# SpikeMS: Deep Spiking Neural Network for Motion Segmentation

Chethan M. Parameshwara, Simin Li, Cornelia Fermüller, Nitin J. Sanket, Matthew S. Evanusa, Yiannis Aloimonos
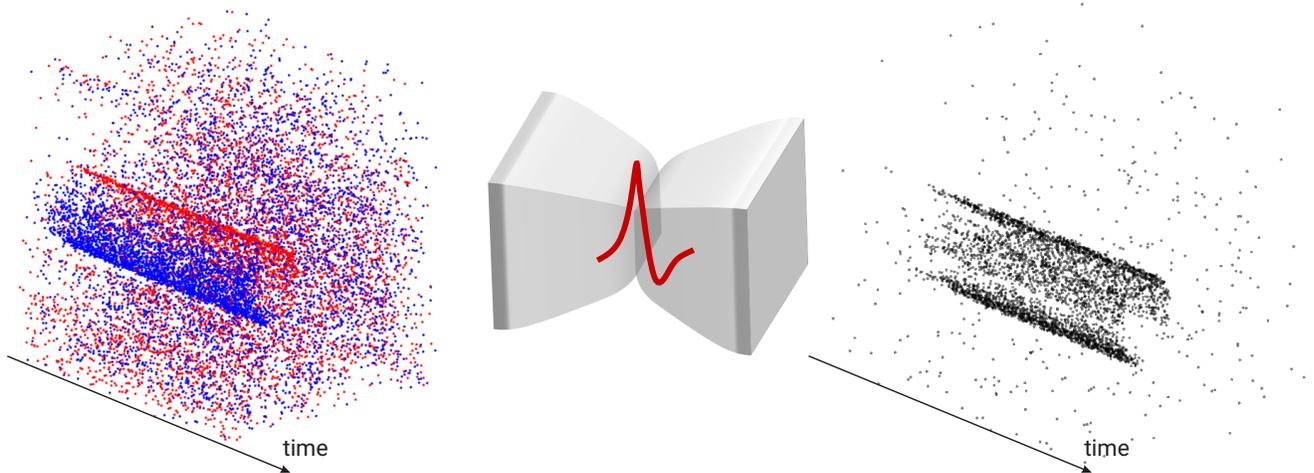
Fig. 1.   Event-based motion segmentation pipeline using a deep spiking neural network. Left to right: Event stream input represented as red (brightness increase) and blue (brightness decrease), representation of the proposed encoder-decoder spiking neural network called *SpikeMS* and the output predicted spike containing only the region of moving object(s). *All the images in this paper are best viewed in color on a computer screen at a zoom of 200%.*

*Abstract*— Spiking Neural Networks (SNN) are the so-called third generation of neural networks which attempt to more closely match the functioning of the biological brain. They inherently encode temporal data, allowing for training with less energy usage and can be extremely energy efficient when coded on neuromorphic hardware. In addition, they are well suited for tasks involving event-based sensors, which match the event-based nature of the SNN. However, SNNs have not been as effectively applied to real-world, large-scale tasks as standard Artificial Neural Networks (ANNs) due to the algorithmic and training complexity. To exacerbate the situation further, the input representation is unconventional and requires careful analysis and deep understanding. In this paper, we propose *SpikeMS*, the first deep encoder-decoder SNN architecture for the real-world large-scale problem of motion segmentation using the event-based DVS camera as input. To accomplish this, we introduce a novel spatio-temporal loss formulation that includes both spike counts and classification labels in conjunction with the use of new techniques for SNN backpropagation. In addition, we show that *SpikeMS* is capable of *incremental predictions*, or predictions from smaller amounts of test data than it is trained on. This is invaluable for providing outputs even with partial input data for low-latency applications and those requiring fast predictions. We evaluated *SpikeMS* on challenging synthetic and real-world sequences from EV-IMO, EED and MOD datasets and achieving results on a par with a comparable ANN method, but using potentially 50 times less power.

*Chethan M. Parameshwara and Simin Li contributed equally to this work. (Corresponding author: Chethan M. Parameshwara)*

All authors are associated with the Perception and Robotics Group, University of Maryland, College Park. Emails: {cmparam9, sli12348, fer, nitin, evanusa, yiannis} @umiacs.umd.edu

## SUPPLEMENTARY MATERIAL

The accompanying video and supplementary material are available at prg.cs.umd.edu/SpikeMS.

## I. INTRODUCTION

The animal brain is remarkable at perceiving motion in complex scenarios with high speed and extreme energy efficiency. Inspired by the animal brain, an alternative version of Artificial Neural Networks (ANNs) called Spiking Neural Networks (SNNs) aim to replicate the *dynamical system* aspects of living neurons. In contrast to standard ANNs which are essentially networks of complex functions, SNNs are comprised of networks of neurons modeled as differential equations, and inherently encode temporal data and offer low power and highly parallelizable computations. Furthermore, they possess the capability to deliver predictions whose confidence scale with the availability of input data [1], [2]. These low-power and low-latency properties are of great use to real-world robotics applications such as self-driving cars or drones, which demand fast responses during navigation in challenging scenarios [3].

Until recently, SNNs have been restricted to simple tasks and small datasets due to instability in learning regimes [4]. Recent development in new spike learning mechanisms [5], [6] has made it possible to design SNNs for real-world robotics applications. This coupled with neuromorphic processors such as Intel's ® Loihi [7] and IBM's TrueNorth [8]) along with neuromorphic sensors such

as DVS [9] and ATIS [10]) have made it possible for producing real-world prototypes, drastically enhancing the appeal of such technologies.

In this work, we propose a deep SNN architecture called *SpikeMS* for the problem of motion segmentation using a monocular event camera. We consider the data from event sensors as they are well-suited for motion segmentation (due to the disparity in event density at object boundaries) and are a natural fit for SNNs (due to their temporal nature). We will now formally define the problem statement and our main contributions.

### A. Problem Formulation and Contributions

We address the following question: *How do you learn to segment the scene into background and foreground (moving objects) using a Spiking Neural Network from the data of a moving monocular event camera?*

Our spiking neural network, *SpikeMS*, takes the event stream as input and outputs predictions of each event's class association as either foreground (moving object) or background (moving due to camera motion).

The model learns to distinguish between the spatio-temporal patterns of moving objects and the background. To the best of our knowledge, this is the first end-to-end deep encoder-decoder SNN. In particular, we evaluate our network on the task of motion segmentation using event input.

The main contributions of the paper are given below:

- A novel end-to-end deep encoder-decoder Spiking Neural Network (SNN) framework for motion segmentation from event-based cameras.
- Demonstration of "early" evaluation of the network (at low latency), which we call *Incremental Predictions*, for imprecise but fast detection of moving objects for variable-sized integration windows.

## II. RELATED WORK

### A. Spiking Neural Network Weight Learning Rules

While the concept of a spiking neuron has been around for a few decades [11], their progress has been bounded by the difficulty in training due to the ubiquitous vanishing gradient problem for deep neural networks. In SNNs, the neurons output pulses that are non-differentiable, rendering attempts at directly applying the backpropagation algorithm non-trivial. Early attempts at training SNNs revolved around more biologically plausible Hebbian-style mechanisms [12] that only involve *local* updates, such as *Spike Time Dependent Plasticity* (STDP) [13], avoiding gradient issues. Work in this field continues to this day, with results [14]–[16] demonstrating utilities of STDP in training deep SNNs. Early attempts at incorporating backpropagtion into SNNs involved first training a traditional ANN, and then transferring the learned weights to an SNN [1]. Recent methods, which we build off of here, allow the SNN to be directly trained through backpropagation by finding a surrogate, continuous value function that roughly correlates for the spike activity [5], [6], [17].

### B. SNNs for Visual Tasks and Event Data

There has been a renewed interest in using SNNs to process data directly from event-based visual sensors, such as the DVS since the sensor produces spike-like activity that fits well with SNN neurons. Applications of SNNs in this domain include classification problems [18] such as digit recognition [19], object recognition [20] gesture recognition [21], and optical flow [22], [23]. Recent development of neuromorphic processors such as the Intel Loihi [7] has lead to the deployment of SNNs on hardware [22], [24], [25].

Closest related to our work, in [4] recently a neural architecture of multiple layers has been designed. A six-layer neural network (five convolutional and one pooling layer) is used to learn with supervision to regress the three parameters of camera rigid rotation. In Lee *et al.* [26] a deep hybrid encoder decoder architecture was designed for self-supervised optic flow estimation. The encoding layers are SNN with the backpropagation learning employing the approximation of [17], and the residual and decoding layers are ANN with the self-supervised loss computed from the images of a combined DVS and image sensor (DAVIS). Our network is the first architecture for the problem of motion segmentation with event data.

Event-based cameras have been recognized as a promising sensor for the problem of segmentation and detection of independently moving objects, as the event stream carries essential information about the movement of object boundaries [27]. Classical approaches [28]–[30] treat motion segmentation as a geometric problem and model it as an artifact of motion compensation of events. In ANN approaches, the input representation is formed by binning the events within a time-interval and convert to an image-like frame based structure [3], [31] the so called "event-frames". Our approach is similar to [32], but rather than creating event-frames, a sampled version of the event stream is fed directly into the network, taking advantage of the SNN's temporal nature in conjunction with the temporal nature of the event stream.

## III. SPIKEMS ARCHITECTURE

### A. Event Camera Input

A traditional camera records frames at a fixed frame rate by integrating the number of photons for the chosen shutter time for all pixels *synchronously* (in a global shutter camera). In contrast, an event camera only records the polarity of logarithmic brightness changes *asynchronously* at each pixel, resulting in asynchronous packets of information containing the pixel location and the time of the change known as an *event*. If the brightness at time $t$ of a pixel at location $\mathbf{x}$ is given by $I_{t,\mathbf{x}}$ an event is triggered when

$$\| \log\left(I_{t+\delta t,\mathbf{x}}\right) - \log\left(I_{t,\mathbf{x}}\right) \|_1 \geq \tau \qquad (1)$$

where $\delta t$ is a small time increment and $\tau$ is a trigger threshold. Each event outputs the following data: $\mathbf{e} = \{\mathbf{x}, t, p\}$, where $p = \pm 1$ denotes the sign of the brightness change. We will denote the event stream in a spatio-temporal
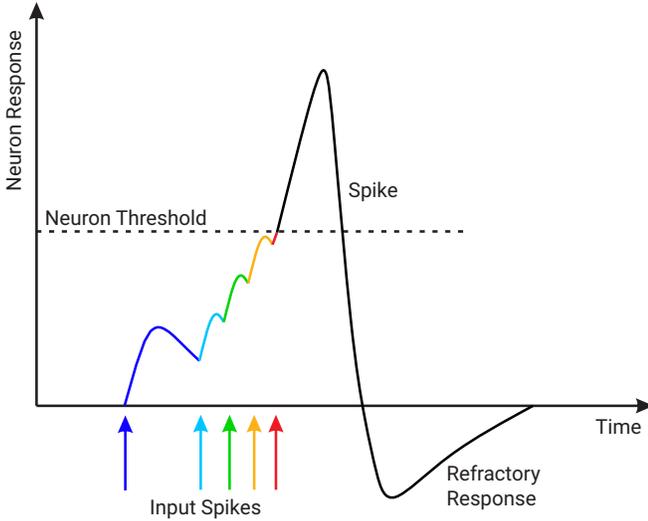
Fig. 2. Depiction of the dynamical activity of a spiking neuron. The neuron receives input coming either from the data or lower layers (shown here as colored arrows), which generate bumps in the membrane voltage; we refer to this voltage in the paper as $u(t)$. If the voltage $u(t)$ exceeds a threshold $\vartheta$, shown here as the dotted line, the neuron outputs a spike, and then enters a refractory phase where it is less likely to fire another spike for a short time. Computationally, this spiking after passing a threshold amounts to feeding $u(t)$ through the spike function $f_s$. The effect that incoming pulses have on the voltage, and the extent of the refractory response, is governed in the Spike Response Model (SRM) [33] via the $\varepsilon$ and $\nu$ kernels respectively (See Section III-B for more detail).

window as $\mathcal{E}(t, t + \delta t) = \{e_i\}_{i=1}^N$ ($N$ is the number of events). We refer to it as event slice/stream/cloud/volume or spike train interchangeably.

### B. Spiking Neuron Model

Spiking Neurons, unlike traditional *rate-encoding* neurons (commonly used neurons in standard ANNs), implicitly encode time in their formulation. They are modeled loosely after neurons in the brain, following the pioneering work by Hodgkin-Huxley [11] which laid the groundwork for differential equation modeling of neuronal activity. We utilize the Spike Response Model (SRM) which similar to all spiking neuron models, sums up incoming voltage from pre-synaptic neurons, but contains two filters: a filter that accounts for the neuron's *self-refractory* response denoted as $\nu$, and a *spike response kernel* that accounts for the integration of incoming pre-synaptic pulses denoted as $\varepsilon$. For a given neuron $i$ at timestep $t$, the update of the neuron's synaptic potential dynamics takes the form of:

$$u_i(t) = \left( \sum_j w_j (\varepsilon * s_j) \right) + (\nu * s)$$
$$= \mathbf{w}^\top \mathbf{a} + (\nu * s) \tag{2}$$

for all incoming weight connections from pre-synaptic neurons $1, ..., j$, where $a(t) = (\varepsilon * s)(t)$, $s_i(t)$ is an input spike train in a neuron and $*$ denotes the convolution operator. An output spike is generated whenever $u(t)$ reaches the spiking threshold $\vartheta$ (the dotted line in Fig. 2)

The motivation for using SRM neuron types is that it inexpensively models the refractory behavior of neurons without having to run multiple differential equation solvers, as seen in other models.

Specific choices of $\varepsilon$ and $\nu$ reduce the SRM equations to a LIF neuron [34]. Here, we use the formulation from [4]:

$$\varepsilon(t) = \frac{t}{\tau_s} e^{1 - \frac{t}{\tau_s}} \mathcal{H}(t) \tag{3}$$

$$\nu(t) = -2\vartheta e^{1 - \frac{t}{\tau_r}} \mathcal{H}(t) \tag{4}$$

where $\mathcal{H}$ is the Heaviside function, and $\tau_s$ and $\tau_r$ are the spike response and refractory time constants.

The activity of the neurons is then propagated forward through the layers of the network, in the same manner as an ANN. The feed-forward weight matrix $\mathbf{W}^{(l)} = [\mathbf{w}_1, ..., \mathbf{w}_{N_{l+1}}]$ for a given layer $l$ with $N_l$ neurons is applied to the activity resulting from the spike response kernel, added to the refractory activity and then thresholded. Thus, for all layers $l$ in the network, the activity is forward-propagated as:

$$a^{(l)}(t) = (\varepsilon_d * s^{(l)})(t) \tag{5}$$

$$u^{(l+1)}(t) = \mathbf{W}^{(l)} a^{(l)}(t) + (\nu * s^{(l+1)}(t)) \tag{6}$$

$$s^{(l+1)}(t) = f_s(u^{(l+1)}(t)) \tag{7}$$

where $f_s$ is the thresholding function, $\mathbf{W}^{(l)}$ is the forward weight matrix for layer $l$, and $\varepsilon_d$ is the spike response kernel with delay, as in [6]. The input to the network, $s^{(0)}$, is the event data over the learning window.

### C. Network Architecture

*SpikeMS* utilizes an end-to-end deep Spiking Neural Network, in contrast to many recent models [26] that use a hybrid combination of spiking and rate-encoding layers. To the best of our knowledge, *SpikeMS* is the first end-to-end spike trained deep encoder-decoder network for large scale tasks such as motion segmentation.

*SpikeMS* consists of a traditional hourglass-shaped layer structure of an autoencoder, with larger layers progressively encoded to smaller representations, which are then decoded back to the original size. We use three encoder layers followed by three decoder layers. The first three convolutional layers perform spatial downsampling with a stride of 2 and kernel size of 3x3. The output of the first encoder layer contains 16 channels, and each encoder layer after doubles the number of channels. The last three decoder layers perform spatial upsampling using transposed convolution, with a stride of 2 and kernel size of 3x3. Decoder layers 4 and 5 each halve the number of channels. The last layer (6) outputs the predicted spikes of the moving object(s) using 2 channels, representing positive and negative event polarities.

### D. Spatio-Temporal Loss

We propose a novel loss formulation which takes full advantage of the spatio-temporal nature of the event data. Our loss function consists of two parts: a binary cross entropy loss $\mathcal{L}_{\text{bce}}$ and spike loss $\mathcal{L}_{\text{spike}}$.

The binary cross-entropy loss $\mathcal{L}_{\text{bce}}$ is computed by comparing the predicted temporal spike projection to the ground truth temporal spike projection.

$$\mathcal{L}_{\text{bce}} = -\left(\mathbb{1}_f \log\left(\hat{\mathcal{E}}_f\right) + \mathbb{1}_b \log\left(\hat{\mathcal{E}}_b\right)\right) \quad (8)$$

Where, the spike projection $\mathcal{E}$ is obtained as: $\mathcal{E}(\mathbf{x}) = \sum_t \mathcal{E}(\mathbf{x})$. Such a projection converts a spike train into a real-valued output, encoding the frequency of spikes. And the groundtruth foreground and background labels are denoted as $\mathbb{1}_f$ and $\mathbb{1}_b$ respectively.

The spike loss $\mathcal{L}_{\text{spike}}$ is derived from the Van-Rossom distance [35] and measures the distance between two binary spike trains [36]. $\mathcal{L}_{\text{spike}}$ preserves the temporal precision of the event stream. The ground truth spike labels are generated by applying a binary mask to the event cloud input $\mathcal{E}$, i.e., masking all non-moving-object events (background events) as 0, and keeping intact the events that correspond to the moving object. $\mathcal{L}_{\text{spike}}$ is given by

$$\mathcal{L}_{\text{spike}} = \sum_{t=0}^{t+\delta t} \left(\hat{\mathcal{E}}(t, t+\delta t) \circ \mathbb{1}_f - \tilde{\mathcal{E}}(t, t+\delta t)\right)^2 dt \quad (9)$$

where $\circ$ denotes the Hadamard product and $\mathbb{1}_f$ denotes the mask of foreground spikes.

The overall loss $\mathcal{L}_{\text{total}}$ is given by

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{bce}} + \lambda \mathcal{L}_{\text{spike}} \quad (10)$$

where $\lambda$ is a weighting factor. The error is backpropagated through the network using SLAYER [6].

### E. Simulation of SNNs on GPU

SNNs are continuous dynamical systems which can process input event streams asynchronously. However, for our experiments, we simulate the SNN network on a GPU. To achieve this, we need to discretize the event data at fixed time steps. To balance the trade-off between accuracy and resource availability [4], we restrict the simulation time step to one millisecond. We train our SNNs with fixed simulation time window/width/steps $\Delta t_{\text{train}}$ of 10ms for all experiments. However, to test the out-of-domain temporal performance, we test our predictions on simulation time steps $\Delta t_{\text{test}}$ of 1ms to 25ms. To fit the event inputs into fixed time steps, the multiple events with the same polarity and spatial location within a timeframe are represented as a single binary event. This downsampling collapses all events within the window into a single event. The simulated network is trained with the publicly available PyTorch implementation of SLAYER [6].
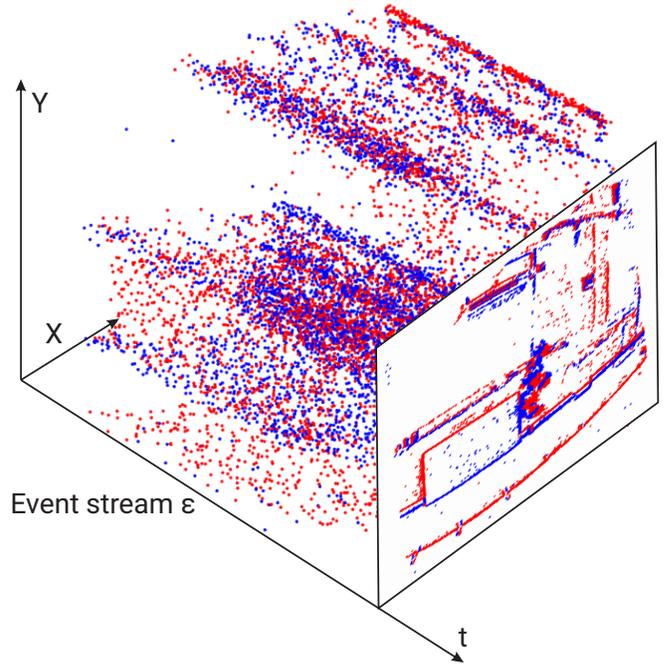


Fig. 3. Representation of event stream $\mathcal{E}$ and its corresponding projection (event frame). Note that, only event streams are fed to *SpikeMS*. The event frame is shown only for clarity purposes.

## IV. Experiments and Results

We evaluate our approach on publicly available synthetic and real datasets. We demonstrate performance of *SpikeMS* both qualitatively and quantitatively by employing the Intersection over Union ($IoU$) and Detection Rate (DR) metrics [30].

### A. Overview of Datasets

We use the publicly available MOD [3] and EV-IMO [31] datasets for training and evaluating the motion segmentation predictions. MOD [3] is a synthetic dataset specifically targeted for learning based motion segmentation approaches. The simulated data contains objects moving in an indoor room-like environment with randomized wall textures, static/dynamic objects and the object/camera trajectories. EV-IMO [31] contains monocular event-based camera data captured in a lab environment with challenging scenarios (multiple objects moving in random trajectories and varying speeds). EV-IMO contains five different sequences (`boxes`, `floor`, `wall`, `table`, and `fast`) which were collected using the DAVIS 346 camera.

### B. Quantitative Results

We compare our method against state-of-the-art ANNs (both 2D and 3D) and the results are given in Table I. In particular, 2D ANNs (EV-IMO [31], EVDodgeNet [3]) are trained with inputs consisting of event-frames computed by accumulating (or projecting) events on a plane. In contrast, 3D ANNs (GConv [32] and PointNet++ [37]) are trained directly on the event cloud $\mathcal{E}$. We evaluate ANN-2D with event frames integrated with a time width $\Delta t$ of 25ms.
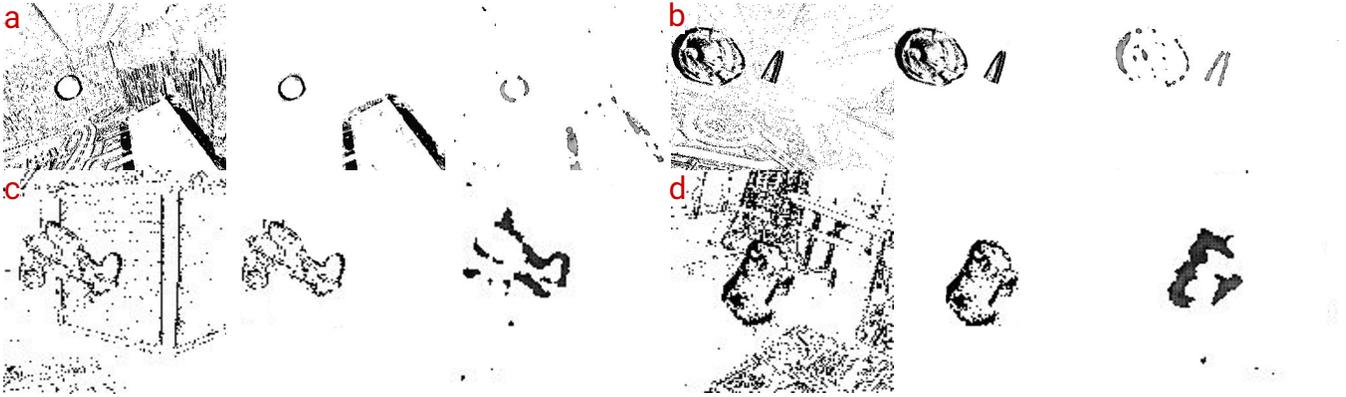
4

Fig. 4. Qualitative Evaluation of our approach on two datasets. Top row (a and b): MOD dataset, Bottom row (c and d): EV-IMO dataset. Each sample includes, (left to right) event stream, groundtruth, and network prediction. Here, we show event projections for clarity purposes but *SpikeMS* predicts spatio-temporal spikes.

TABLE I

QUANTITATIVE EVALUATION USING IoU (%) ↑ METRIC ON EV-IMO AND MOD DATASETS.

| Method | EV-IMO | | | | | | | | | | MOD | |
| | boxes | | floor | | wall | | table | | fast | | | |
| | 100 | 20 | 100 | 20 | 100 | 20 | 100 | 20 | 100 | 20 | 100 | 20 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| EV-IMO[†] [31] | 70±5 | | **59±9** | | 78±5 | | 79±6 | | 67±3 | | - | |
| EVDodgeNet [3] | **67±8** | | 61±6 | | 72±9 | | 70±8 | | 60±10 | | 75±12 | |
| GConv[†] [32] | 81±8 | 60±18 | 79±7 | 55±19 | 83±4 | 80±7 | 57±14 | **51±16** | 74±17 | **39±19** | - | - |
| PointNet++ [37] | 71±22 | 80±15 | 68±18 | 76±10 | 75±19 | 74±20 | 62±28 | 68±23 | 24±10 | 20±6 | 74±13 | **67±15** |
| Ours ($\mathcal{L}_{bce}$) | 57±11 | 59±7 | 56±9 | 46±12 | 62±8 | 62±9 | 51±12 | 45±12 | 42±13 | 36±13 | 62±11 | 63±7 |
| Ours ($\mathcal{L}_{spike}$) | 45±4 | 52±7 | 49±8 | 44±6 | 53±15 | 47±11 | 43±15 | 37±4 | 41±6 | 35±4 | 55±11 | 55±8 |
| Ours ($\mathcal{L}_{bce} + \mathcal{L}_{spike}$) | 61±7 | **65±8** | **60±5** | 53±16 | 65±7 | 63±6 | 52±13 | **50±8** | 45±11 | **38±10** | 68±7 | **65±5** |

[†]Results taken directly from [32]

ANN-3D approaches are evaluated with two time widths $\Delta t$ of 20ms and 100ms similar to [32].

During evaluation, *SpikeMS* is tested at $\Delta t_{test}$ = 100ms and $\Delta t_{test}$ = 20ms (trained at $\Delta t_{train}$=10ms) for a fair comparison with ANNs-3D. Table I provides the mean $IoU$ results on multiple sequences of the EV-IMO and MOD datasets. We observe that the performance of *SpikeMS* is comparable to ANN-2D and ANN-3D approaches in all cases. However, note that the ANN-2D and ANN-3D perform better in the domain they are trained in as compared to *SpikeMS* and we speculate this is because of more stable training procedures in ANNs. This points to a direction of future work for SNNs of proposing better training methodologies.

In Table I, we also compare our results when trained on different loss functions. We observe that the proposed spatio-temporal loss formulation performs better than just using the spike loss or crossentropy loss as it utilizes the information from both spatial and time domains together.

Finally, we also compare *SpikeMS* with classical hand-crafted methods in Table II and we see that, our SNN approach outperforms most hand-crafted methods whilst being deployable directly on neuromorphic hardware. This would lead to huge power savings when deployed on a robot.

### C. Incremental Prediction

We test the network's capability to perform *incremental predictions* evaluating the network at different testing discretized window sizes, with $\Delta t_{test}$ ranging from 1ms to 25 ms, while keeping the training window fixed at $\Delta t_{train}$ of 10ms. This experiment tests for the out-of-domain generalization performance of *SpikeMS*.

All predictions made at $\Delta t < 10$ms can be considered as *incremental predictions* (See Fig. 5), since the testing window is smaller than the training window. This is particularly important for robotics applications since one can filter these incremental predictions to get close to the accuracy of the model with long time predictions but with a lower latency. For example, we can filter predictions (we use a linear Kalman filter [38] for filtering) of 3ms to obtain up to ∼64% of the accuracy of 10ms predictions, but with 70% less latency which might be required for time-critical controllers. We also experiment with values greater than 10ms, where we examine whether longer integration windows yield more accurate results.

Fig. 5 shows the plot of accuracy (IoU) versus the duration of input spikes during simulation, considered for prediction. We observe that the prediction accuracy increases over time with the occurrence of more spikes, but critically, that the SNN is able to output reasonable predictions from less spikes. As shown in SNNs outperform ANN-2D and ANN-3D at early stages with less amount of data. We observe that ANN-3D outperforms ANN-2D since it is trained with temporal augmentation techniques as proposed in [32]. Note that the SNN does not rely on temporal
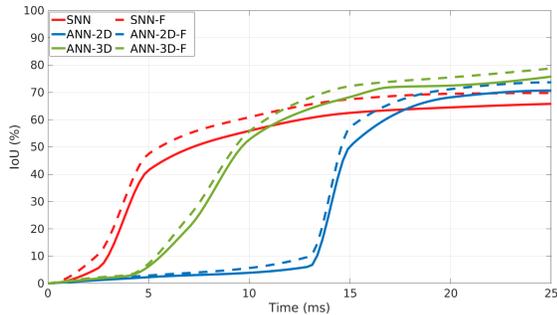
Fig. 5. Incremental Prediction: Segmentation accuracy vs. input spike window length in milliseconds for various simulation time width $\Delta t$. SpikeMS is able to achieve good accuracy significantly faster than ANNs, given smaller input data. The dashed lines represent accuracy improvement after employing a filtering (See Sec. IV-C).

augmentations for incremental predictions rather utilizes dynamic nature of spiking architecture. This demonstrates how well *SpikeMS* generalize outside the temporal domain.

### D. Qualitative Results

Fig. 4 shows qualitative results of our approach on the two datasets. For each example the input, moving object groundtruth, and network prediction are shown. Note that, we show the event projections for clarity purposes but the network input and outputs are the event cloud/spikes. We can observe that the network output predictions are similar to the ground truth for the moving objects in the presence of significant background variation and motion dynamics.

Fig. 6 shows the performance of *SpikeMS* on real-world event streams, again with significant background variations and patterns. These results demonstrate the capability of *SpikeMS* to generalize to different environments without any retraining or fine tuning of the network.

### E. Power Efficiency

We further analyze the benefits of *SpikeMS* compared to a fully ANN architecture with respect to power consumption. It is important to note that the main power consumption benefits occur when the SNN is implemented on a neuromorphic hardware such as the Intel® Loihi [7], where the network *only* consumes power when there is a spike. Hence, the power consumption depends on the mean spike activity of the incoming data and the number of synaptic operations. In contrast, ANNs perform dense matrix operations without exploiting the event sparsity. Thus, in anticipation of deploying *SpikeMS* on the new generation of neurmorphic chips, we demonstrate the power savings by comparing the number of operations by a metric proposed in [26].

Table III provides the average number of synaptic operations in SNNs along with a conservative estimate of the energy benefit when compared to an ANN-2D. We can observe that SNNs require a significantly lower number of synaptic operations and power as compared to ANNs.
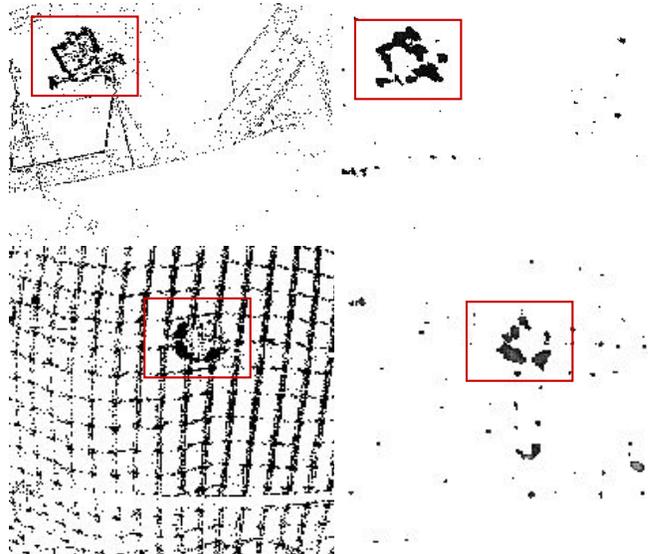


Fig. 6. Results showing motion segmentation generalization without fine-tuning or re-training on real world data. *SpikeMS* is able to segment the moving object from the scene even in the presence of substantial background noise. The objects in the red bounding box are the true moving objects. Top row: A fast drone approaching a moving event camera. Bottom row: Moving object behind netted background.

TABLE II

COMPARISON WITH STATE-OF-THE-ART CLASSICAL APPROACHES FOR EED, MOD, EV-IMO DATASETS.

| Method | Detection Rate (%) ↑ | | |
|---|---|---|---|
| | EED | MOD | EV-IMO |
| Mitrokhin *et al.* [28] | 88.93 | 70.12 | 48.79 |
| Stoffregen *et al.* [29] | 93.17 | - | - |
| 0-MMS [30] | **94.2** | **82.35** | **81.06** |
| Ours | **91.5** | **68.82** | **65.14** |

## V. CONCLUSION

We presented a first deep encoder-decoder Spiking Neural Network for a large-scale problem and demonstrated our architecture on the task of motion segmentation using data from a monocular event camera. Our novel spatio-temporal loss formulation takes full advantage of the spatio-temporal nature of the event data. We demonstrated the unique ability of our network, *SpikeMS*, for incremental prediction and showed its capability to generalize across a range of temporal intervals without explicit augmentation. A comprehensive qualitative and quantitative evaluation was provided using synthetic and real-world sequences from the EV-IMO, EED and MOD datasets. It was shown that *SpikeMS* achieves performance comparable to an ANN method, but with $50\times$ less power consumption.

TABLE III

PERFORMANCE METRICS FOR EV-IMO AND MOD DATASETS.

| Method | EV-IMO | | | | | MOD |
|---|---|---|---|---|---|---|
| | boxes | floor | wall | table | fast | |
| Num. Operations ($\times10^8$) | 0.42 | 0.34 | 0.52 | 0.38 | 0.53 | 0.83 |
| Energy benefit ($\times$) | 116.19 | 143.53 | 93.84 | 128.42 | 92.07 | 58.80 |

REFERENCES

[1] P. U. Diehl, D. Neil, J. Binas, M. Cook, S.-C. Liu, and M. Pfeiffer, "Fast-classifying, high-accuracy spiking deep networks through weight and threshold balancing," in *2015 International Joint Conference on Neural Networks (IJCNN)*, 2015, pp. 1–8.

[2] M. Pfeiffer and T. Pfeil, "Deep learning with spiking neurons: opportunities and challenges," *Frontiers in Neuroscience*, vol. 12, p. 774, 2018.

[3] N. J. Sanket, C. M. Parameshwara, C. D. Singh, A. V. Kuruttukulam, C. Fermüller, D. Scaramuzza, and Y. Aloimonos, "EVDodgeNet: deep dynamic obstacle dodging with event cameras," 2019.

[4] M. Gehrig, S. B. Shrestha, D. Mouritzen, and D. Scaramuzza, "Event-based angular velocity regression with spiking networks," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2020, pp. 4195–4202.

[5] F. Zenke and S. Ganguli, "Superspike: Supervised learning in multilayer spiking neural networks," *Neural Computation*, vol. 30, no. 6, pp. 1514–1541, 2018.

[6] S. B. Shrestha and G. Orchard, "Slayer: Spike layer error reassignment in time," *arXiv preprint arXiv:1810.08646*, 2018.

[7] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain *et al.*, "Loihi: A neuromorphic manycore processor with on-chip learning," *IEEE Micro*, vol. 38, no. 1, pp. 82–99, 2018.

[8] F. Akopyan, J. Sawada, A. Cassidy, R. Alvarez-Icaza, J. Arthur, P. Merolla, N. Imam, Y. Nakamura, P. Datta, G.-J. Nam *et al.*, "Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 34, no. 10, pp. 1537–1557, 2015.

[9] P. Lichtsteiner, C. Posch, and T. Delbruck, "A $128 \times 128$, 120 db $15 \mu$ s latency asynchronous temporal contrast vision sensor," *IEEE Journal of Solid-State Circuits*, vol. 43, no. 2, pp. 566–576, 2008.

[10] C. Posch, D. Matolin, and R. Wohlgenannt, "A qvga 143 db dynamic range frame-free pwm image sensor with lossless pixel-level video compression and time-domain cds," *IEEE Journal of Solid-State Circuits*, vol. 46, no. 1, pp. 259–275, 2010.

[11] A. L. Hodgkin and A. F. Huxley, "A quantitative description of membrane current and its application to conduction and excitation in nerve," *The Journal of Physiology*, vol. 117, no. 4, pp. 500–544, 1952.

[12] T. J. Sejnowski and G. Tesauro, "The Hebb rule for synaptic plasticity: algorithms and implementations," in *Neural Models of Plasticity*. Elsevier, 1989, pp. 94–103.

[13] B. Nessler, M. Pfeiffer, and W. Maass, "Stdp enables spiking neurons to detect hidden causes of their inputs," *Advances in neural information processing systems*, vol. 22, pp. 1357–1365, 2009.

[14] S. R. Kheradpisheh, M. Ganjtabesh, S. J. Thorpe, and T. Masquelier, "Stdp-based spiking deep convolutional neural networks for object recognition," *Neural Networks*, vol. 99, pp. 56–67, 2018.

[15] M. Evanusa, C. Fermüller, and Y. Aloimonos, "A deep 2-dimensional dynamical spiking neuronal network for temporal encoding trained with STDP," *arXiv preprint arXiv:2009.00581*, 2020.

[16] C. Lee, P. Panda, G. Srinivasan, and K. Roy, "Training deep spiking convolutional neural networks with stdp-based unsupervised pre-training followed by supervised fine-tuning," *Frontiers in Neuroscience*, vol. 12, p. 435, 2018.

[17] J. H. Lee, T. Delbruck, and M. Pfeiffer, "Training deep spiking neural networks using backpropagation," *Frontiers in neuroscience*, vol. 10, p. 508, 2016.

[18] D. Neil, M. Pfeiffer, and S.-C. Liu, "Learning to be efficient: Algorithms for training low-latency, low-compute deep spiking neural networks," in *Proceedings of the 31st Annual ACM Symposium on Applied Computing*, 2016, pp. 293–298.

[19] G. Orchard, A. Jayawant, G. K. Cohen, and N. Thakor, "Converting static image datasets to spiking neuromorphic datasets using saccades," *Frontiers in Neuroscience*, vol. 9, p. 437, 2015.

[20] G. Orchard, C. Meyer, R. Etienne-Cummings, C. Posch, N. Thakor, and R. Benosman, "Hfirst: A temporal approach to object recognition," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 37, no. 10, pp. 2028–2040, 2015.

[21] A. Amir, B. Taba, D. Berg, T. Melano, J. McKinstry, C. Di Nolfo, T. Nayak, A. Andreopoulos, G. Garreau, M. Mendoza *et al.*, "A low power, fully event-based gesture recognition system," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7243–7252.

[22] G. Haessig, A. Cassidy, R. Alvarez, and G. Benosman, R.and Orchard, "Spiking optical flow for event-based sensors using IBM's TrueNorth neurosynaptic system," pp. 860–870, 2018.

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

[24] G. Haessig, X. Berthelon, S.-H. Ieng, and R. Benosman, "A spiking neural network model of depth from defocus for event-based neuromorphic vision," *Scientific reports*, vol. 9, no. 1, pp. 1–11, 2019.

[25] A. Renner, M. Evanusa, and Y. Sandamirskaya, "Event-based attention and tracking on neuromorphic hardware," in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2019, pp. 1709–1716.

[26] C. Lee, A. K. Kosta, A. Z. Zhu, K. Chaney, K. Daniilidis, and K. Roy, "Spike-flownet: event-based optical flow estimation with energy-efficient hybrid neural networks," in *European Conference on Computer Vision*. Springer, 2020, pp. 366–382.

[27] F. Barranco, C. Fermüller, and E. Ros, "Real-time clustering and multi-target tracking using event-based sensors," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2018, pp. 5764–5769.

[28] A. Mitrokhin, C. Fermüller, C. Parameshwara, and Y. Aloimonos, "Event-based moving object detection and tracking," in *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2018, pp. 1–9.

[29] T. Stoffregen, G. Gallego, T. Drummond, L. Kleeman, and D. Scaramuzza, "Event-based motion segmentation by motion compensation," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 7244–7253.

[30] C. M. Parameshwara, N. J. Sanket, C. Deep Singh, C. Fermüller, and Y. Aloimonos, "0-mms: Zero-shot multi-motion segmentation with a monocular event camera," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2021.

[31] A. Mitrokhin, C. Ye, C. Fermüller, Y. Aloimonos, and T. Delbruck, "EV-IMO: motion segmentation dataset and learning pipeline for event cameras," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2019.

[32] A. Mitrokhin, Z. Hua, C. Fermüller, and Y. Aloimonos, "Learning visual motion segmentation using event surfaces," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 414–14 423.

[33] W. Gerstner, "Time structure of the activity in neural network models," *Physical review E*, vol. 51, no. 1, p. 738, 1995.

[34] ——, "Spike-response model," *Scholarpedia*, vol. 3, no. 12, p. 1343, 2008, revision #91800.

[35] M. van Rossum, "A novel spike distance," *Neural Computation*, vol. 13, pp. 751–763, 04 2001.

[36] T. Kreuz, "Measures of spike train synchrony," *Scholarpedia*, vol. 6, no. 10, p. 11934, 2011, revision #190333.

[37] C. R. Qi, L. Yi, H. Su, and L. J. Guibas, "Pointnet++: Deep hierarchical feature learning on point sets in a metric space," 2017.

[38] R. Kalman, "A new approach to linear filtering and prediction problems," *Journal of Basic Engineering*, vol. 82, no. 1, pp. 35–45, 1960.