER) format which incorporates a tuple  $\{x, y, t, p\}$ , comprising the pixel address ( $x$  and  $y$  locations), timestamp ( $t$ ) and polarity of the intensity change ( $p$ ). Here, each ON/OFF polarity corresponds to the increase or decrease in intensity of the pixel, respectively.

Event cameras may not be generally suited for vision applications which need precise intensity information. However, their high temporal resolution and high dynamic range in addition to having low power consumption, make them ideal for usage on resource constrained platforms operating in challenging environments. Optical flow estimation is one such task which heavily suffers in such environments when realized using standard frame-based cameras and can greatly benefit with the usage of event-cameras.

In our work, the raw event stream is transformed into two groups (former and latter) of discretized event frames and are passed as inputs to the SNN part of Fusion-FlowNet (Sec. 3.3). The input to the SNN encoder-branch consists of a sequence of event frames with four channels, each from the ON/OFF polarity of event frames from the former and the latter groups as illustrated in Fig. 1. This representation preserves the spatio-temporal information in the event stream while displaying superior algorithmic performance and high energy-efficiency.

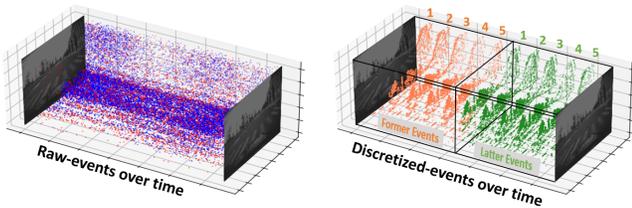


Figure 1. (left) Asynchronous raw event stream between two consecutive frame-based images. (right) Discretized event-frames between two consecutive frame-based images to shape the former and latter groups of events.

### 3.1.3 Sensor-fusion

Interestingly, numerous available sensors, including the Dynamic and Active Vision Sensor (DAVIS) [4], are capable of simultaneously generating the asynchronous events as well as synchronous grayscale frames, simplifying the hardware costs of sensor-fusion. In addition, since there is a single camera coordinate system for both data modalities,

the requirements for any expensive transformation and synchronization between multiple coordinate systems are eliminated. For this purpose, we employ the DAVIS sensor in this work.

In our work, the frame-based images serve two objectives. First, they are provided as network inputs and allow for dense optical flow predictions. Second, they are used for constructing the unsupervised loss required for training. On the other hand, the event streams are only provided as network inputs and enable accurate optical flow prediction in challenging environments as discussed previously. The proposed sensor fusion framework would thus allow to accurately estimate dense optical flow.

## 3.2. Neuron Models

The primary difference between ANN and SNN operations is the notion of time. While ANNs feed-forward the dense analog-valued inputs at once, SNNs process the sparse binary inputs (spikes) as a function of time. Accordingly, different neuron models are employed in ANNs and SNNs.

### 3.2.1 LeakyReLU Model

In ANNs, LeakyReLU neuron [31] replaces the negative part of the popular ReLU model by a linear function with a relatively small slope as below:

$$y = \begin{cases} x, & \text{if } x > 0 \\ \alpha x, & \text{otherwise} \end{cases} \quad (2)$$

where  $\alpha$  is typically set to 0.01-0.1. Note, ReLU has a “dead neuron” problem that some neurons could get stuck in the negative side and play no role in discriminating between inputs. LeakyReLU addresses this problem by having a non-zero slope in the negative direction. This makes it useful especially for hard regression tasks such as motion estimation and predicting pixel-wise and high resolution outputs. In our work, LeakyReLU is employed for the ANN part of Fusion-FlowNet.

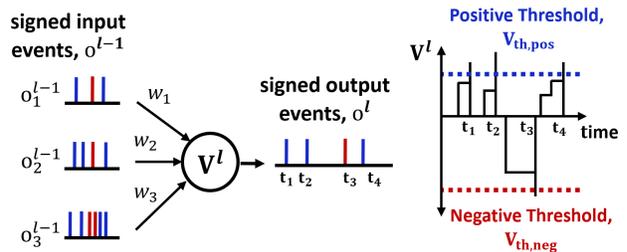


Figure 2. Dynamics of Signed Integrate-and-Fire (SIF) neuron model. Whenever the membrane potential crosses either positive- or negative-threshold, the neuron fires a signed spike output and resets its membrane potential.

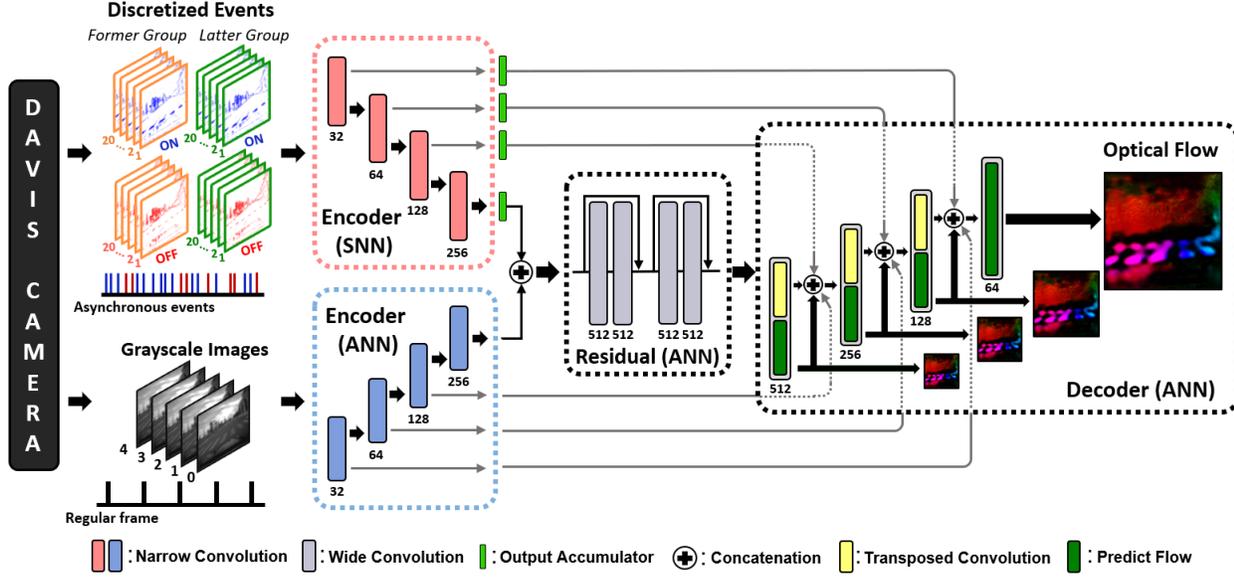


Figure 3. Detailed illustration of Fusion-FlowNet<sub>Early</sub>. The network contains the SNN- and ANN-based encoder-branches to extract features from asynchronous event streams and synchronized grayscale images, respectively. The rest of networks, involving residual and decoder blocks, are composed of ANN layers. The colors represent the types of layers. Best viewed in color.

### 3.2.2 Signed Integrate-and-Fire (SIF) Model

Spiking neurons are inspired by biological models for emulating the efficient event-based operations in the human brain. In the literature, the Integrate-and-Fire (IF) neuron model [5] is widely used for building SNNs because of its simplicity. In an IF neuron, input spikes are modulated by weight ( $w$ ) and accumulated in an internal state of the neuron, called membrane potential over time. In the discrete time model, whenever the membrane potential ( $v$ ) crosses a firing threshold, the neuron emits a binary output (1 or 0) and resets the membrane potential as follows,

$$v^l[n+1] = v^l[n] + w^l o^{l-1}[n] \quad (3)$$

where  $o^{l-1}[n]$  indicates the spike output from previous layer at time-step  $n$ . However, the IF neuron would also suffer from the “dead neuron” problem. To address this issue, we propose a Signed Integrate-and-Fire (SIF) neuron model that can generate signed spike outputs. The SIF neuron is equipped with positive and negative thresholds that enable the generation of positive- and negative-valued spike outputs, respectively. This operation is illustrated in Fig. 2 and formulated as follows:

$$o^l = \begin{cases} +1, & \text{if } v^l > v_{th,pos} \\ -1, & \text{elif } v^l < v_{th,neg} \\ 0, & \text{otherwise} \end{cases} \quad (4)$$

However, the discontinuous and non-differentiable spike generation function of SIF model poses a critical challenge

for conventional gradient-based learning. To overcome this challenge, we propose a surrogate gradient method for the SIF neuron to enable end-to-end backpropagation (discussed in Sec. 3.5).

### 3.3. Fusion-FlowNet Architecture

The Fusion-FlowNet incorporates a deep fused network architecture that supports an end-to-end learning. It is built upon the U-Net architecture [25] that contains four encoder layers, two residual blocks and four decoder layers. The distinctions in our work involve the addition of dual pathways starting at the encoder, namely the SNN- and ANN-based branches. Each branch is composed of narrow convolution layers (similar to grouped convolutions used in AlexNet [14]) containing half the number of intermediate feature maps, compared to the original wide convolution layers. This is possible because of the usage of different modalities of input data, leading to reduction in network parameters without compromising on qualitative performance.

In the SNN-based encoder-branch, the four-channeled input event frames sequentially pass through the narrow convolution layers consisting of SIF neurons over time while being downsampled at each layer. At every time-step, the weighted spike outputs from each layer are integrated into the corresponding output accumulator. After passing all consecutive event images, the accumulated output is passed on ahead to subsequent layers.

In the ANN-based encoder-branch, the consecutive frame-based images in the time window pass through the

narrow ANN layers in a single time step. Each ANN layer comprises of a convolution, batch-norm [11] and a LeakyReLU layer. Here too, the feature maps are down-sampled at each layer.

After completing the forward propagation in both encoder-branches, the outputs are fused together before passing through the rest of the network. This is achieved by concatenating the intermediate activations from both the SNN and ANN branches at the same spatial locations. The fused activations from the last encoder layer pass through the residual blocks while the fused intermediate encoder outputs serve as input to corresponding layers in the decoder block. The four layers in the decoder block perform upsampling using transposed convolutions as well as produce multi-scale optical flow predictions. The multi-scale flow predictions, the transposed convolution outputs and the corresponding activations from the encoder layers are all concatenated together to construct the input for the next decoder layer. Finally, a full-scale optical flow prediction having the same dimension as the input frames is made at the final decoder layer. Fig. 3 showcases the proposed Fusion-FlowNet architecture, illustrating the discussed operations.

### 3.4. Unsupervised Training Method

Due to the limited availability of event-camera datasets containing ground-truth labels, we adopt an unsupervised approach to train the optical flow estimation network [13]. Fusion-FlowNet is trained using unlabeled sequences, utilizing frame-based images for computing the loss. The overall loss function is composed of two parts:

$$l_{total} = l_{photo} + \lambda l_{smooth} \quad (5)$$

where  $l_{photo}$  and  $l_{smooth}$  represent photometric loss and smoothness loss respectively, and  $\lambda$  denotes the loss weight factor.

#### 3.4.1 Photometric Loss

Photometric loss helps to realize the object motion over time by tracking the pixel intensities between images. It is computed by using the start and end-frame grayscale images ( $I_t(x, y)$ ,  $I_{t+dt}(x, y)$ ) as well as the predicted optical flow. A spatial transformer [12] inversely warps the end-frame image ( $I_{t+dt}(x, y)$ ) using the current estimated optical flow ( $u, v$ ) to obtain an image prediction ( $I_{t+dt}(x + u, y + v)$ ). Then, the photometric loss ( $l_{photo}$ ) aims to minimize the discrepancy between the start-frame image ( $I_t(x, y)$ ) and the image prediction ( $I_{t+dt}(x + u, y + v)$ ). The computation is as follows:

$$l_{photo} = \sum_{x,y} \rho(I_t(x, y) - I_{t+dt}(x + u, y + v)) \quad (6)$$

where  $I_t$  ( $I_{t+dt}$ ) indicates the pixel intensity of the first (last) frame-based image,  $u, v$  are the flow estimates in the  $x, y$  directions,  $\rho$  is the robust Charbonnier loss  $\rho(x) = (x^2 + \eta^2)^r$  used for outlier rejection [27]. We set  $r = 0.45$  and  $\eta = 1e^{-3}$  as they show optimal results in prior works [32, 15].

#### 3.4.2 Smoothness Loss

Smoothness loss ( $l_{smooth}$ ) is applied to reduce the optical flow deviations between neighboring pixels by adding a regularizing effect on the predicted flow. It is computed as follows:

$$l_{smooth} = \sum_j \sum_i (\|u_{i,j} - u_{i+1,j}\| + \|u_{i,j} - u_{i,j+1}\| + \|v_{i,j} - v_{i+1,j}\| + \|v_{i,j} - v_{i,j+1}\|) \quad (7)$$

where  $u_{i,j}$  and  $v_{i,j}$  are the flow estimates at pixel location  $(i, j)$  in the  $x$  and  $y$  directions, respectively.

### 3.5. Backpropagation in Fusion-FlowNet

After forward propagation, the final loss ( $l_{total}$ ) is evaluated and used to perform the backward propagation of the gradients. In ANN layers, the LeakyReLU is a differentiable activation that can be represented by the linear functions where the slope differs in positive and negative parts of input as shown in *left* of Fig. 4. The derivative of LeakyReLU activation ( $\frac{\partial f(x)}{\partial x}$ ) is unity when input is positive,  $\alpha$  when input is negative, and zero otherwise. Hence, standard backpropagation can calculate the gradient of the loss function with respect to each weight using chain rule. The parameter updates for the  $l^{th}$  ANN layer are described as follows:

$$\Delta w_{ANN}^l = \frac{\partial loss}{\partial f(x^l)} \frac{\partial f(x^l)}{\partial o^l} \frac{\partial o^l}{\partial w^l} \quad (8)$$

By contrast, the spike generation mechanism of SIF neuron results in a hard threshold function, making it discontinuous and non-differentiable. Hence, standard backpropagation cannot be directly applied to SNNs in its native form as illustrated in *right* of Fig. 4. To overcome this impediment, we present a surrogate gradient method for approximately estimating the spike generation function of SIF neuron. The surrogate gradient of SIF model is herein computed as follows:

$$\frac{\partial o[n]}{\partial v[n]} = \begin{cases} \frac{1}{V_{th,pos}}, & \text{if } v^l > v_{th,pos} \\ \frac{1}{V_{th,neg}}, & \text{if } v^l < v_{th,neg} \\ 0, & \text{otherwise.} \end{cases} \quad (9)$$

where each threshold ( $V_{th,pos}, V_{th,neg}$ ) accounts for the change in the signed spike outputs with respect to the inputs.

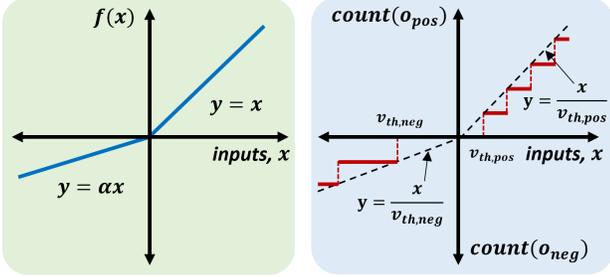


Figure 4. Illustration of activation functions (*left*) LeakyReLU neuron (*right*) Signed Integrate-and-Fire (SIF) neuron.

During the backward pass, the errors ( $\frac{\partial l_{total}}{\partial o^l}$ ) are back-propagated through the SNN layers using the surrogate gradient method and BackPropagation Through Time (BPTT) [30]. In BPTT, the network is unrolled for all time-steps and the weight update is assessed as the sum of gradients over each time-step. The parameter updates of the  $l^{th}$  SNN layer are described as follows:

$$\Delta w_{SNN}^l = \sum_n \frac{\partial \text{loss}}{\partial o^l[n]} \frac{\partial o^l[n]}{\partial v^l[n]} \frac{\partial v^l[n]}{\partial w^l} \quad (10)$$

## 4. Experiments

### 4.1. Dataset and Training Details

We validate Fusion-FlowNet on the MVSEC dataset [33] which contains events as well as grayscale frame sequences recorded using the DAVIS346 camera [4] in multiple indoor and outdoor environments. We use three indoor\_flying sequences and two outdoor\_driving sequences. The indoor\_flying sequences were collected using a drone flying in a closed room containing a variety of objects and are used mainly for evaluation. The outdoor\_day sequences were recorded from a car driving on public roads. We employ the outdoor\_day2 sequence for training and outdoor\_day1 sequence for evaluation. The training and evaluation are performed for two different time-window lengths (i.e,  $dt=1$  and  $dt=4$ ). Every consecutive pair of grayscale images encapsulate an event volume where  $dt=1$  corresponds to constructing inputs using one such event volume while  $dt=4$  corresponds to using four such volumes.

The event streams and frame-based images from left-camera are used for training. They are pre-processed by randomly cropping to  $256 \times 256$  size and flipping horizontally and vertically (with 0.5 probability). The learning rate is scaled by 0.7 every 5 epochs until 20 epoch, and every 10 epochs thereafter. The number of event frames in each group are set to 5 for the  $dt = 1$  case and 20 for the  $dt = 4$  case. In ANN layers, LeakyReLU model is employed with an  $\alpha$  of 0.1. In SNN layers, the positive and negative thresholds of the SIF neuron are set to 0.75 and 7.5, respectively. The loss weight factor  $\lambda$  is set to 0.0003.

Table 1.  $AEE_{event}$  comparison with previous works

$AEE_{event}$	dt=1 frame				dt=4 frame			
	ind1	ind2	ind3	out1	ind1	ind2	ind3	out1
Zhu et al.'19	0.58	1.02	0.87	<b>0.32</b>	2.18	3.85	3.18	1.30
EV-FlowNet	1.03	1.72	1.53	0.49	2.25	4.05	3.45	1.23
Spike-FlowNet	0.84	1.28	1.11	0.49	2.24	3.83	3.18	<b>1.09</b>
Fusion-FlowNet	<b>0.56</b>	<b>0.95</b>	<b>0.76</b>	0.59	<b>1.68</b>	<b>3.24</b>	<b>2.43</b>	1.17

### 4.2. Evaluation of Optical Flow

For evaluation, the center cropped images of  $256 \times 256$  size are taken from indoor\_flying1,2,3 and outdoor\_day1 sequences. For indoor\_flying sequences, events and grayscale frames corresponding to the entire sequences are used for evaluation. However, for outdoor\_day1 sequence, 800 grayscale frames and the associated event streams are used for evaluation as suggested in [15, 32]. For quantitative results, we calculate the standard Average End-point Error (AEE) which is the mean Euclidean distance between the estimated flow ( $y_{estim}$ ) and the provided ground-truth ( $y_{gt}$ ). In our work, we measure the two types of AEE results: (1) over all pixels ( $AEE_{all}$ ) and (2) over pixels where events are present within the time-window ( $AEE_{event}$ ).

$$AEE = \frac{1}{m} \sum_m \|(u, v)_{estim} - (u, v)_{gt}\|_2 \quad (11)$$

where  $m$  indicates the count of active pixels in the event frames for  $AEE_{event}$  and every pixels of images for  $AEE_{all}$ .

### 4.3. Results

We compare Fusion-FlowNet with previous state-of-the-art works [32, 34, 15] in terms of the performance of optical flow prediction. As listed in Table 1, only  $AEE_{event}$  results are compared here since other works do not provide AEE values for dense optical flow estimation. We observe that Fusion-FlowNet outperforms other implementations in almost all scenarios. The outdoor\_day1 sequence is known to have suffered from certain issues with its grayscale images during dataset creation, leading to anomalous results for  $AEE_{event}$  as well as  $AEE_{all}$ . We report the results for it to maintain completeness in terms of comparison with previous works. Fig. 5 visualizes the predicted flow for this work and compares it with previous state-of-the-art methods.

### 4.4. Ablation studies

#### 4.4.1 Architectural Variations

We perform an ablation study to analyze the effect of architectural variations on model performance and efficiency. We investigate a second architecture where the dual pathway branches are extended to residual blocks. As shown in Fig. 6, we denote the first architecture as

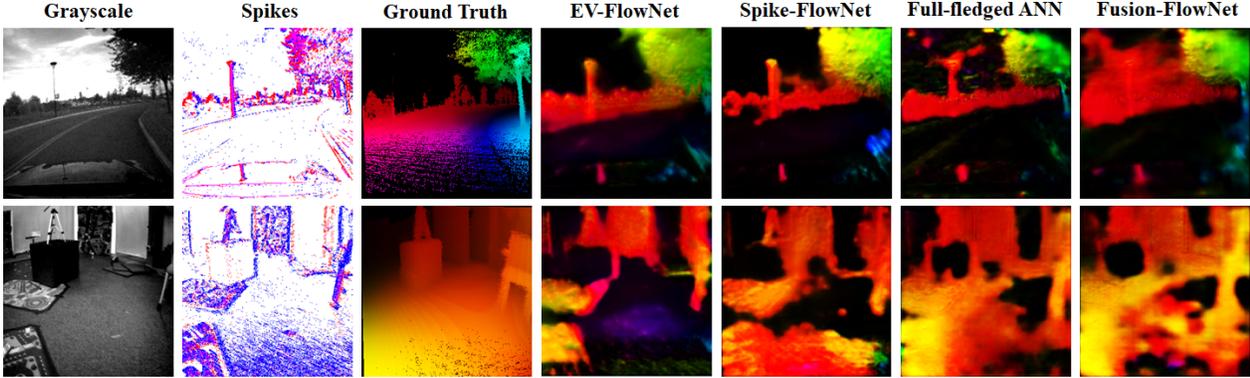


Figure 5. Predicted optical flow compared with other state-of-the-art methods. EV-FlowNet [32] and Spike-FlowNet [15] use only the event stream as input. Full-fledged ANN uses only grayscale images as input. Fusion-FlowNet uses both the event stream as well as grayscale images. The samples are taken from (top) *outdoor\_day1* and (bottom) *indoor\_flying2*. Best viewed in color.

Table 2. Average Endpoint Error (AEE) results for ablation studies

	indoor1		indoor2				indoor3				outdoor1					
	dt=1		dt=4		dt=1		dt=4		dt=1		dt=4		dt=1		dt=4	
	event	all														
Fusion-FlowNet <sub>Early</sub>	<b>0.56</b>	<b>0.62</b>	<b>1.68</b>	<b>1.81</b>	<b>0.95</b>	<b>0.89</b>	<b>3.24</b>	<b>2.90</b>	<b>0.76</b>	<b>0.85</b>	<b>2.43</b>	<b>2.46</b>	0.59	1.02	1.17	<b>3.06</b>
Fusion-FlowNet <sub>Late</sub>	0.57	0.63	1.71	1.89	0.99	0.92	3.26	2.93	0.79	0.87	2.46	2.54	0.55	1.00	1.34	3.48
Fusion <sub>Early</sub> [IF model]	<b>0.56</b>	<b>0.62</b>	1.72	1.93	0.97	0.90	3.36	3.07	0.78	0.87	2.51	2.63	0.58	1.04	1.37	3.52
Fusion <sub>Late</sub> [IF model]	0.57	0.64	1.71	1.90	1.00	0.93	3.41	3.08	0.80	0.88	2.56	2.64	0.55	<b>0.99</b>	1.38	3.53
Spike-FlowNet	0.84	0.91	2.24	2.94	1.28	1.23	3.83	4.09	1.11	1.20	3.18	3.92	<b>0.49</b>	1.42	<b>1.09</b>	3.28
Full-fledged ANN	0.60	0.68	1.73	1.90	1.00	0.97	3.35	3.03	0.83	0.97	2.52	2.62	0.83	1.53	1.27	3.19

Fusion-FlowNet<sub>Early</sub> and the second architecture as Fusion-FlowNet<sub>Late</sub>. Rows 1–2 in Table 2 highlight the optical flow prediction capability of both the architectures. We find that Fusion-FlowNet<sub>Early</sub> outperforms Fusion-FlowNet<sub>Late</sub> in predicting accurate optical flow outputs. Fusion-FlowNet<sub>Early</sub> contains comparatively larger number of parameters and fuses the intermediate features from the ANN/SNN branches in early layers, leading to better AEE results. On the other hand, Fusion-FlowNet<sub>Late</sub> performs the fusion at later layers leading to promising advantages in further reducing the network parameters and computational cost, as shown in Table 3.

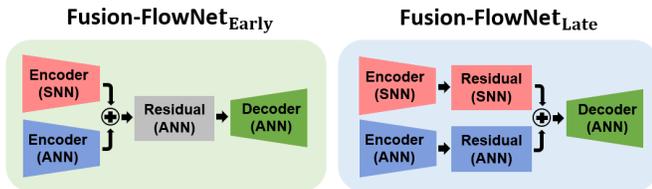


Figure 6. Architectures of (left) Fusion-FlowNet<sub>Early</sub> and (right) Fusion-FlowNet<sub>Late</sub>. Best viewed in color.

#### 4.4.2 Neuron Model Choice

For investigating the benefits of proposed SIF neuron, we compare the variations of Fusion-FlowNet with the SNN

blocks composed of SIF and IF neuron models. Rows 3–4 in Table 2 provide the AEE results for Fusion-FlowNet with IF neurons in the SNN layers. A comparison with results in rows 1–2 show that networks using SIF model can predict more accurate flow outputs compared to networks using IF model. This establishes the benefit of the SIF model towards mitigating the “dead neuron” problem in deep SNN layers.

#### 4.4.3 Sensor Fusion

We study the usefulness of sensor fusion approach against single sensor approaches using inputs as either the event streams or frame-based images. For the event only approach, we investigate Spike-FlowNet [15], a hybrid neural architecture where the initial layers are composed of SNNs and the deeper layers are composed of ANNs. Note, Spike-FlowNet utilizes the similar event-based input representation scheme and unsupervised learning method, providing a fair comparison. For the frame-based image only approach, we implement a custom full-fledged ANN architecture that resembles the U-Net [25] architecture, and train it with the equivalent unsupervised method as Fusion-FlowNet.

Rows 5–6 of Table 2 summarize the results for the single sensor approaches. Unsurprisingly, both Fusion-FlowNet<sub>Early</sub> and Fusion-FlowNet<sub>Late</sub> achieve better AEE

Table 3. Comparison of number of parameters and computational energy cost for different architectures for dt=1 and dt=4 cases. Lowest numbers highlighted in bold. (results averaged over all indoor and outdoor1 sequences)

	#Parameters ( $\times 10^6$ )		#OPS <sub>ANN</sub> ( $\times 10^9$ )		Spiking Activity (%)		#OPS <sub>SNN</sub> ( $\times 10^6$ )		E <sub>Total</sub> (mJ)		Improvement	
	dt=1	dt=4	dt=1	dt=4	dt=1	dt=4	dt=1	dt=4	dt=1	dt=4	dt=1	dt=4
Full-fledged ANN	13.044	13.046	5.339	5.367	–	–	–	–	24.536	24.666	1.00×	1.00×
Spike-FlowNet	13.039	13.039	4.409	4.409	0.480	1.008	15.81	195.99	20.296	20.458	1.21×	1.21×
Fusion-FlowNet_Early	12.269	12.270	4.648	4.648	0.173	<b>0.174</b>	<b>1.03</b>	<b>4.18</b>	21.381	21.384	1.15×	1.15×
Fusion-FlowNet_Late	<b>7.549</b>	<b>7.550</b>	<b>2.849</b>	<b>2.849</b>	<b>0.147</b>	0.179	5.24	6.44	<b>13.113</b>	<b>13.114</b>	<b>1.87×</b>	<b>1.88×</b>

performances in dt=1 and dt=4 scenarios compared to single sensor approaches. This verifies that the proposed fusion approach benefits from utilizing the complementary characteristics of event- and frame-based images, leading to better performance in both slow- and fast-motion scenarios. Furthermore, in comparison to prior works as listed in Table 1, both fusion options provide superior AEE results.

#### 4.5. Computational Efficiency

We validate the efficiency of Fusion-FlowNet in terms of the number of network parameters and computational energy cost for inference. Table 3 provides a detailed analysis on computational efficiency along with the comparison with previously discussed alternate architectures.

We observe that both Fusion-FlowNet<sub>Early</sub> and Fusion-FlowNet<sub>Late</sub> contain fewer number of parameters compared to a full-fledged ANN architecture and Spike-FlowNet. This is due to the usage of narrow convolution layers which greatly reduce the number of parameters and computations. In particular, Fusion-FlowNet<sub>Late</sub> contains the least number of network parameters ( $\sim 58\%$  compared to full-fledged ANN) as the residual blocks contain the majority of the parameters and utilizing narrow convolutional layers for them helps reduce the total network parameters drastically.

For estimating the computational energy cost for different architectures, we first describe how computations in SNNs and ANNs differ from each other. Conceptually, SNNs perform highly sparse asynchronous ACcumulate (AC) operations over time. These synaptic operations are executed only at the arrival of input spikes due to the nature of binary-valued inputs. In contrast, ANNs perform expensive Multiply-and-ACcumulate (MAC) operations for computing dense Matrix-Vector Multiplications (MVMs). Based on the findings in [10], a MAC operation requires a total of  $E_{MAC}=4.6pJ$  of energy while an AC operation requires only  $E_{AC}=0.9pJ$  for a 32-bit floating-point computation (45nm CMOS technology). This leads to the AC operation being  $5.1\times$  more energy-efficient compared to the MAC operation. These findings coupled with the number of synaptic operations are commonly used to benchmark the computational energy cost of SNNs. [21, 26, 15].

Next, we calculate the total number of synaptic operations for every layer. In the SNN layers, the number of synaptic operations are obtained by multiplying the pre-

spike activities, the number of synaptic connections and the number of time-steps. Also, the computational energy of AC and MAC computations are taken into consideration for SNNs and ANNs, respectively. The total computational energy cost can be formalized as:

$$\#OPS_{SNN} = N \sum_l M_l C_l F_l, \#OPS_{ANN} = \sum_l M_l C_l \quad (12)$$

$$E_{Total} = \#OPS_{SNN} \times E_{AC} + \#OPS_{ANN} \times E_{MAC} \quad (13)$$

where  $M$  is the number of neurons,  $C$  is the number of synaptic connections,  $F$  represents mean spiking activity,  $N$  is the number of timesteps,  $\#OPS_{SNN}/\#OPS_{ANN}$  indicate the number of operations for SNN/ANN portions, and  $E_{Total}$  denotes the total computational energy cost.

The last column in Table 3 provides the overall improvement in computational energy cost. We observe that Fusion-FlowNet<sub>Late</sub> demonstrates the highest improvement in energy ( $\sim 1.88\times$ ) compared to full-fledged ANN. This is because more layers, including encoder and residual blocks, utilize narrow convolutions, leading to reduction in the number of parameters and consequently reduction in the total computational energy cost. Furthermore, the SNN pathway contributes negligibly to the total computational energy cost compared to ANN pathway.

## 5. Conclusion

We propose a sensor/architecture fusion framework for accurately estimating optical flow in challenging environments. We leverage the complementary characteristics of event- and frame-based sensors as well as ANNs and SNNs. Our framework (Fusion-FlowNet) reports state-of-the-art optical flow prediction results, while substantially reducing network parameters and computational energy cost. This work contributes two different deep fused architectures (Fusion-FlowNet<sub>Early</sub> and Fusion-FlowNet<sub>Late</sub>), having different applications of interest. Fusion-FlowNet<sub>Early</sub> provides highly accurate dense optical flow, proving to be appropriate for safety-critical applications. While, Fusion-FlowNet<sub>Late</sub> promises immense benefits in terms of computational efficiency, making it suitable for the edge applications on resource-constrained hardware.

## References

- [1] Myo Tun Aung, Rodney Teo, and Garrick Orchard. Event-based plane-fitting optical flow for dynamic vision sensors in fpga. In *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 1–5, May 2018.
- [2] Steven S. Beauchemin and John L. Barron. The computation of optical flow. *ACM computing surveys (CSUR)*, 27(3):433–466, 1995.
- [3] Ryad Benosman, Charles Clercq, Xavier Lagorce, Sio-Hoi Ieng, and Chiara Bartolozzi. Event-based visual flow. *IEEE Transactions on Neural Networks and Learning Systems*, 25(2):407–417, Feb 2014.
- [4] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A  $240 \times 180$  130 db  $3 \mu\text{s}$  latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, Oct 2014.
- [5] Anthony N Burkitt. A review of the integrate-and-fire neuron model: I. homogeneous synaptic input. *Biological cybernetics*, 95(1):1–19, 2006.
- [6] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautam Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, Yuyun Liao, Chit-Kwan Lin, Andrew Lines, Ruokun Liu, Deepak Mathaikutty, Steven McCoy, Arnab Paul, Jonathan Tse, Guruguhannathan Venkataramanan, Yi-Hsin Weng, Andreas Wild, Yoonseok Yang, and Hong Wang. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, 38(1):82–99, January 2018.
- [7] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 2758–2766, 2015.
- [8] Guillermo Gallego, Tobi Delbruck, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi, Stefan Leutenegger, Andrew Davison, Jörg Conradt, Kostas Daniilidis, et al. Event-based vision: A survey. *arXiv preprint arXiv:1904.08405*, 2019.
- [9] Guillermo Gallego, Henri Rebecq, and Davide Scaramuzza. A unifying contrast maximization framework for event cameras, with applications to motion, depth, and optical flow estimation. *CoRR*, abs/1804.01306, 2018.
- [10] Mark Horowitz. 1.1 computing’s energy problem (and what we can do about it). In *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, pages 10–14. IEEE, 2014.
- [11] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [12] Max Jaderberg, Karen Simonyan, Andrew Zisserman, and Koray Kavukcuoglu. Spatial transformer networks. In C. Cortes, N. Lawrence, D. Lee, M. Sugiyama, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 28. Curran Associates, Inc., 2015.
- [13] J Yu Jason, Adam W Harley, and Konstantinos G Derpanis. Back to basics: Unsupervised learning of optical flow via brightness constancy and motion smoothness. In *European Conference on Computer Vision*, pages 3–10. Springer, 2016.
- [14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*, pages 1097–1105, 2012.
- [15] Chankyu Lee, Adarsh Kosta, Alex Zihao Zhu, Kenneth Chaney, Kostas Daniilidis, and Kaushik Roy. Spike-flownet: Event-based optical flow estimation with energy-efficient hybrid neural networks. In *European Conference on Computer Vision*, pages 366–382. Springer, 2020.
- [16] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A  $128 \times 128$  120 db  $15 \mu\text{s}$  latency asynchronous temporal contrast vision sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, Feb 2008.
- [17] Pengpeng Liu, Michael Lyu, Irwin King, and Jia Xu. Self-low: Self-supervised learning of optical flow. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4571–4580, 2019.
- [18] Bruce D. Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. In *Proceedings of the 7th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI’81*, pages 674–679. San Francisco, CA, USA, 1981. Morgan Kaufmann Publishers Inc.
- [19] Misha Mahowald. The silicon retina. In *An Analog VLSI System for Stereoscopic Vision*, pages 4–65. Springer, 1994.
- [20] Iacopo Masi, Yue Wu, Tal Hassner, and Prem Natarajan. Deep face recognition: A survey. In *2018 31st SIBGRAPI conference on graphics, patterns and images (SIBGRAPI)*, pages 471–478. IEEE, 2018.
- [21] Paul A Merolla, John V Arthur, Rodrigo Alvarez-Icaza, Andrew S Cassidy, Jun Sawada, Filipp Akopyan, Bryan L Jackson, Nabil Imam, Chen Guo, Yutaka Nakamura, et al. A million spiking-neuron integrated circuit with a scalable communication network and interface. *Science*, 345(6197):668–673, 2014.
- [22] Manjunath Narayana, Allen Hanson, and Erik Learned-Miller. Coherent motion segmentation in moving camera videos using optical flow orientations. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1577–1584, 2013.
- [23] Liyuan Pan, Miaomiao Liu, and Richard Hartley. Single image optical flow estimation with an event camera. *arXiv preprint arXiv:2004.00347*, 2020.
- [24] Zhe Ren, Junchi Yan, Bingbing Ni, Bin Liu, Xiaokang Yang, and Hongyuan Zha. Unsupervised deep learning for optical flow estimation. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [25] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. *CoRR*, abs/1505.04597, 2015.
- [26] Bodo Rueckauer, Iulia-Alexandra Lungu, Yuhuang Hu, Michael Pfeiffer, and Shih-Chii Liu. Conversion of continuous-valued deep networks to efficient event-driven networks for image classification. *Frontiers in neuroscience*, 11:682, 2017.

- [27] Deqing Sun, Stefan Roth, and Michael J. Black. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *Int. J. Comput. Vision*, 106(2):115–137, Jan. 2014.
- [28] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8934–8943, 2018.
- [29] Heng Wang, Alexander Kläser, Cordelia Schmid, and Cheng-Lin Liu. Action recognition by dense trajectories. In *CVPR 2011*, pages 3169–3176. IEEE, 2011.
- [30] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990.
- [31] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.
- [32] Alex Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Ev-flownet: Self-supervised optical flow estimation for event-based cameras. In *Proceedings of Robotics: Science and Systems*, Pittsburgh, Pennsylvania, June 2018.
- [33] Alex Zihao Zhu, Dinesh Thakur, Tolga Özaslan, Bernd Pfrommer, Vijay Kumar, and Kostas Daniilidis. The multivehicle stereo event camera dataset: An event camera dataset for 3d perception. *IEEE Robotics and Automation Letters*, 3(3):2032–2039, 2018.
- [34] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised event-based learning of optical flow, depth, and egomotion. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 989–997, 2019.