# Event-based Video Reconstruction via Potential-assisted Spiking Neural Network

Lin Zhu[1,2], Xiao Wang[2], Yi Chang[2], Jianing Li[1,2], Tiejun Huang[1], Yonghong Tian[1,2,*]

Peking University[1], Peng Cheng Laboratory[2]

## Abstract

*Neuromorphic vision sensor is a new bio-inspired imaging paradigm that reports asynchronous, continuously per-pixel brightness changes called 'events' with high temporal resolution and high dynamic range. So far, the event-based image reconstruction methods are based on artificial neural networks (ANN) or hand-crafted spatiotemporal smoothing techniques. In this paper, we first implement the image reconstruction work via deep spiking neural network (SNN) architecture. As the bio-inspired neural networks, SNNs operating with asynchronous binary spikes distributed over time, can potentially lead to greater computational efficiency on event-driven hardware. We propose a novel Event-based Video reconstruction framework based on a fully Spiking Neural Network (EVSNN), which utilizes Leaky-Integrate-and-Fire (LIF) neuron and Membrane Potential (MP) neuron. We find that the spiking neurons have the potential to store useful temporal information (memory) to complete such time-dependent tasks. Furthermore, to better utilize the temporal information, we propose a hybrid potential-assisted framework (PA-EVSNN) using the membrane potential of spiking neuron. The proposed neuron is referred as Adaptive Membrane Potential (AMP) neuron, which adaptively updates the membrane potential according to the input spikes. The experimental results demonstrate that our models achieve comparable performance to ANN-based models on IJRR, MVSEC, and HQF datasets. The energy consumptions of EVSNN and PA-EVSNN are 19.36× and 7.75× more computationally efficient than their ANN architectures, respectively. The code and pretrained model are available at https://sites.google.com/view/evsnn.*

## 1. Introduction

Event cameras [2, 5] are bio-inspired vision sensors that pose a paradigm shift in the way visual information is acquired. Compared with standard cameras, event cameras have high temporal resolution, high dynamic range (140 dB vs. 60 dB of standard cameras), and low power consumption. Event cameras work asynchronously, recording the stream of events $(t, x, y, p)$ which includes the timestamp, pixel location and polarity of the brightness changes.

Despite the advantages of the event data, it is not friendly to human vision and traditional computer vision [48, 53]. As a solution, image reconstruction bridges the gap between human visualization and events, giving us an intuition of the rich information encoded by events. On other hand, image is a useful representation for conventional frame-based computer vision [42]. Reconstructing images from asynchronous events has been explored in various researches. Early works attempt to recover the intensity of an image from events based on hand-crafted priors [46, 47, 3, 30]. Recently, deep neural network based reconstruction models [52, 42, 43, 48, 50, 34, 55] have demonstrated impressive performance. The events usually be transformed in to time-surfaces, event images or voxel grids as the input of convolutional neural network. However, large artificial neural networks (ANN) can be memory and computationally intensive [48], consuming power and hampering the low latency of event cameras.

In fact, the sparse event data can be effectively combined with neuromorphic hardware for low-power spiking neural network (SNN) applications [14]. Compared with ANN, SNN is more biologically realistic and its neurons communicate with each other via discrete spikes instead of continuous-valued activations. Visual systems [32, 1] constructed with SNN and event cameras have demonstrated their capacity in solving visual tasks as well as prominent energy-efficiency. However, most of the SNN work has so far been focused on problems like classification [10, 59, 62], optical estimation [35, 15], motion segmentation [33], and angular velocity regression [11]. To the best of our knowledge, we are the first to attempt image reconstruction task based on a deep SNN architecture.

In this paper, we propose a novel Event-based Video reconstruction framework based on a fully Spiking Neural Network (EVSNN), which utilizes the Leaky-Integrate-and-Fire (LIF) neurons and a Membrane Potential (MP) neuron. To better extract the temporal information, we propose a hybrid potential-assisted framework (PA-EVSNN) using the

---

*Corresponding author.

membrane potential of spiking neurons. The main contributions of this paper are summarized as follows:

1) We first explore a fully spiking neural network (EVSNN) architecture on event-based image reconstruction, which utilizes LIF neuron and MP neuron. This is also the first attempt to develop a deep SNN for image reconstruction task.

2) We propose a hybrid potential-assisted SNN (PA-EVSNN), which uses adaptive membrane potential (AMP) neurons to improve the temporal receptive field of EVSNN. AMP neurons can adjust the membrane time constant according to the input spike to adapt to various reconstruction scenes.

3) The experiments on public datasets demonstrate that the proposed models have comparable performance to existing ANN-based models, while the energy consumptions of EVSNN and PA-EVSNN are $19.36\times$ and $7.75\times$ more computationally efficient than their ANN architectures, respectively. Compared to E2VID, the proposed EVSNN and PA-EVSNN achieve $24.15\times$ and $8.76\times$ more computationally efficient improvement, respectively.

## 2. Related work

**Spiking Neural Network** Supervised learning of SNNs was first proposed by SpikeProp [4], it used a linear approximation to overcome the non-differentiable threshold-triggered firing mechanism of SNNs, backpropagation was utilized to update weight. Some works applied to single-layer SNN optimization appeared, including Tempotron [13], Re-SuMe [37], and SPAN [28]. Recently, the surrogate gradient method provides an effective solution for training multi-layer SNN [25, 18, 57, 49, 24, 19]. It utilized surrogate derivatives to define the derivative of the threshold-triggered firing mechanism. Therefore, SNN can be optimized by gradient descent algorithm like ANNs, which makes the training of deep SNN possible.

Most of the learning-based SNN work has so far been focused on problems like classification [10, 59, 62], optical estimation [35, 15, 15], motion segmentation [33], and angular velocity regression [11]. There are also some unsupervised SNNs [67, 63] proposed for image reconstruction based on spike camera [66]. Among them, [15] and [38] utilized deep SNNs for optical and depth estimation, respectively. In addition, Lee *et al.* [23] proposed an ANN-SNN hybrid architecture for optical estimation, using SNN as encoder and ANN as decoder and residual block. Zhang *et al.* [61] proposed the ANN-SNN hybrid network for event-based synthetic aperture imaging.

**Event-based Video Reconstruction** Video reconstruction is an important topic in event-based vision field. Early reconstruction works are based on hand-crafted features to estimate intensity from events, e.g. optimization [3], regularization [30] and temporal filtering [46, 47]. Some
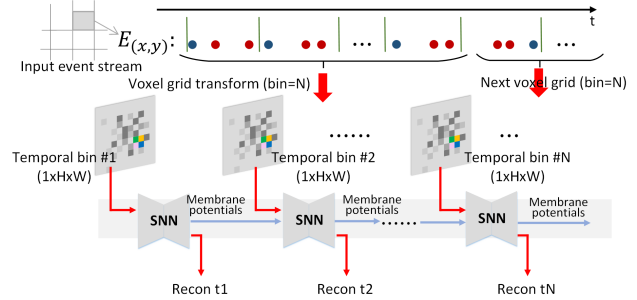


Figure 1. **The event representation and work flow of our framework.** The event stream (red/blue dots represent on/off events, respectively) is split into multiple windows and transformed into continuous voxel grids. Each voxel grid includes $N$ temporal bins with different information. Our SNN recurrent uses the current single channel temporal bin and last membrane potentials of each spiking neuron to generate new reconstructions at each moment.

works [7, 20, 41] also applied SLAM to estimate the brightness. Recently, deep learning methods have shown impressive performance on event-based video reconstruction. Wang *et al.* [52] utilized generative adversarial network (GAN) to reconstruct intensity with real grayscale frames. Rebecq *et al.* [42, 43] proposed an effective E2VID model which based on a U-Net [44] model. The network was trained in a supervised manner with a synthetic dataset generated from ESIM [40]. Scheerlinck *et al.* [48] proposed a light-weight framework to achieve fast inference speed with only a minor drop in accuracy. Stoffregen *et al.* [50] proposed to use more complex synthetic dataset to train the network, bringing a large performance boost on real datasets. Federico *et al.* [34] proposed a novel self-supervised learning method for image reconstruction, getting rid of training data. Weng *et al.* [55] presented a hybrid CNN-transformer network for image reconstruction.

In this work, different from the above ANN-based models, we first propose to use the energy-efficient deep SNN models for reconstructing videos from event stream.

## 3. Method

### 3.1. Input Representations

To process the asynchronous event with SNN, the event data is needed to be converted into an event representation which includes the temporal information. In this work, we use the continuous voxel grid [65] to train and test our model, which is defined as: $E(x, y, t_n) = \sum_i p_i \max(0, 1 - |t_n - t_i^*|)$, where $t_i^* = \frac{B-1}{\Delta T}(t_i - t_0)$, $t_i^*$ is the normalized event timestamp. As shown in Fig. 1, the event stream can be adaptively divided into continuous bins of voxel grid.

### 3.2. Spiking Neurons

ANN and SNN can model the same types of network topologies, but SNN replaces the artificial neuron model

with a spiking neuron model. The artificial neuron model operates on a weighted sum of inputs, and passing the result through a sigmoid or ReLU nonlinearity. In SNN, the weighted sum of inputs contributes to the membrane potential of the spiking neuron. If the membrane potential of the spiking neuron reaches a threshold, then the neuron will emit a spike to its subsequent connections. The information in SNN is propagated by discrete spikes,thus spiking neuron is the basic computing unit.

**LIF Neurons** The Leaky Integrate-and-Fire (LIF) model [12] is a widely used neuron model in SNN, which is more biologically realistic than the Integrate-and-Fire (IF) neuron model. The subthreshold dynamics of LIF neuron is defined as $\tau \frac{\mathrm{d}V(t)}{\mathrm{d}t} = -(V(t) - V_{rest}) + X(t)$, where $V(t)$ represents the membrane potential of the neuron at time $t$, $X(t)$ represents the input to neuron, $\tau$ is the membrane time constant. A spike fires if $V(t)$ exceeds the threshold $V_{th}$, $V_{rest}$ is the resting potential after firing. For better representation, we rewrite the above equation as the discrete form:

$$\begin{cases} V_t = V_{t-1} + \frac{1}{\tau}(-(V_{t-1} - V_{rest}) + X_t) \\ S_t = H(V_t - V_{th}) \end{cases} \quad (1)$$

where $V_t$ denotes the membrane potential after neuronal dynamics at $t$. $S_t$ denotes the spike output at $t$, $H(\cdot)$ denotes the Heaviside step function which is defined as $H(x) = 1$ for $x \geq 0$ and $H(x) = 0$ for $x < 0$. We set $V_{rest} = V_{reset}$ in our work. LIF neuron can extract temporal information during the integration and firing process, however, its output is binary spikes which can only represent limited information. Moreover, after each firing process, $M_t$ is reset thus the temporal information is also partially lost. Based on above, we introduce the membrane potential neurons.

**Membrane Potential Neurons** The membrane potential neurons (MP neurons) are non-spiking neurons which output membrane potential instead of spikes [51, 58]. In our image reconstruction task, MP neurons can extract more useful temporal information hidden in the neurons. The dynamics of the MP neurons is same as LIF neuron. For MP neurons, $M_t$ is equal to $V_t$ since there is no spike fire and $V_t$ reset process. If we set $V_{rest} = 0$, Eq. 1 can be written as:

$$\begin{cases} V_t = (1 - \frac{1}{\tau})V_{t-1} + \frac{1}{\tau}X_t \\ O_t = V_t \end{cases} \quad (2)$$

where $O_t$ denotes the output of the neuron at $t$. Eq. 2 is similar to the function of recurrent neural networks. The membrane time constant $\tau$ controls the balance between remembering $X_t$ and forgetting $V_{t-1}$. Thus it can be considered as a simple version of Long Short-Term Memory (LSTM) module [16].
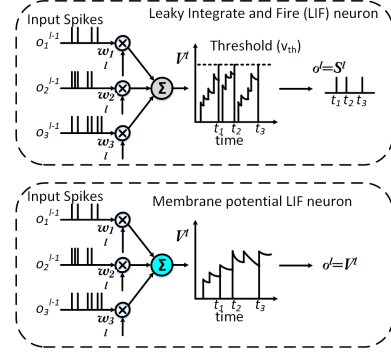


Figure 2. **The dynamics of LIF neuron and MP_LIF neuron.** For LIF neurons, if the membrane potential reaches a threshold, then the neuron will emit a spike to its subsequent connections and reset to resting state. At each time step, MP_LIF neuron outputs its membrane potential as the weighted sum of input spikes.

### 3.3. The Proposed SNN Model

In this paper, we propose two SNN architectures for event-based reconstruction, namely EVSNN and PA-EVSNN. EVSNN is a fully spiking neural network, all synaptic operations in the network are SNN operations. PA-EVSNN shares the same spiking encoder and decoder architecture, with the additional MP neurons to improve the performance. Both models are fully convolutional networks, the architecture is shown in Fig. 3.

**EVSNN (A fully spiking neural network)** Our EVSNN is a variant of the U-shaped model [44]. First, the event data is transformed into event voxels. For each time step, a $1 \times W \times H$ event voxel is fed in to EVSNN and transformed as the size of $N_c \times W_1 \times H_1$, followed by $N_e$ encoder layers, $N_r$ residual blocks, $N_d$ decoder layers, and a final image prediction layer. The number of channels is doubled after each encoder layer. All spiking neurons in encoder layers, decoder layers, and residual blocks are LIF neurons, which enables computationally efficiency. To ensure a fully SNN architecture, EVSNN utilizes concatenate as spike skip connection. In the final image prediction layer, MP_LIF neuron is introduced to integrate all spikes and predict the gray scale image. An ablation of each network component can be found in Sec. 4.4. We use $N_c = 32$, $N_e = N_d = 3$ and $N_r = 1$. EVSNN can handle most scenes in existing datasets while the computationally efficiency is 19.36 times than ANN architecture.

**PA-EVSNN (Potential-assisted EVSNN)** EVSNN is a fully SNN with very low energy consumption. However, the reconstruction performance is limited by the binary spikes (e.g., the gray scale of the image is not rich enough). Based on EVSNN, we further propose a potential-assisted EVSNN model. MP neuron is introduced in each encoder and decoder layer to help extract the temporal information hidden in the spikes. We also propose an adaptive membrane potential (AMP) neuron, which greatly enhances the tempo-
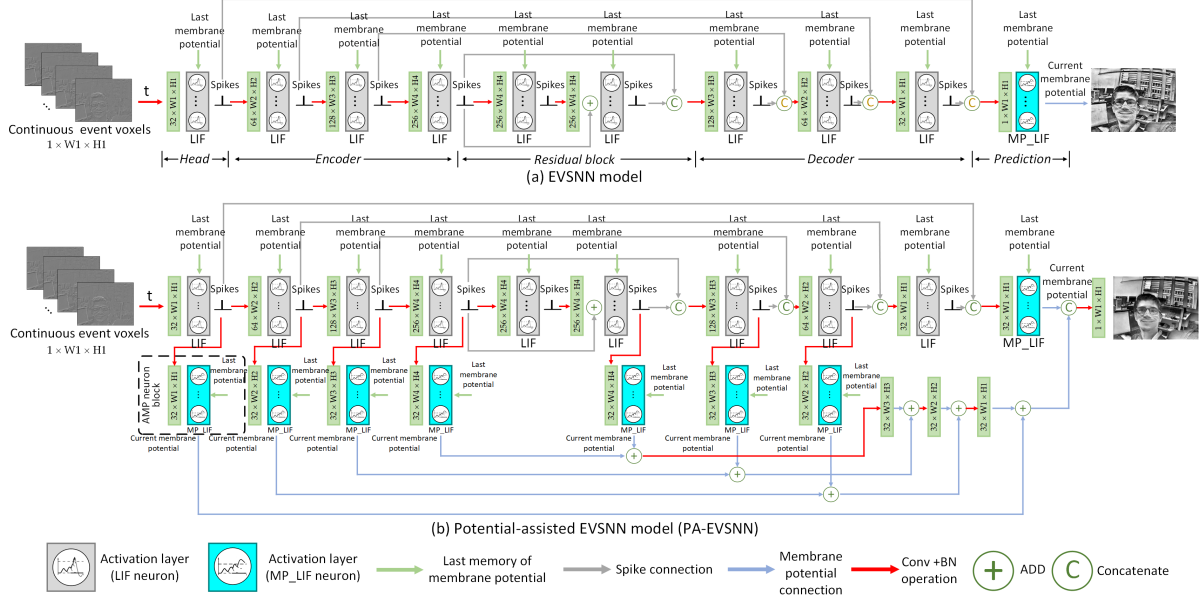
Figure 3. **The proposed spiking neural network architecture.** EVSNN is a fully spiking neural network, composed of head, encoder, residual block, decoder, and prediction layers. Based on EVSNN as backbone, PA-EVSNN introduces MP neurons to further improve the performance. MP_LIF denotes MP neurons with the dynamics of LIF. As floating-point multiplication operations are introduced by MP neurons, we consider PA-EVSNN as a hybrid network. The energy consumptions of EVSNN and PA-EVSNN are 19.36 × and 7.75 × more computationally efficient than their ANN architectures. Please refer to our supplementary material for details of network architecture.

ral receptive field of the network. Notice that although the backbone of PA-EVSNN is SNN architecture, the introduction of MP neuron brings non binary spikes in the network (about 8.4% ANN floating-point operations), thus we consider PA-EVSNN as a hybrid network. Compared to existing ANN models, PA-EVSNN still has great advantages in energy consumption (7.75 × more efficient) while achieving comparable performance. More detailed analysis of the SNN and ANN operations can be found in Sec. 4.5.

**Adaptive Membrane Potential Neurons** As analyzed in Sec. 3.2, the membrane time constant $\tau$ in MP neuron plays an analogous role as the gates in LSTM module. [10] proposed the parametric LIF neuron by introducing a learnable $\tau$ in classification task. Inspired by this, we propose an adaptive membrane potential neuron (AMP neuron). Different from the fixed $\tau$ of parametric LIF neuron learned from the training dataset, AMP neuron can adjust $\tau$ according to the input spike to adapt to various reconstruction scenes.

According to Eq. 2, ideally, when the light changes fast, the network should choose a large $\tau$ to remember more new information while forgetting more last memory, and vice versa. Due to the event measures the change of light intensity, the average spike firing rate reflects the global motion of the scene to a certain extent, which is useful for estimating a proper $\tau$. The spike firing rate of each channel in $l$-th layer can be estimated by $F = \mathrm{AvgPool}(S_l)$, where $\mathrm{AvgPool}(\cdot)$ denotes the average pooling operation, $S_l$ is the spike tensor of $l$-th layer. Then the local motion intensity of input spikes can be estimated by $I = \mathrm{MaxPool}(\mathrm{Conv}(S_l))$,
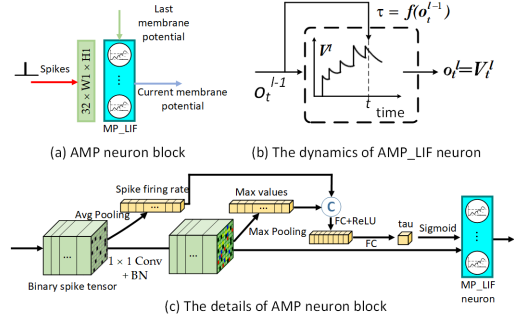


Figure 4. **The adaptive membrane potential (AMP) neuron block.** (a) An AMP neuron block in Fig. 3. (b) The dynamics of AMP_LIF neuron. The membrane time constant $\tau$ is adjusted by the input $o_t^{l-1}$. (c) The details of AMP neuron block.

where $\mathrm{MaxPool}(\cdot)$ denotes the max pooling operation. Finally, the membrane time constant is updated by

$$\tau = \frac{1}{\mathcal{S}(\mathrm{Linear}([F, I]))} \tag{3}$$

where $\mathcal{S}(\cdot)$ denotes the sigmoid activation function, $\mathrm{Linear}(\cdot)$ is the full connection layer shown in Fig. 4.

**Loss Functions** We use LPIPS loss and temporal consistency loss: $\mathcal{L}_{total} = \sum_{k=0}^{L} \mathcal{L}_k^R + \lambda \sum_{k=L_0}^{L} \mathcal{L}_k^{TC}$, where $\mathcal{L}_k^R$ is the LPIPS loss [60], $\mathcal{L}_k^{TC}$ is temporal consistency loss [22, 43].

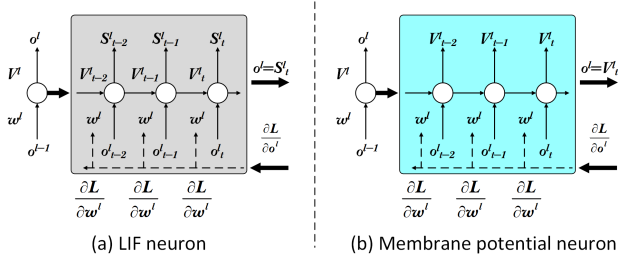**(a) LIF neuron** **(b) Membrane potential neuron**

Figure 5. **The backpropagation of spiking neurons.** For LIF neurons, we use ArcTan as surrogate function to calculate the derivative of spiking function. For MP neurons, the gradients can be directly computed by Eq. 4.

### 3.4. Training Details of SNN

During the training process, we set $L$ in loss function same as the training sequence length (i.e., 40 - 60), and $L_0$ is set as 2. In each time step, an event voxel with the size of $1 \times H_1 \times W_1$ is fed into the network. According to Eq. 4, the backpropagated errors pass through the spiking neuron layer and MP neuron layer using BackPropagation Through Time (BPTT) [56]. In BPTT, the network is unrolled for all discrete time-steps. The loss is calculated every 5 time steps, and the weight update is computed as the sum of gradients from each time-step as follows:

$$\Delta w^l = \sum_n \frac{\partial \mathcal{L}_{total}}{\partial o_t^l} \frac{\partial o_t^l}{\partial V_t^l} \frac{\partial V_t^l}{\partial w^l} \quad (4)$$

$$\text{where} \quad \frac{\partial o_t^l}{\partial V_t^l} = \begin{cases} H_1'(V_t - V_{th}) & \text{if } o_t^l = S_t^l \\ 1 & \text{if } o_t^l = V_t^l \end{cases}$$

where $o_t^l$ is the output of the neuron at time $t$, $\frac{\partial o_t^l}{\partial V_t^l}$ denotes the derivative of spike with respect to the membrane potential after charging at time step $t$. Since $\frac{\partial o_t^l}{\partial V_t^l}$ is not differentiable, we adopt surrogate gradient method [31] to calculate it. The shifted ArcTan function $H_1(x) = \frac{1}{\pi} \arctan(\pi x) + \frac{1}{2}$ is utilized as the surrogate function of the Heaviside step function $H(\cdot)$. If the neuron is a spiking neuron, we have . Otherwise, if the neuron is a MP neuron with no spiking output, then $o_t^l = V_t^l$, we have $\frac{\partial o_t^l}{\partial V_t^l} = 1$ which is similar to an ANN activation function.

## 4. Experiments

### 4.1. Experimental setup

For fair comparison to ANN-based reconstruction methods, we use the exact same synthetic data from E2VID [43] to train our SNN. The dataset is generated by ESIM, an event simulator, and consists of 950 training sequences and 50 validation sequences. MS-COCO images [26] are mapped to a 3D plane and random 6-DOF camera motions are used to trigger events. During training, the nonzero values of event tensors are normalized as mean and standard deviation is 0 and 1, respectively. The events and images are randomly cropped to $128 \times 128$ to augment the data.

Our models are implemented with SpikingJelly [9], an open-source deep learning framework for SNNs based on PyTorch [36]. An NVIDIA TITAN Xp GPU is used to train our model. We adopt the batch size of 8 and Adam optimizer [21] in the training process. The network is trained for 100 epochs, with a learning rate of 0.002. The weight $\lambda$ of temporal consistency loss is set as 1. The reset value $V_{reset}$ of all neurons is set to 0, and the membrane time constant $\tau$ of LIF neurons is set to 2.

### 4.2. Evaluation on Public Datasets

We evaluate our model on three public datasets IJRR [29], MVSEC [64], and HQF [50]. Following [43] and [48], to ensure the intensity values lay within a similar range, we apply histogram normalization to both the output and groundtruth frames. Moreover, to make the timestamps of the reconstruction and groundtruth strictly consistent, we use the events between two adjacent frames to generate each reconstruction. We compare our models with four state-of-the-art-methods E2VID [43], FireNet [48], SPADE-E2VID [6], and E2VID+ [50]. All results are generated by the pre-trained model from the original paper. We compared reconstructed images against groundtruths using the metrics: mean squared error (MSE), structural similarity (SSIM) [54] and perceptual similarity (LPIPS) [60].

The main quantitative results are presented in Table 1. Notice that E2VID+ is trained on a more challenging synthetic dataset, while the other five methods are trained on the same data as [43]. To sum up, E2VID+ performs best in most datasets. SPADE-E2VID performs well on MVSEC dataset, but the SSIM and LPIPS are lower than E2VID on IJRR and HQF datasets. The results show that EVSNN can handle these scenes. Our PA-EVSNN achieves comparable performance to ANN-based models such as E2VID and FireNet. Please refer to our supplementary material for additional quantitative and qualitative results.

### 4.3. Temporal Component Ablation

Inspired by [43], we design an experiment to measure the effective size of the temporal receptive field of SNN and ANN. As shown in Fig.7, four different settings are tested: ANN w/o recurrent, ANN + LSTM (E2VID), SNN + LIF (EVSNN), and SNN + LIF + AMP_LIF. These networks are all based on U-Net architecture with three encoders. To verify the ability of spiking neurons in temporal information extraction, at initialization phase (T = 1-50), the states of temporal components (e.g., LSTM and spiking neuron) are initialized at zero. The images at each moment are reconstructed by continuous event input. Then we artificially stop the events at T=50. In subsequent iterations after T = 50, we feed the empty event tensors to the network and reconstruct images to test the effective size of temporal receptive field.

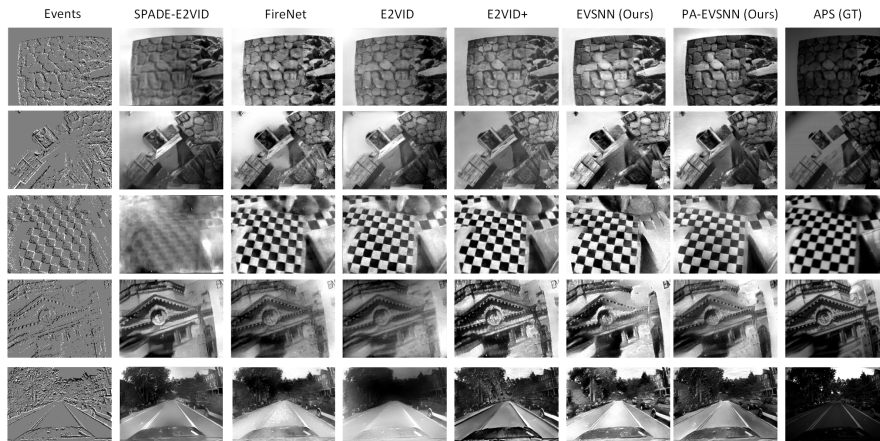| Events | SPADE-E2VID | FireNet | E2VID | E2VID+ | EVSNN (Ours) | PA-EVSNN (Ours) | APS (GT) |

Figure 6. **Qualitative comparison with state-of-the-art ANN-based methods.** We compare our SNN models with four ANN-based models (SPADE-E2VID, FireNet, E2VID, and E2VID+) on IJRR (Row 1-2), HQF (Row 3-4), and MVSEC (Row 5) datasets. The results show that the proposed EVSNN and PA-EVSNN perform comparably to most ANN-based models, and the energy consumptions are 24.15 times and 8.76 times lower than E2VID, respectively (see Table 7). More qualitative results can be found in our supplementary material.

Table 1. Comparison on IJRR, HQF, and MVSEC Datasets.

| Method | IJRR | | | MVSEC | | | HQF | | |
|---|---|---|---|---|---|---|---|---|---|
| | MSE ↓ | SSIM ↑ | LPIPS ↓ | MSE ↓ | SSIM ↑ | LPIPS ↓ | MSE ↓ | SSIM ↑ | LPIPS ↓ |
| *E2VID | 0.059 | **0.643** | <u>0.338</u> | 0.138 | 0.377 | 0.651 | 0.081 | <u>0.545</u> | 0.406 |
| *FireNet | 0.060 | 0.602 | 0.340 | 0.105 | 0.361 | 0.600 | 0.065 | 0.542 | <u>0.391</u> |
| *SPADE-E2VID | 0.063 | 0.572 | 0.365 | <u>0.095</u> | **0.443** | 0.556 | 0.080 | 0.512 | 0.424 |
| *[1]E2VID+ | **0.043** | 0.618 | **0.321** | **0.088** | <u>0.427</u> | **0.490** | **0.047** | **0.560** | **0.338** |
| †EVSNN (Ours) | 0.061 | 0.570 | 0.362 | 0.104 | 0.389 | <u>0.538</u> | 0.086 | 0.482 | 0.433 |
| †PA-EVSNN (Ours) | <u>0.046</u> | <u>0.626</u> | 0.367 | 0.107 | 0.403 | 0.566 | <u>0.061</u> | 0.532 | 0.416 |

[1] E2VID+ is trained on the simulated dataset proposed in [50] , while other five models are all trained on the simulated dataset from [43].

* ANN model. † SNN model. Notice that the energy consumption of SNN is much lower than that of ANN, see Table 7 for detail.

To better analyze the results, we randomly pick 50 event sequences from dynamic_6dof of IJRR dataset, and plot the average values of MSE, SSIM, LPIPS in Fig. 8. The spike firing rates of EVSNN and PA-EVSNN are also reported. The results shown in Fig 7 and 8 show that E2VID, EVSNN, and PA-EVSNN can complete initialization in 10 iterations. As T increases, the quantitative results of E2VID, EVSNN, and PA-EVSNN continue to improve, which shows the effectiveness of temporal component. In contrast, the quantitative scores of ANN w/o recurrent do not change significantly, since it has no temporal component. After T=50, the quantitative score of E2VID, EVSNN, and PA-EVSNN decreases slowly, indicating that our SNN model has temporal receptive field similar to ANN + LSTM. Assisted by the membrane potential, PA-EVSNN performs better than EVSNN on quantitative scores. To sum up, our SNN structure has the capacity of temporal information extracting although it may be weaker than ANN+LSTM.

### 4.4. Spiking Neural Network Architecture

We investigate different SNN network architectures. All the experiments are conducted on IJRR dataset.
**Spiking Neurons** In the first ablation study, we explore the effect of different spiking neurons on the reconstruction performance. We test three types of spiking neurons

in EVSNN: IF neuron, LIF neuron, and PLIF (parametric LIF) neuron. Since IF neuron simply integrates inputs and lacks the decay mechanism, its performance is worse than the other two neurons. For the other two spiking neurons, a membrane time constant controls the decay. As shown in the upper part of Table 2, LIF neurons perform slightly better than PLIF neurons.

**Membrane Potential Neurons** Based on the EVSNN-LIF architecture, we further analyze the effects of MP neurons. Fig. 7 and 8 show that MP neurons can improve the reconstruction quality of SNN. We test four types of MP neurons: MP_IF, MP_LIF, MP_PLIF, and AMP_LIF. These neurons are non-spiking neurons which output membrane potential instead of spikes. PLIF neurons can learn a fixed membrane time constant based on the training dataset. However, the fixed membrane time constant learned from synthetic data may be not suitable for complex scenes. As shown in the lower part of Table 2, our AMP_LIF performs best because it can adaptively adjust the decay rate by the input spikes.

**Spike Skip Connection** Spike skip connection gathers the spike outputs of the encoder and decoder. An effective connection operation can greatly improve the performance of SNN. Based on the EVSNN-LIF and PA-EVSNN-AMP_LIF architectures, we study four types of spike connections. As shown in Table 3, ADD performs best because
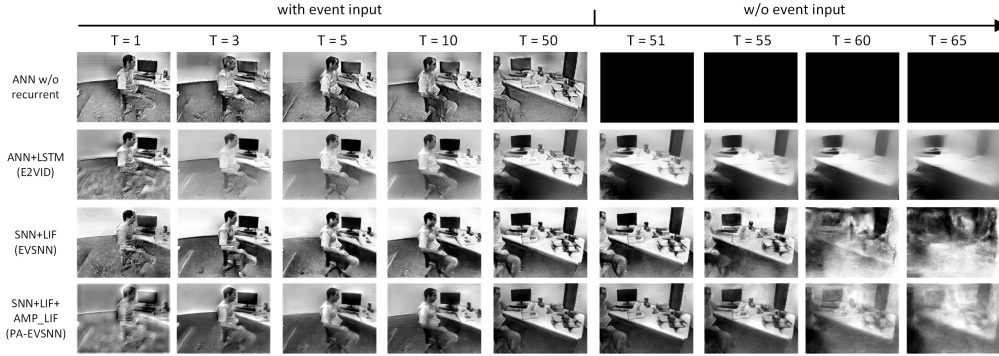
Figure 7. **Comparison on the different temporal components of SNN and ANN.** This figure shows image reconstructions of different ANN and SNN variants at initialization and end phases. At initialization phase (T = 1-50), the states of temporal components (e.g., LSTM and spiking neuron) are initialized at zero, all models are fed continuous event tensors to test the reconstruction of each moment. In subsequent iterations after T = 50, the models are fed empty event tensors to test the effective size of temporal receptive field.
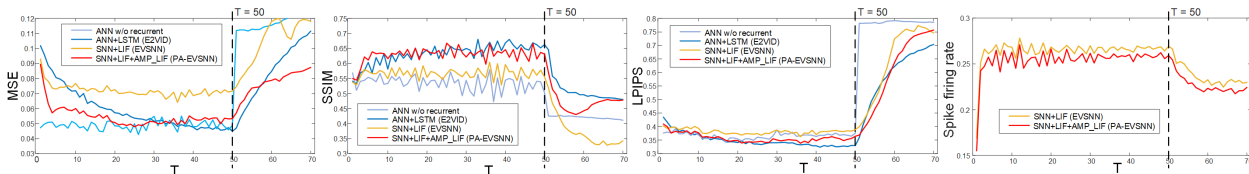


Figure 8. **Quantitative analysis of the temporal components.** This figure shows the MSE, SSIM, LPIPS, and spike firing rate at each iteration. The experiment setting is same as Fig.7, four ANN and SNN variants (ANN w/o recurrent, ANN + LSTM, SNN + LIF, and SNN + LIF + AMP_LIF) are tested. The results show that spiking neuron can improve the ability of temporal information extracting.

Table 2. Ablation studies of spiking neurons and MP neurons.

| Model | MSE↓ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| EVSNN-IF | 0.108 | 0.341 | 0.608 |
| EVSNN-PLIF | 0.063 | 0.569 | 0.367 |
| EVSNN-LIF | **0.061** | **0.570** | **0.362** |
| PA-EVSNN-MP_IF | 0.121 | 0.362 | 0.741 |
| PA-EVSNN-MP_LIF | 0.056 | 0.597 | 0.388 |
| PA-EVSNN-MP_PLIF | 0.053 | 0.599 | 0.378 |
| PA-EVSNN-AMP_LIF | **0.042** | **0.632** | **0.376** |

Table 3. Ablation studies of different spike skip connections.

| Model | MSE↓ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| [1]EVSNN-ADD | **0.049** | **0.586** | **0.350** |
| [2]EVSNN-OR | 0.063 | 0.534 | 0.395 |
| [3]EVSNN-IAND | <u>0.051</u> | 0.557 | <u>0.357</u> |
| [4]EVSNN-CONCAT | 0.061 | <u>0.570</u> | 0.362 |
| [1]PA-EVSNN-ADD | **0.041** | **0.635** | <u>0.388</u> |
| [2]PA-EVSNN-OR | 0.064 | 0.591 | 0.436 |
| [3]PA-EVSNN-IAND | 0.055 | 0.602 | 0.410 |
| [4]PA-EVSNN-CONCAT | <u>0.046</u> | <u>0.626</u> | **0.367** |

Defining the connect operation as $g(A_l, B_l)$, where $A_l \in \{0,1\}$ and $B_l \in \{0,1\}$ denote the spike output of $l$-th encoder and decoder, respectively. The different connections can be implemented as
[1] ADD: $g_{ADD}(A_l, B_l) = A_l + B_l$
[2] OR: $g_{OR}(A_l, B_l) = \max(A_l, B_l)$
[3] IAND: $g_{IAND}(A_l, B_l) = (1 - A_l) \cdot B_l$
[4] CONCAT: $g_{CON.}(A_l, B_l) = [A_l, B_l]$

it retains more information by adding spikes of the encoding and decoding layers. However, it brings non-spike output, e.g., the addition of two spikes will output 2, which breaks the fully SNN architecture and adds additional power consumption. In contrast, OR, IAND, and CONCAT all output spikes. As shown in Table 3, CONCAT performs best while the IAND also performs well. Although the number of parameters of CONCAT is twice that of IAND, we choose CONCAT in our architecture for better performance.

**Number of Encoders and Residual blocks** Finally, we search the number of encoders and residual blocks. The results are shown in Table 4, e.g., EVSNN-e3-res1 means EVSNN with three encoders and one residual block. Considering both the performance and complexity, we choose EVSNN-e3-res1 and PA-EVSNN-e3-res1 as our models.

## 4.5. Energy Consumption and Limitation Analysis

**Energy Comparison of SNN and ANN** Typically, the number of synaptic operations is used as a metric for benchmarking the computational energy of neuromorphic hard-

ware [27]. In ANN, each operation computes a dot-product involving one floating-point (FP) multiplication and one FP addition as a multiply-accumulate (MAC) computation. In contrast, the computations in SNN implemented on neuromorphic hardware are event-driven. Therefore, in the absence of spikes, there are no computations and no active energy is consumed [8]. Thus, in SNN, each operation is only one FP addition due to binary spikes. The low consumption of SNN synapse operation combined with activation sparsity provides large improvements in computational efficiency.

To compare the consumption between SNN and ANN architectures, the evaluation should be conducted on the same structure [39]. Thus, we compute the energy consumption between our SNN models and their ANN versions (e.g., re-

Table 4. Ablation studies of different network architectures.

| Model | MSE↓ | SSIM↑ | LPIPS↓ |
|---|---|---|---|
| EVSNN-e2-res1 | **0.060** | 0.569 | 0.364 |
| EVSNN-e3-res1 | 0.061 | 0.570 | 0.362 |
| EVSNN-e4-res1 | 0.061 | **0.576** | 0.379 |
| EVSNN-e3-res0 | 0.061 | 0.569 | **0.360** |
| EVSNN-e3-res2 | 0.067 | 0.570 | 0.371 |
| PA-EVSNN-e2-res1 | 0.050 | **0.628** | 0.376 |
| PA-EVSNN-e3-res1 | **0.046** | 0.626 | 0.367 |
| PA-EVSNN-e4-res1 | 0.048 | 0.618 | 0.379 |
| PA-EVSNN-e3-res0 | 0.058 | 0.599 | 0.413 |
| PA-EVSNN-e3-res2 | **0.046** | 0.615 | **0.361** |

Table 5. Spike firing rate of EVSNN and PA-EVSNN.

| Layer | Spiking Neuron Num. | Neuron Type | Spike Firing Rate | |
|---|---|---|---|---|
| | | | EVSNN | PA-EVSNN |
| Head | 32×H×W | LIF | 0.2479 | 0.2444 |
| Down1 | 64×H×W | LIF | 0.2459 | 0.2308 |
| Down2 | 128×H×W | LIF | 0.1352 | 0.1339 |
| Down3 | 256×H×W | LIF | 0.1174 | 0.1183 |
| Res1-1 | 256×H×W | LIF | 0.1241 | 0.1098 |
| Res1-2 | 256×H×W | LIF | 0.1308 | 0.1200 |
| Up1 | 128×H×W | LIF | 0.1905 | 0.1983 |
| Up2 | 64×H×W | LIF | 0.3338 | 0.3573 |
| Up3 | 32×H×W | LIF | 0.3580 | 0.3081 |
| Overall spike firing rate | | | 0.2642 | 0.2511 |

Table 6. Comparison of compute energy between ANN and SNN.

| | EVSNN | PA-EVSNN |
|---|---|---|
| [1](a) Normalized $\#OP_{ANN}$ | 1 | 1 |
| [2](b) Normalized $\#OP_{SNN}$ | 0.264 | 0.251 |
| [3](c) Normalized $\#OP_{MP\_layer}$ | 0 | 0.084 |
| [4](d) ANN/SNN Energy | 19.36 | 7.75 |

[1] $\#OP_{ANN}$ is the total number of ANN operations if all spiking neurons are replaced with an ANN activation function (e.g., ReLU).
[2] $\#OP_{SNN} = SpikeRate \times \#OP_{ANN}$.
[3] MP layer contains FP multiplications and additions, thus its consumption is considered same as ANN.
[4] Each operation in ANN (SNN) consumes 4.6pJ (0.9pJ). ANN/SNN Energy can be calculated by $\frac{(a) \times 4.6}{(c) \times 4.6 + (1-(c)) \times (b) \times 0.9}$.

Table 7. Energy comparison of E2VID and our models.

| | E2VID-LSTM | E2VID-GRU | EVSNN | PA-EVSNN |
|---|---|---|---|---|
| Para. Num. | 10.71M | 9.16 M | 4.41M | 4.62M |
| Spike Rate | - | - | 0.264 | 0.251 |
| $\#OP_{ANN}$ | 20.07G | 17.63G | 0 | 1.49G |
| $\#OP_{SNN}$ | 0 | 0 | 16.12G | 16.35 G |
| Energy ($10^{-3}$J) | 92.32 | 81.10 | 3.83 | 10.55 |
| [1]Normalized Energy | 1 | 0.8783 | 0.0414 | 0.1142 |

[1] The energy consumptions of EVSNN and PA-EVSNN are 24.15× and 8.76× more computationally efficient than E2VID, respectively.

place the spiking neurons with ReLU). In most technologies, the addition operation is much cheaper than the multiplication operation. We compute the energy cost/operation for ANNs and SNNs in 45nm CMOS technology. The energy cost for 32-bit ANN MAC operation is 5.1 more than SNN addition operation (4.6pJ vs. 0.9pJ) [17].

The number of synaptic operations in SNN can be calculated by multiplying $\#OP_{ANN}$[1] by the spike firing rate. For example, a spike rate of 1 (every neuron fired) implies that the number of operations for ANN and SNN are the same (though operations are MAC in ANN while addition in SNNs). Lower spike rates denote more sparsity in spike events and higher energy-efficiency. As shown in Table 5, we count the average spike firing rate of EVSNN and PA-EVSNN on IJRR dataset. The comparison results are shown in Table 6. Notice that our models do not require multiple time-steps simulation, which brings a great advantage in energy consumption. For EVSNN, all operations are SNN operations, the average spike firing rate is 26.4%, its energy consumption is 19.36 × lower than ANN. Since there are 8.4% MAC operations in PA-EVSNN, the average spike firing rate of spiking neurons is 25.1%, it costs 7.75 × lower energy consumption compared to its fully ANN version.

**Energy Comparison with E2VID** Here we compare the energy consumption of our models with E2VID. Table 7 reports the energy comparsion[2] with 180×240 input size. Each ANN operation consumes 4.6 pJ, brings

20.07G×4.6pJ = $9.232 \times 10^{-2}$J energy consumption. Compared to LSTM, GRU is a recurrent module with fewer parameters. Our EVSNN has 16.12G SNN operation with 26.4% spike firing rate, which costs 16.12G×26.4%×0.9pJ = $3.83 \times 10^{-3}$J. For PA-EVSNN, we consider the 1.49G operations which come from MP neurons as ANN operations. Thus, the overall energy cost of PA-EVSNN is 1.49G×4.6pJ+16.35G×25.1%×0.9pJ = $1.055 \times 10^{-2}$J. In summary, the energy consumptions of EVSNN and PA-EVSNN are 24.15 × and 8.76 × more computationally efficient than E2VID, respectively.

**Limitation** To make SNN training faster and more stable, we add the batch normalization (BN) after convolution (CONV) layer. Notice that BN can be folded in a CONV layer after training [45]. However, BN is not unbiased. If there is no spike input, BN will also produces non-zero values, which may activate spiking neurons. This will increase the spike rate of SNN, thereby increasing energy consumption (see Fig. 8, the spike rate > 0 when no event input (T > 50)). Reducing the spike rate may be a future direction.

## 5. Conclusion

In this paper, we have presented EVSNN and PA-EVSNN models, the event-based video reconstruction models based on SNN architecture. We show that the spiking neurons have the capability of extracting temporal information, and SNN can achieve large scale regression tasks such as event-based video reconstruction. Compared to E2VID, the proposed EVSNN and PA-EVSNN have 24.15× and 8.76× more computationally efficient improvement, which shows great potential of SNN for low consumption applications. We believe that the development of energy-efficient

---

[1] In ANN-based model, the number of ANN operations (MAC) is defined by $\#OP_{ANN} = \sum k_w \times k_h \times c_{in} \times h_{out} \times w_{out} \times c_{out}$, where $k_w$ and $k_h$ are kernel size, $c_{in}$ and $c_{out}$ are the number of input and channels, $h_{out}$ and $w_{out}$ are output feature map size, and $f_{in}$ and $f_{out}$ is the number of input (output) features.

[2] Energy = $\#OP_{ANN} \times 4.6$pJ $+ \#OP_{SNN} \times 0.9$pJ $\times SpikeRate$. Notice that $\#OP_{SNN}$ must be operated on binary spikes (i.e., 0 or 1).

SNN models for large-scale regression tasks is promising.

# References

[1] Arnon Amir, Brian Taba, David Berg, Timothy Melano, Jeffrey McKinstry, Carmelo Di Nolfo, Tapan Nayak, Alexander Andreopoulos, Guillaume Garreau, Marcela Mendoza, et al. A low power, fully event-based gesture recognition system. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7243–7252, 2017. 1

[2] Patrick Lichtsteiner an Christoph Posch and Tobi Delbruck. A $128 \times 128$ 120 *db* 15 $\mu s$ latency asynchronous temporal contrast vision sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, 2008. 1

[3] Patrick Bardow, Andrew J Davison, and Stefan Leutenegger. Simultaneous optical flow and intensity estimation from an event camera. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 884–892, 2016. 1, 2

[4] Sander M Bohte, Joost N Kok, and Han La Poutre. Error-backpropagation in temporally encoded networks of spiking neurons. *Neurocomputing*, 48(1-4):17–37, 2002. 2

[5] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A $240 \times 180$ 130 *db* 3 $\mu s$ latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014. 1

[6] Pablo Rodrigo Gantier Cadena, Yeqiang Qian, Chunxiang Wang, and Ming Yang. Spade-e2vid: Spatially-adaptive denormalization for event-based video reconstruction. *IEEE Transactions on Image Processing*, 30:2488–2500, 2021. 5

[7] Matthew Cook, Luca Gugelmann, Florian Jug, Christoph Krautz, and Angelika Steger. Interacting maps for fast visual interpretation. In *The 2011 International Joint Conference on Neural Networks*, pages 770–776. IEEE, 2011. 2

[8] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *Ieee Micro*, 38(1):82–99, 2018. 7

[9] Wei Fang, Yanqi Chen, Jianhao Ding, Ding Chen, Zhaofei Yu, Huihui Zhou, and Yonghong Tian. other contributors. spikingjelly, 2020. 5

[10] Wei Fang, Zhaofei Yu, Yanqi Chen, Timothée Masquelier, Tiejun Huang, and Yonghong Tian. Incorporating learnable membrane time constant to enhance learning of spiking neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2661–2671, 2021. 1, 2, 4

[11] Mathias Gehrig, Sumit Bam Shrestha, Daniel Mouritzen, and Davide Scaramuzza. Event-based angular velocity regression with spiking networks. In *2020 IEEE International Conference on Robotics and Automation*, pages 4195–4202. IEEE, 2020. 1, 2

[12] Wulfram Gerstner, Werner M Kistler, Richard Naud, and Liam Paninski. *Neuronal dynamics: From single neurons to networks and models of cognition*. Cambridge University Press, 2014. 3

[13] Robert Gütig and Haim Sompolinsky. The tempotron: a neuron that learns spike timing–based decisions. *Nature Neuroscience*, 9(3):420–428, 2006. 2

[14] Germain Haessig, Andrew Cassidy, Rodrigo Alvarez, Ryad Benosman, and Garrick Orchard. Spiking optical flow for event-based sensors using ibm's truenorth neurosynaptic system. *IEEE Transactions on Biomedical Circuits and Systems*, 12(4):860–870, 2018. 1

[15] Jesse Hagenaars, Federico Paredes Valles, and Guido De Croon. Self-supervised learning of event-based optical flow with spiking neural networks. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021. 1, 2

[16] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997. 3

[17] 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*, pages 10–14. IEEE, 2014. 8

[18] Yingyezhe Jin, Wenrui Zhang, and Peng Li. Hybrid macro/micro level backpropagation for training deep spiking neural networks. *arXiv preprint arXiv:1805.07866*, 2018. 2

[19] Jacques Kaiser, Hesham Mostafa, and Emre Neftci. Synaptic plasticity dynamics for deep continuous local learning (decolle). *Frontiers in Neuroscience*, 14:424, 2020. 2

[20] Hanme Kim, Stefan Leutenegger, and Andrew J Davison. Real-time 3d reconstruction and 6-dof tracking with an event camera. In *European Conference on Computer Vision*, pages 349–364. Springer, 2016. 2

[21] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5

[22] Wei-Sheng Lai, Jia-Bin Huang, Oliver Wang, Eli Shechtman, Ersin Yumer, and Ming-Hsuan Yang. Learning blind video temporal consistency. In *Proceedings of the European Conference on Computer Vision*, pages 170–185, 2018. 4

[23] Chankyu Lee, Adarsh Kumar Kosta, Alex Zihao Zhu, Kenneth Chaney, Kostas Daniilidis, and Kaushik Roy. Spikeflownet: event-based optical flow estimation with energy-efficient hybrid neural networks. In *Proceedings of the European Conference on Computer Vision*, pages 366–382. Springer, 2020. 2

[24] Chankyu Lee, Syed Shakib Sarwar, Priyadarshini Panda, Gopalakrishnan Srinivasan, and Kaushik Roy. Enabling spike-based backpropagation for training deep neural network architectures. *Frontiers in Neuroscience*, 14:119, 2020. 2

[25] Jun Haeng Lee, Tobi Delbruck, and Michael Pfeiffer. Training deep spiking neural networks using backpropagation. *Frontiers in Neuroscience*, 10:508, 2016. 2

[26] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision*, pages 740–755. Springer, 2014. 5

[27] 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. 7

[28] Ammar Mohemmed, Stefan Schliebs, Satoshi Matsuda, and Nikola Kasabov. Span: Spike pattern association neuron for learning spatio-temporal spike patterns. *International Journal of Neural Systems*, 22(04):1250012, 2012. 2

[29] Elias Mueggler, Henri Rebecq, Guillermo Gallego, Tobi Delbruck, and Davide Scaramuzza. The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and slam. *The International Journal of Robotics Research*, 36(2):142–149, 2017. 5

[30] Gottfried Munda, Christian Reinbacher, and Thomas Pock. Real-time intensity-image reconstruction for event cameras using manifold regularisation. *International Journal of Computer Vision*, 126(12):1381–1393, 2018. 1, 2

[31] Emre O Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal Processing Magazine*, 36(6):51–63, 2019. 5

[32] Garrick Orchard, Cedric Meyer, Ralph Etienne-Cummings, Christoph Posch, Nitish Thakor, and Ryad Benosman. Hfirst: A temporal approach to object recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(10):2028–2040, 2015. 1

[33] Chethan M Parameshwara, Simin Li, Cornelia Fermüller, Nitin J Sanket, Matthew S Evanusa, and Yiannis Aloimonos. Spikems: Deep spiking neural network for motion segmentation. *arXiv preprint arXiv:2105.06562*, 2021. 1, 2

[34] Federico Paredes-Vallés and Guido CHE de Croon. Back to event basics: Self-supervised learning of image reconstruction for event cameras via photometric constancy. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3446–3455, 2021. 1, 2

[35] Federico Paredes-Vallés, Kirk YW Scheper, and Guido CHE de Croon. Unsupervised learning of a hierarchical spiking neural network for optical flow estimation: From events to global motion perception. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(8):2051–2064, 2019. 1, 2

[36] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in Neural Information Processing Systems*, 32:8026–8037, 2019. 5

[37] Filip Ponulak and Andrzej Kasiński. Supervised learning in spiking neural networks with resume: sequence learning, classification, and spike shifting. *Neural Computation*, 22(2):467–510, 2010. 2

[38] Ulysse Rançon, Javier Cuadrado-Anibarro, Benoit R Cottereau, and Timothée Masquelier. Stereospike: Depth learning with a spiking neural network. *arXiv preprint arXiv:2109.13751*, 2021. 2

[39] Nitin Rathi and Kaushik Roy. Diet-snn: A low-latency spiking neural network with direct input encoding and leakage and threshold optimization. *IEEE Transactions on Neural Networks and Learning Systems*, 2021. 7

[40] Henri Rebecq, Daniel Gehrig, and Davide Scaramuzza. Esim: an open event camera simulator. In *Conference on Robot Learning*, pages 969–982. PMLR, 2018. 2

[41] Henri Rebecq, Timo Horstschäfer, Guillermo Gallego, and Davide Scaramuzza. Evo: A geometric approach to event-based 6-dof parallel tracking and mapping in real time. *IEEE Robotics and Automation Letters*, 2(2):593–600, 2016. 2

[42] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. Events-to-video: Bringing modern computer vision to event cameras. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3857–3866, 2019. 1, 2

[43] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. High speed and high dynamic range video with an event camera. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2019. 1, 2, 4, 5, 6

[44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 234–241. Springer, 2015. 2, 3

[45] 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. 8

[46] Cedric Scheerlinck, Nick Barnes, and Robert Mahony. Continuous-time intensity estimation using event cameras. In *IEEE Asian Conference on Computer Vision (ACCV)*, pages 308–324. Springer, 2018. 1, 2

[47] Cedric Scheerlinck, Nick Barnes, and Robert Mahony. Asynchronous spatial image convolutions for event cameras. *IEEE Robotics and Automation Letters*, 4(2):816–822, 2019. 1, 2

[48] Cedric Scheerlinck, Henri Rebecq, Daniel Gehrig, Nick Barnes, Robert Mahony, and Davide Scaramuzza. Fast image reconstruction with an event camera. In *IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 156–163, 2020. 1, 2, 5

[49] Sumit Bam Shrestha and Garrick Orchard. Slayer: Spike layer error reassignment in time. *arXiv preprint arXiv:1810.08646*, 2018. 2

[50] Timo Stoffregen, Cedric Scheerlinck, Davide Scaramuzza, Tom Drummond, Nick Barnes, Lindsay Kleeman, and Robert Mahony. Reducing the sim-to-real gap for event cameras. In *European Conference on Computer Vision (ECCV)*, pages 534–549. Springer, 2020. 1, 2, 5, 6

[51] Beck Strohmer, Rasmus Karnøe Stagsted, Poramate Manoonpong, and Leon Bonde Larsen. Integrating non-spiking interneurons in spiking neural networks. *Frontiers in Neuroscience*, 15:184, 2021. 3

[52] Lin Wang, Yo-Sung Ho, Kuk-Jin Yoon, et al. Event-based high dynamic range image and very high frame rate

video generation using conditional generative adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10081–10090, 2019. 1, 2

[53] Xiao Wang, Jianing Li, Lin Zhu, Zhipeng Zhang, Zhe Chen, Xin Li, Yaowei Wang, Yonghong Tian, and Feng Wu. Visevent: Reliable object tracking via collaboration of frame and event flows, 2021. 1

[54] Zhou Wang, Alan C Bovik, Hamid R Sheikh, Eero P Simoncelli, et al. Image quality assessment: from error visibility to structural similarity. *IEEE Transactions on Image Processing*, 13(4):600–612, 2004. 5

[55] Wenming Weng, Yueyi Zhang, and Zhiwei Xiong. Event-based video reconstruction using transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2563–2572, 2021. 1, 2

[56] Paul J Werbos. Backpropagation through time: what it does and how to do it. *Proceedings of the IEEE*, 78(10):1550–1560, 1990. 5

[57] Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in Neuroscience*, 12:331, 2018. 2

[58] Zhenzhi Wu, Hehui Zhang, Yihan Lin, Guoqi Li, Meng Wang, and Ye Tang. Liaf-net: Leaky integrate and analog fire network for lightweight and efficient spatiotemporal information processing. *IEEE Transactions on Neural Networks and Learning Systems*, 2021. 3

[59] Yannan Xing, Gaetano Di Caterina, and John Soraghan. A new spiking convolutional recurrent neural network (scrnn) with applications to event-based hand gesture recognition. *Frontiers in Neuroscience*, 14:1143, 2020. 1, 2

[60] Richard Zhang, Phillip Isola, Alexei A Efros, Eli Shechtman, and Oliver Wang. The unreasonable effectiveness of deep features as a perceptual metric. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 586–595, 2018. 4, 5

[61] Xiang Zhang, Wei Liao, Lei Yu, Wen Yang, and Gui-Song Xia. Event-based synthetic aperture imaging with a hybrid network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14235–14244, 2021. 2

[62] Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li. Going deeper with directly-trained larger spiking neural networks. *arXiv preprint arXiv:2011.05280*, 2020. 1, 2

[63] Yajing Zheng, Lingxiao Zheng, Zhaofei Yu, Boxin Shi, Yonghong Tian, and Tiejun Huang. High-speed image reconstruction through short-term plasticity for spiking cameras. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6358–6367, 2021. 2

[64] 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. 5

[65] 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/CVF Conference on Computer Vision and Pattern Recognition*, pages 989–997, 2019. 2

[66] Lin Zhu, Siwei Dong, Tiejun Huang, and Yonghong Tian. A retina-inspired sampling method for visual texture reconstruction. In *2019 IEEE International Conference on Multimedia and Expo (ICME)*, pages 1432–1437. IEEE, 2019. 2

[67] Lin Zhu, Siwei Dong, Jianing Li, Tiejun Huang, and Yonghong Tian. Retina-like visual image reconstruction via spiking neural model. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1438–1446, 2020. 2