# MDOE: A Spatiotemporal Event Representation Considering the Magnitude and Density of Events

Fuqiang Gu<sup>®</sup>, *Member, IEEE*, Yong Lee, Yuan Zhuang<sup>®</sup>, *Member, IEEE*, You Li<sup>®</sup>, Jingbin Liu<sup>®</sup>, Fangwen Yu<sup>®</sup>, *Member, IEEE*, Ruiyuan Li<sup>®</sup>, and Chao Chen<sup>®</sup>

*Abstract*—Event-based sensors (e.g., DVS cameras) are capable of higher dynamic range, higher temporal resolution, lower time latency, and better power efficiency compared to conventional devices (e.g., RGB cameras). However, learning from these sensors remains challenging; event-based sensors output a stream of asynchronous events, which cannot be directly used by stateof-the-art convolutional neural networks (CNNs). In this paper, we present a novel event-based representation called MDOE that considers both the magnitude and density of events. Compared to existing representations, which discard one or more types of information about event polarity, temporal information, and/or density, MDOE contains richer information about events. It has two benefits: (i) it is a conceptually-simple generic representation that is task-independent; (ii) it achieves superior performance relative to existing representations on a variety of event-based datasets.

*Index Terms*—Event-based learning, event sensors, deep learning, object recognition, convolutional neural networks.

Manuscript received 7 February 2022; accepted 14 June 2022. Date of publication 27 June 2022; date of current version 6 July 2022. This letter was recommended for publication by Associate Editor Y. Bekiroglu and Editor M. Vincze upon evaluation of the reviewers' comments. This work was supported in part by the National Natural Science Foundation of China under Grants 42174050, 41874031, 62172066, 61976168, and 42111530064, in part by the Special Fund of Hubei Luojia Laboratory under Grant 220100007, in part by the Shenzhen Science and Technology Program under Grant JCYJ20210324123611032, in part by the Excellent Youth Foundation of Hubei Scientific Committee under Grant 2021CFA040, in part by the Guangdong Basic and Applied Basic Research Foundation under Grant 2021A1515110343, in part by the Chongqing Innovative and Entrepreneurial Program of Returned Overseas Chinese Scholars under Grant cx2021047. (*Corresponding author: Fugiang Gu.*)

Fuqiang Gu, Ruiyuan Li, and Chao Chen are with the College of Computer Science, Chongqing University, Chongqing 400044, China (e-mail: gufq@ cqu.edu.cn; liruiyuan@cqu.edu.cn; cschaochen@cqu.edu.cn).

Yong Lee is with the School of Mechanical and Electronic Engineering, Wuhan University of Technology, Wuhan 430070, China (e-mail: yonglee@ whut.edu.cn).

Yuan Zhuang is with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan, China, the Wuhan Institute of Quantum Technology, and Wuhan University Shenzhen Research Institute, Wuhan 430079, China (e-mail: yuan.zhuang@whu.edu.cn).

You Li is with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University and the Hubei Luojia Laboratory, Wuhan 430079, China (e-mail: liyou@whu.edu.cn).

Jingbin Liu is with the State Key Laboratory of Information Engineering in Surveying, Mapping and Remote Sensing, Wuhan University, Wuhan 430079, China (e-mail: jingbin.liu@whu.edu.cn).

Fangwen Yu is with the Department of Precision Instrument, Tsinghua University, Beijing 100084, China (e-mail: yufangwen@tsinghua.edu.cn).

Digital Object Identifier 10.1109/LRA.2022.3186523



Fig. 1. Illustration of event-based representation, which converts asynchronous events into grid-like voxels.

#### I. INTRODUCTION

**WENT-BASED** sensors, such as dynamic vision sensors (DVS) [1], [2] and NeuTouch tactile sensors [3], are bioinspired devices that mimic the efficient event-driven communication mechanisms of the brain. Different from conventional sensors (e.g., RGB cameras), which synchronously capture the scene at a fixed rate, event-based sensors asynchronously report the changes (called events) of the scene. For example, DVS cameras capture brightness changes for each pixel independently rather than intensity images as RGB cameras do. Compared to conventional sensors, event-based sensors are advantageous of higher dynamic range, higher temporal resolution, lower time latency, and higher power efficiency [4].

While event-based sensors are attractive for many applications (e.g., self-driving vehicles, security and surveillance, industrial automation), they can not be directly used by state-ofthe-art CNNs since their outputs are a stream of asynchronous events. Single event alone contains little information about the scene. To make use of event-based data, certain methods have been proposed to learn useful representations that can be exploited by CNNs. Fig. 1 gives an example of converting events into grid-like voxels.

Among these representations, of particular interest for this work are Event Frame [5], Surface of Active Events (SAE) [6], Event Count [7], Voxel Grid [8], Hierarchy of Time Surfaces (HOTS) [9], Histograms of Averaged Time Surfaces (HATS) [10], Event Spike Tensor (EST) [4], and Matrix-LSTM [11]. These representations have been used for object recognition [4], [9], [10], optical flow estimation [4], [11], visual-inertial odometry [12], etc. However, existing representations discard one or more types of information about event polarity, density, and/or temporal information, which may degrade the performance of downstream applications.

2377-3766 © 2022 IEEE. Personal use is permitted, but republication/redistribution requires IEEE permission. See https://www.ieee.org/publications/rights/index.html for more information.

Dimension	Polarity	Time	Density
$H \times W$	×	×	×
$2 \times H \times W$	$\checkmark$	×	×
$2 \times H \times W$	$\checkmark$	×	$\times$
-	×	×	×
$2 \times H \times W$	$\checkmark$	×	×
mps $4 \times H \times W$	$\checkmark$	×	×
$C \times H \times W$	×	$\checkmark$	×
$2 \times C \times H \times W$	$\checkmark$	$\checkmark$	$\times$
$C \times H \times W$	×	$\checkmark$	×
$4 \times C \times H \times W$	$\checkmark$	$\checkmark$	$\checkmark$
	$\begin{tabular}{ c c c c } \hline \textbf{Dimension} \\ \hline H \times W \\ 2 \times H \times W \\ 2 \times H \times W \\ \hline \hline & 2 \times H \times W \\ \hline & 2 \times H \times W \\ 0 \times H \times W \\ 2 \times C \times H \times W \\ 2 \times C \times H \times W \\ 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 4 \times C \times H \times W \\ \hline \hline & 5 \times C \times H \times W \\ \hline \hline \hline & 5 \times C \times H \times W \\ \hline \hline \hline & 5 \times C \times H \times W \\ \hline \hline \hline \hline & 5 \times C \times H \times W \\ \hline \hline$	$\begin{tabular}{ c c c c } \hline \textbf{Dimension} & \textbf{Polarity} \\ \hline H \times W & \checkmark \\ 2 \times H \times W & \checkmark \\ 2 \times H \times W & \checkmark \\ 2 \times H \times W & \checkmark \\ \hline 2 \times H \times W & \checkmark \\ C \times H \times W & \checkmark \\ 2 \times C \times H \times W & \checkmark \\ 2 \times C \times H \times W & \checkmark \\ 4 \times C \times H \times W & \checkmark \\ \hline \end{array}$	$\begin{tabular}{ c c c c c } \hline \textbf{Dimension} & \textbf{Polarity} & \textbf{Time} \\ \hline H \times W & \times & \times & \\ 2 \times H \times W & \checkmark & \times & \\ 2 \times H \times W & \checkmark & \times & \\ 2 \times H \times W & \checkmark & & \\ 2 \times H \times W & \checkmark & \checkmark & \\ 2 \times H \times W & \checkmark & \checkmark & \\ 2 \times C \times H \times W & \checkmark & \checkmark & \\ 2 \times C \times H \times W & \checkmark & \checkmark & \\ 2 \times C \times H \times W & \checkmark & \checkmark & \\ 4 \times C \times H \times W & \checkmark & \checkmark & \\ 4 \times C \times H \times W & \checkmark & \checkmark & \\ \hline \end{tabular}$

 TABLE I

 COMPARISON OF DIFFERENT REPRESENTATIONS FOR EVENT-BASED DATA

C denotes the number of temporal bins, H and W represent the height and width of image, respectively.  $\checkmark$  denotes a representation contains such information, while  $\times$  means a representation discards such information.

This paper presents MDOE (Magnitude and Density Of Events), a novel representation for asynchronous event-based data. MDOE is obtained by concatenating the Magnitude of Events (MOE) and Density of Events (DOE). Compared to existing event-based representations, the proposed MDOE contains richer information about the polarity, temporal information, as well as density of events. Our representation has two benefits: (i) it is a conceptually-simple generic representation that is task-independent; (ii) it achieves superior performance relative to existing representations on a variety of event-based datasets.

Table I compares popular representations in terms of output dimension, and whether to contain polarity information, temporal information, and event density. From the table, one can find that the proposed MDOE contains more information about the scene than existing representations, and hence is more likely to give better performance. This is also justified by the higher test accuracy achieved by the proposed method on several public event-based datasets.

# II. RELATED WORK

Event-based sensors have been successfully applied for a variety of tasks such as object recognition [4], [9], [10], [14], gesture recognition [15], [16], optical flow estimation [4], [11], and visual-inertial odometry [12]. Despite their advantages in terms of dynamic range, temporal resolution, and time latency, event-based sensors are still not in the mainstream due mainly to the difficulty of learning good representations from sparse and asynchronous stream of events. In this section, we provide a brief overview of related works on event-based representations.

Handcrafted Representations: A naive way to represent asynchronous events is to use the number of events occurred at each pixel [5], but such method discards the temporal information and polarity information of events, which may degrade the accuracy of downstream applications. A similar representation is also used in [7], but the work in [7] considers event polarity information in addition to event counts. Instead of counting events for each pixel, SAE [6] uses the most recent timestamps as event-based representation. In HOTS [9], the concept of *Time Surface* is introduced to describe the recent history of events in the spatial neighborhood of an event. However, HOTS suffers from high latency since computing time surfaces is relatively expensive, and is sensitive to noisy events as it considers only the timestamps of the last events in the neighborhood of an event. A popular extension to HOTS is HATS [10], which uses a local memory mechanism. Instead of using the last event in the neighborhood of an event, HATS computes time surfaces by considering all past events in a spatiotemporal window. In this way, noisy events have less impact on the task. Another variant of time surface is introduced [13], in which a four-channel grid is built containing the number of positive and negative events occurred at each pixel as well as the most recent timestamps.

Spike-based Representations: Due to their event-driven property, SNNs have been applied in several event-based tasks such as edge detection [17], object classification [14], [18], and gesture recognition [15], [19]. SNNs are more biologically plausible than deep learning models and can be run on power-efficient neuromorphic processors such as the IBM TrueNorth [20] and Intel Loihi [21]. However, SNNs are non-differentiable and hence cannot directly use well-established back-propagation methods to train the models, limiting the usability of SNNs in real-world scenarios. A typical method to address this problem is training a conventional neural network with frame-based data and converting the learned parameters to a SNN for event-based data [22]. In recent years, several works tried to approximate the derivative of the spike function [23], [24]. However, the achieved accuracy of spike-based methods is still not comparable to state-of-the-art deep learning methods.

*Learning-based Representations:* In recent years, several learning-based representations have been proposed to deal with asynchronous events. In [4], a grid-based representation called EST is proposed to learn end-to-end features directly from asynchronous event-based data through differentiable kernel convolution and quantization. A multi-layer perceptron is used to learn a trilinear kernel that generates voxel-grid features, which can be tuned to the downstream task. In Matrix-LSTM [11], a grid of Long Short-Term Memory (LSTM) cells are used to learn end-to-end an event-based representation, which can capture local temporal features while retaining spatial structures. In [25], a variant of LSTM network called PhaseLSTM that uses a time gate is proposed to learn the precise timings of incoming events, which extracts relevant features with a word embedding

layer. While PhaseLSTM achieves good performance on simple datasets, it may not capture general features well due to the limited representation ability of its embedding layer [26]. In [16], the event stream is treated as a set of 3D points in space-time, and a neural network called PointNet++, which is originally developed for recognizing 3D point clouds, is used to extract event-based features. In [27], a general frame is proposed to deal with event data by using the asynchronous sparse convolutional network that is converted from a synchronous model.

Most closely related to this paper are methods that map raw asynchronous events into voxel or 4D representations, mainly including Voxel Grid [8] and EST [4]. The Voxel Grid can capture the full spatiotemporal distribution of the events by discretizing the time domain and then accumulating the events in a linear way, but it discards the polarity information about events. EST [4] extends the Voxel Grid by separating the polarity and building a 4D representation with a learnable trilinear kernel. Nevertheless, EST does not consider all past events nor the density of events. In addition, EST is an end-to-end representation, which needs to be tuned to downstream tasks. Compared to these representations, our representation contains richer information and is a more generic representation that is task-independent.

#### III. PROPOSED METHOD

The output of event-based sensors is a stream of asynchronous events that encode the location, time, polarity (sign) of the changes. It can be described by a sequence:

$$\varepsilon = \{e_i\}_{i=1}^{I} = \{\mathbf{x}_i, t_i, p_i\}_{i=1}^{I}$$
(1)

where  $\mathbf{x}_i$  is the location (for DVS cameras,  $\mathbf{x}_i = (x_i, y_i)$  is the coordinates of the pixel triggering the event),  $t_i$  is the timestamp when the event is generated, and  $p_i$  is the polarity of the event. The polarity takes two values: 1 and -1, representing ON and OFF events, respectively. I is the number of events in the sequence.

To process these asynchronous events using feature extractors (e.g., CNN), we need to convert the event sequence into a fixedsize representation. Although there are already many representation methods proposed in the literature, such as HATS [10] and EST [4], they discard certain information about the events (e.g., temporal information, polarity, or density), which may affect the performance of downstream applications (e.g., object recognition, optical flow estimation).

In this paper, we propose a novel representation called MDOE, which contains richer information about the events. MDOE is the combination of the MOE and DOE. The MOE can be considered as an unbounded membrane potential that accumulates all past events from the same pixel. By contrast, DOE reflects the event density information over a temporal bin. Combining MOE and DOE allows to capture the information about all past events and the events occurred during each temporal bin.

Specifically, the MOE is inspired by the dynamics of neurons in a spiking neural network (SNN) where a post-neuron integrates the inputs from all the pre-neurons and fires when its membrane potential surpasses a threshold. The dynamics of



Fig. 2. (a) The dynamics of a neuron's membrane potential in a SNN, where  $u_T$  and  $u_R$  are the firing threshold and reset value. When the membrane potential of a neuron surpasses the threshold  $u_T$ , it will fire (output a spike) and its membrane potential will be reset to  $u_R$ ; (b) The magnitude of events of a pixel is similar to the membrane potential, but it has no firing process and refractory period.

neurons' membrane potential in SNNs [28] can be described as:

$$\tau \frac{du(t)}{dt} = -u(t) + \sum_{i} w_i x'_i \tag{2}$$

where u(t) denotes the internal membrane potential of a neuron at time t,  $\sum_i w_i x'_i$  is the weighted summations of the inputs from pre-neurons, and  $\tau$  is a time constant. Fig. 2(a) visualizes the dynamic of a neuron's membrane potential. It can be seen that during a period of time, the denser the spikes, the larger the membrane potential. Fig. 2(b) shows the proposed MOE for the same event sequence. The main difference is that in the MOE, there is no firing process and hence the magnitude will not be reset. The motivation behind this is that resetting the membrane potential (or magnitude) may lose certain information about the event stream. Thus, by sampling the MOE at regular time interval for each polarity, we can obtain the representation of dimension  $2 \times C \times H \times W$ , where H and W are the spatial height and width of the event image, respectively, and C is the number of temporal bins.

Apart from the MOE, we also integrate the DOE of each temporal window into the MDOE representation. An earlier work [7] has considered the total number of triggered events for each pixel, which may not be informative enough since it can not reflect the event distribution. Instead of considering the total number of events during a time period, we divide the time period into different time windows and count the events over each temporal window pixel by pixel, resulting a vector of size C for each pixel. By doing this for each polarity, we can obtain the DOE representation with dimension  $2 \times C \times H \times W$ .

Both MOE and DOE are 4-dimensional (namely,  $2 \times C \times H \times W$ ), concatenating them results in a 4-dimensional MDOE

(3)

representation  $(4 \times C \times H \times W)$ . Let *E* be the MOE, we describe it as follows:

$$E_{\pm}(x_l, y_m, c_n) = \lambda E_{\pm}(x_l, y_m, c_{n-1})$$
$$+ \sum_{e_i \in \varepsilon_{\pm} \{t_{n-1} \le t_i \le t_n\}} \delta(x_l - x_i) \delta(y_m - y_i) k(t_n - t_i)$$

$$t_n = t_1 + (c_n + 1)\Delta T \tag{4}$$

$$\delta(a) = \begin{cases} 1, \text{ if } a = 0\\ 0, \text{ otherwise} \end{cases}$$
(5)

$$k(a) = \max(0, 1 - |a|) \tag{6}$$

where  $\varepsilon_+$  and  $\varepsilon_-$  are event sequences with positive polarity and negative polarity, respectively.  $\lambda$  is the decaying factor that determines how much past information is considered, and it ranges from 0 to 1. k(a) is a bilinear kernel that emulates the neuron dynamics,  $(x_l, y_m, c_n)$  is the spatiotemporal coordinates on a voxel grid, and  $x_l \in \{0, 1, \dots, W-1\}, y_m \in \{0, 1, \dots, H-1\}$ , and  $c_n \in \{0, 1, \dots, C-1\}$ .  $\Delta T$  is temporal bin size. The MOE can be considered as an unbounded membrane potential that accumulates all past events from the same pixel. The closer the events distributed during a time period, the larger the value of MOE.

Similarly, we can write the DOE (denoted by V) as:

$$V_{\pm}(x_l, y_m, c_n) = \sum_{e_i \in \varepsilon_{\pm}\{t_{n-1} < t_i \le t_n\}} \delta(x_l - x_i) \delta(y_m - y_i) f(e_i)$$
(7)

where  $f(e_i)$  is the indicator function that ensures to count only events falling between the specified time period.  $f(e_i)$  is written as

$$f(e_i) = \begin{cases} 1, \text{ if } e_i \in \varepsilon_{\pm} \{ t_{n-1} < t_i \le t_n \} \\ 0, \text{ otherwise} \end{cases}$$
(8)

To avoid overflow when the MOE makes sum operation, we normalize the timestamps events to the range of [0, 1] before calculation. For real-time applications, one can convert the timestamps into seconds. DOE reflects the number of events over a time window. The more events triggered during the time period, the larger the value of DOE.

To reduce the effect of sensor noises (e.g., reporting different number of events for the exact scene at different time), the resulting MOE ( $E_{\pm}$ ) and DOE ( $V_{\pm}$ ) are normalized via dividing them by their respective maximum values across pixels and temporal bins. The MDOE (denoted by S) is the concatenation of MOE and DOE along the polarity, namely

$$S(p, x_l, y_m, c_n) = Cat_{\pm}(E_{\pm}(x_l, y_m, c_n), V_{\pm}(x_l, y_m, c_n))$$
(9)

where *Cat* represents the concatenation operation, p indicates the polarity of MOE and DOE, and  $p \in \{0, 1, 2, 3\}$ , where p = 0and p = 1 represent the ON and OFF polarity for MOE, while p = 2 and p = 3 correspond to the ON and OFF polarity for DOE. For better clarity, we provide the pseudo-code of calculating the MDOE in Algorithm 1.

**Input** : A stream of events  $\varepsilon = \{\mathbf{x}_i, t_i, p_i\}_{i=1}^{I}$ , the spatiotemporal dimension of output dimension (C, H, W), and decay factor  $\lambda$ **Output:** MDOE representation with dimension  $4 \times C \times H \times W$ 1 Separate the event sequence by their polarity and obtain  $\varepsilon_+$  and  $\varepsilon_-$ ; 2 Calculate the temporal bin by  $\Delta T = \frac{t_I - t_1}{C}$ ; 3 for  $x_l \in \{0, \cdots, W-1\}, y_m \in \{0, \cdots, H-1\}, c_n \in \{0, 1, \cdots, C-1\}$ do 4 Initialize  $E_{\pm}(x_l, y_m, c_n)$  and  $V_{\pm}(x_l, y_m, c_n)$  to 0; 5 end 6 for  $x_l \in \{0, \cdots, W-1\}, y_m \in \{0, \cdots, H-1\}, c_n \in \{0, 1, \cdots, C-1\}$ do  $t_{n-1} = t_1 + c_n \Delta T$ ; //Compute  $t_{n-1}$ 7 8  $t_n = t_1 + (c_n + 1)\Delta T$ ; //Compute  $t_n$ for  $e_i \in \varepsilon_+$  do 9 if  $t_i \in (t_{n-1}, t_n]$  then 10  $E_{+}(x_{l}, y_{m}, c_{n}) = \lambda E_{+}(x_{l}, y_{m}, c_{n-1}) +$ 11  $\delta(x_l - x_i)\delta(y_m - y_i)k(t_n - t_i);$  $V_+(x_l, y_m, c_n) \coloneqq \delta(x_l - x_i)\delta(y_m - y_i)f(e_i);$ 12 13 end 14 end 15 for  $e_i$  $\in \varepsilon_{-}$  do if  $t_i \in (t_{n-1}, t_n]$  then 16  $E_{-}(x_{l}, y_{m}, c_{n}) = \lambda E_{-}(x_{l}, y_{m}, c_{n-1}) +$ 17  $\delta(x_l - x_i)\delta(y_m - y_i)k(t_n - t_i);$  $V_{-}(x_l, y_m, c_n) \coloneqq \delta(x_l - x_i)\delta(y_m - y_i)f(e_i);$ 18 19 end 20 end 21 end 22 for  $x_l \in \{0, \cdots, W-1\}, y_m \in \{0, \cdots, H-1\}, c_n \in \{0, 1, \cdots, C-1\}$ do  $E_{\pm}(x_l, y_m, c_n) = E_{\pm}(x_l, y_m, c_n) / max(E_{\pm});$ 23 24  $V_{\pm}(x_l, y_m, c_n) = V_{\pm}(x_l, y_m, c_n) / max(V_{\pm});$ 25 end 26  $S = Cat_{\pm}(E_{\pm}, V_{\pm})$ ; //Concatenate MOE and DOE to generate MDOE; 27 Return S;

Algorithm 1: Calculation of MDOE.



Fig. 3. Sample examples of the used three public event-based datasets. For better visualization, the Event Frame is used to demonstrate the sample of N-Caltech101 and N-Cars datasets. For EvTouch-Object dataset, which is an event-based tactile dataset, we demonstrate the images of objects used in the data collection.

#### **IV. EXPERIMENTS AND RESULTS**

## A. Datasets

We evaluate the proposed method using three publicly available event-based datasets: N-Caltech101 [29], N-Cars [10], and EvTouch-Objects [3], [14]. N-Caltech101 (Neuromorphic-Caltech101) is an event-based version of the popular Caltech101 dataset [30]. To convert the images to event sequences, an ATIS event camera was installed on a motorized pan-tilt unit, and automatically moved while pointing at images from the original dataset (Caltech101) that were shown on a LCD monitor. N-Cars (Neuromorphic-Cars) is a real-world event-based dataset for recognizing whether a car is present in a scene. It was recorded using an ATIS camera that was mounted behind the windshield of a car. EvTouch-Objects dataset is an event-driven tactile dataset for object recognition. It was collected using a 7-DoF

Representation	N-Caltech101 (Std)	N-Cars (Std)
MOE	84.9 (0.9)	93.5 (0.5)
DOE	85.6 (0.4)	92.9 (0.5)
MDOE	85.9 (0.1)	94.2 (0.4)

Franka Emika Panda arm equipped with a Robotiq 2F-140 gripper, equipped with a NeuTouch event-based tactile sensor [3] and an ACES decoder [31] to decode the sensor signals into spikes. This dataset comprises tactile data from 36 object classes. For each object class, 20 samples were collected, yielding a total of 720 samples. Fig. 3 gives some sample examples of the used datasets. For better visualization, the Event Frame method is used to convert events into images.

## B. Implementation Details

We evaluate the proposed method using a ResNet-34 architecture [32] for easier comparison with baseline methods that share the same network architecture. We also evaluate our method using other state-of-the-art networks including VGG-19 [33], MobileNet-V2 [34], and Inception-V3 [35]. All the networks are pretrained on ImageNet [36]. Since the number of input channels and output classes for our case is different from these pre-trained models, we adopt the approach used in [4], [7] and replace the first and last layer of the pre-trained models with random weights, and then fine-tune all the parameters on the task. Note that the original ResNet-34 is for 2D images, while our event presentation is a 4D tensor. To adapt to the ResNet-34, we concatenate the event-based representation along the polarity and temporal dimension as channels.

The Adam optimizer [37] is used to train the model by minimizing the cross-entropy loss. The initial learning rate is set to  $1 \times 10^{-4}$  and is reduced by a factor of 0.5 every 10 iterations. The total number of iterations is set to 200. We use a batch size of 4 for all the datasets. To conduct a robust evaluation, we run the model on each dataset for multiple rounds with different random seeds, and report the mean and standard deviation values. We perform early stopping [38] on a validation set using the splits provided by the EST [4] on N-Caltech101 and 20% of the training data on N-Cars. For EvTouch-Objects, there is no validation dataset provided online or in publications, we report the best test accuracy during the 200 epochs. The default value of decay factor in MDOE is set to 0.1.

We have implemented the MDOE in PyTorch and the source code will be available at *https://github.com/uisc-cqu/MDOE*.

## C. Ablation Study

*Contribution of MOE and DOE:* To evaluate the contribution of each component of the proposed MDOE, we first compare the object recognition accuracy of using MDOE, MOE, and DOE, respectively. We use ResNet-34 as the classifier, and let all the methods share the same configuration (e.g., training epochs, learning rate). A temporal bin of 9 is used for all the methods. Table II demonstrates that MDOE performs better than MOE and DOE. Specifically, the MOE achieves a higher accuracy on

TABLE III THE TEST ACCURACY (%) OF USING MDOE WITH DIFFERENT NUMBER OF TEMPORAL BINS

Temporal Bin	N-Caltech101 (Std)	N-Cars (Std)
1	83.4 (0.4)	93.6 (0.3)
9	85.9 (0.1)	94.2 (0.4)
16	86.2 (0.7)	94.1 (0.4)
32	86.6 (0.8)	94.2 (0.4)
64	86.4 (0.6)	94.3 (0.4)

TABLE IV THE TEST ACCURACY (%) OF USING MDOE ACROSS DIFFERENT MODELS

Model	N-Caltech101 (Std)	N-Cars (Std)
VGG-19	81.9 (0.8)	91.9 (0.6)
Inception-V3	84.9 (0.6)	93.6 (0.8)
MobileNet-V2	85.3 (0.6)	94.8 (0.8)
ResNet-34	85.9 (0.1)	94.2 (0.4)

N-Cars but a lower accuracy on N-Caltech101 than the DOE. The proposed method that combines the MOE and DOE outperforms the MOE and DOE alone on both datasets. This might because MDOE contains richer information about the events, including not only the polarity and temporal information, but also the density. This makes it more robust when conducted on different datasets.

Analysis of different temporal bins: We then analyze the effect of different temporal bins on the performance of our method with the number of temporal bins varying from 1 to 64. Table III shows that the number of temporal bins has an impact on the proposed method. The achieved accuracy increases from 83.4% to 85.9% on N-Caltech101 and from 93.6% to 94.2% on N-Cars when the temporal bin rises from 1 to 9. The test accuracy continues to rise until the temporal bin reaches 32 for N-Caltech101. However, it is not the same case for N-Cars dataset where the number of temporal bins has a negligible impact on the proposed method after it reaches to 9. In the following, we set the temporal bin to 9 to analyze the effect of other factors and compare with baseline methods.

Performance across different models: We also evaluate the proposed method across four state-of-the-art deep learning architectures, namely ResNet-34 [32], VGG-19 [33], MobileNet-V2 [34], and Inception-V3 [35]. Table IV shows that the test accuracy of the proposed method using different network architectures on N-Caltech101 and N-Cars datasets. It can be seen that MobileNet-V2 and ResNet-34 can result in a higher accuracy than VGG-19 and Inception-V3 on both datasets. Specifically, ResNet-34 performs the best on N-Caltech101, which is followed by MobileNet-V2 and Inception-V3. By contrast, MobileNet-V2 achieves the highest accuracy on N-Cars, followed by ResNet-34 and Inception-V3. VGG-19 witnesses the lowest accuracy on both datasets among the four architectures.

Analysis of time latency: Time latency is the consumed time period to accumulate evidence in making a decision on the object class [10], which is crucial for the applications requiring fast reaction time. We compare the time latency of the proposed method with Event Frame [5], Event Count [7], and Voxel Grid [8] on N-Caltech101. As shown in Fig. 4, our method

Authorized licensed use limited to: Tsinghua University. Downloaded on May 03,2025 at 04:03:44 UTC from IEEE Xplore. Restrictions apply.



Fig. 4. Time latency as a function on N-Caltech101.



Fig. 5. Effect of using different ratios of training data (N-Caltech101).

outperforms all the baseline methods and can achieve about 82% accuracy using only the first 30 ms of the samples. This makes the proposed method a good match for applications with highly real-time demand such as autonomous navigation. We can also find that a larger time latency does not necessarily mean a higher accuracy. This is especially true for the Event Frame and Event Count methods where the achieved accuracy degrades when the time latency reaches up to 60 ms and 100 ms respectively. This might be that Event Frame and Event Count do not consider the temporal information of events, and a larger time latency may introduce more noise.

Ratio of used training data: We analyze the effect of using different ratios of training data on the performance of the proposed method, and compare it with Event Frame [5], Event Count [7], and Voxel Grid [8] on N-Caltech101. For fair comparison, we let our method and all the baseline methods share the same configuration (e.g., ResNet-34 network, random seed, and learning rate). As shown in Fig. 5, the classification accuracy of all the methods increases as the ratio of used training data goes, and the proposed method is consistently more accurate than the baselines. It can be seen that our method achieves an accuracy of about 63% when using only 10% training data, which is much higher than Voxel Grid (about 59%), Event Count (55%) and Event Frame (53%). This implies that the proposed method can better deal with small datasets, and such feature is very useful for cases where collecting a large amount of data is troublesome or impractical.

*Impact of Decay Factor:* The impact of the decay factor on the performance of the proposed method is also investigated. The considered value of the decay factor ranges from 0 to 1,



Fig. 6. The average test accuracy with standard deviation using different values of the decay factor in MDOE.

TABLE V THE STORAGE AND TIME COMPLEXITY OF DIFFERENT METHODS.

Method	Storage	Time
Event Frame	$\mathcal{O}(BHW)$	$\mathcal{O}(N)$
Event Count Image	$\mathcal{O}(BHW)$	$\mathcal{O}(N)$
Voxel Grid	$\mathcal{O}(BCHW)$	$\mathcal{O}(CN)$
EST	$\mathcal{O}(BCHW)$	$\mathcal{O}(CN)$
MDOE	$\mathcal{O}(BCHW)$	$\mathcal{O}(CN)$

where 0 represents the MOE at the current time window does not use the MOE from previous time windows, and 1 means that all the MOE from past time windows are considered at the current time window. Fig. 6 shows that the decay factor has an impact on the performance of MDOE. The best performance of about 85.9% can be achieved when the decay factor is set to 0.1, 0.2, 0.5, and 0.9, and the smallest standard deviation is witnessed at the value of 0.1. It is also observed that discarding (the corresponding decay factor is 0) or keeping (decay factor is 1) all the information from past time windows does not result in the best performance. This implies that we should keep the information from past windows to some extent to achieve better performance.

Complexity Analysis: The storage cost and computational cost of the proposed method are also analyzed. We compare the proposed method with Event Frame [5], Event Count Image [7], Voxel Grid [8], and EST [4]. The results are given in the Table V, where B is the batch size, C is the number of time windows (or temporal bins), N is the number of events, H and W are the height and width of events, respectively. From the table, we can see that the storage and computational complexity of the MDOE is similar to Voxel Grid and EST, all of which consider the temporal information of events. However, the achieved accuracy of MDOE is much higher than other methods.

#### D. Results on Visual Object Recognition

We compare the proposed method with state-of-the-art methods on two visual object recognition datasets (N-Caltech101 and N-Cars). The baseline methods include HATS [9], Matrix-LSTM [11], E2VID [12], Event Frame [5], Event Count Image [7], Voxel Grid [8], and EST [4]. All the baseline methods except E2VID (that uses ResNet-18) use ResNet-34 as the classifier. For HATS, Matrix-LSTM, E2VID, and EST, we directly

TABLE VI THE TEST ACCURACY (%) COMPARISON OF THE PROPOSED MDOE WITH BASELINE METHODS. BOTH OUR METHOD AND BASELINE METHODS USE RESNET-34 AS THE CLASSIFIER BY DEFAULT. NOTE THAT FOR FAIR COMPARISON WITH E2VID, WE CONDUCT A SEPARATE EXPERIMENT USING RESNET-18 CLASSIFIER AND DIRECTLY REPORT THE BEST TEST ACCURACY WITHOUT USING A VALIDATION DATASET FOR EARLY STOPPING

Method	Classifier	N-Caltech101 (Std)	N-Cars (Std)
HATS		69.1	90.9
EST		81.7	92.5
Matrix-LSTM		83.5 (1.2)	92.2 (0.7)
Event Frame	ResNet-34	80.8 (0.6)	93.4 (0.5)
Event Count		81.1 (0.4)	93.5 (0.6)
Voxel Grid		83.9 (0.3)	92.9 (0.6)
<b>MDOE</b> (Ours)		85.9 (0.1)	94.2 (0.4)
E2VID	DecNet 19	86.6	91.0
<b>MDOE</b> (Ours)	Resinet-18	88.2 (0.4)	94.6 (0.3)

report the recognition accuracy provided in the corresponding reference. For other methods including Event Frame, Event Count, and Voxel Grid, we implement these methods and report the mean accuracy with standard deviation after running the model multiple times. To conduct a fair comparison, we use the same number of channels in the event representations (Voxel Grid and EST) and data augmentation procedures (random shifts and random horizontal flips).

Table VI summarizes the results of these methods. Our method outperforms all the baseline methods on both datasets. In particular, we can see that the proposed method can achieve an accuracy of 85.9% on N-Caltech101 and 94.2% on N-Cars, which is much higher than Voxel Grid (83.9% and 92.9% respectively), Event Count Image (81.1% and 93.5%), and Matrix-LSTM (83.5% and 92.2%). It can be also seen that there is not a baseline method that consistently exceed the other baseline methods on both datasets. For the comparison with E2VID, we conduct a separate experiment using ResNet-18 classifier and directly report the best test accuracy without using a validation dataset for early stopping. We can see that our method outperforms E2VID by about 1.6% and 3.6% on N-Caltech101 and N-Cars, respectively. It can also be observed that the achieved accuracy of both our method and E2VID is much higher than other methods. This is because the best test accuracy is reported without using early stopping via a validation dataset.

# E. Results on Tactile Sensing

We also evaluate the performance of the proposed MDOE on the EvTouch-Objects dataset. Since the EvTouch-Objects data were collected using irregularly-placed NeuTouch taxels, we cannot directly use state-of-the-art CNNs to deal with them. To make use of CNNs, we organize the tactile data in a grid structure of  $11 \times 11$ . Fig. 7 shows that the spatial distribution of NeuTouch taxels, and the corresponding grid-based representation.

We compare the proposed MDOE with a spiking-based method called TactileSGNet [14], and four frame-like methods, namely Event Frame [5], Event Count [7], Voxel Grid [8], and EST [4]. All the frame-like methods and the proposed MDOE use ResNet-34 as the classifier and share the same configuration (e.g., training epochs, learning rate). Table VII shows the



Fig. 7. (a) The spatial distribution of 39 NeuTouch taxels; (b) Grid-based representation of irregular tactile data.

TABLE VII THE TEST ACCURACY (%) OF USING DIFFERENT REPRESENTATIONS FOR TACTILE OBJECT CLASSIFICATION

Method	Average Accuracy (Std)
TactileSGNet	89.4 (0.6)
Event Frame	62.9 (1.1)
Event Count	62.9 (1.6)
Voxel Grid	94.3 (1.5)
EST	93.1 (1.6)
MDOE (Our)	95.8 (1.1)

test accuracy of using different methods for event-based tactile object recognition. From the table, we can see that the proposed method significantly outperforms all the baseline methods. The test accuracy achieved by MDOE is about 95.8%, which is much higher than Voxel Grid (94.3%), EST (93.1%), and TactileSGNet (89.4%). It is a bit surprising to see that Event Frame and Event Count can only achieve an accuracy of 62.9%, which is about 33% lower than our method. This might be that both Event Frame and Event Count discard the temporal information of events, which significantly affect the achieved accuracy.

## V. CONCLUSION AND DISCUSSION

We have proposed a novel representation for asynchronous event-based data which contains not only the polarity and temporal information of events, but also the event density. While the proposed method has demonstrated superior classification accuracy across different deep learning models, there are still some open issues or challenges remained. First, we integrate events into MDOE pixel-by-pixel, and the resulted representation might be computationally expensive for high-resolution event-based data. Actually, to improve the computational and storage efficiency, one can combine events from multiple pixels into a group, and generate representation group-by-group. Second, we have not evaluated the performance of the proposed method on regression tasks, such as flow estimation, simultaneous localization and mapping, and depth estimation. In the future, we will evaluate our work to these regression tasks.

#### REFERENCES

- P. Lichtsteiner, C. Posch, and T. Delbruck, "A 128×128 120dB 15μs latency 480 Q2 asynchronous temporal contrast vision sensor," *IEEE J. Solid-State Circuits*, vol. 43, no. 2, pp. 566–576, Feb. 2008.
- [2] T. Delbrück, B. Linares-Barranco, E. Culurciello, and C. Posch, "Activitydriven, event-based vision sensors," in *Proc. IEEE Int. Symp. Circuits Syst.*, 2010, pp. 2426–2429.
- [3] T. Taunyazov et al., "Event-driven visual-tactile sensing and learning for robots," in Proc. Robotics: Sci. Syst., 2020, p. 020.
- [4] D. Gehrig, A. Loquercio, K. G. Derpanis, and D. Scaramuzza, "End-to-end learning of representations for asynchronous event-based data," in *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, 2019, pp. 5633–5643.
- [5] H. Rebecq, T. Horstschaefer, and D. Scaramuzza, "Real-time visualinertial odometry for event cameras using keyframe-based nonlinear optimization," in *Proc. Brit. Mach. Vis. Conf.*, 2017, pp. 1–8.
- [6] R. Benosman, C. Clercq, X. Lagorce, S.-H. Ieng, and C. Bartolozzi, "Event-based visual flow," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 25, no. 2, pp. 407–417, Feb. 2014.
- [7] A. I. Maqueda, A. Loquercio, G. Gallego, N. García, and D. Scaramuzza, "Event-based vision meets deep learning on steering prediction for selfdriving cars," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 5419–5427.
- [8] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "Unsupervised eventbased learning of optical flow, depth, and egomotion," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 989–997.
- [9] X. Lagorce, G. Orchard, F. Galluppi, B. E. Shi, and R. B. Benosman, "HOTS: A hierarchy of event-based time-surfaces for pattern recognition," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 39, no. 7, pp. 1346–1359, Jul. 2017.
- [10] A. Sironi, M. Brambilla, N. Bourdis, X. Lagorce, and R. Benosman, "HATS: Histograms of averaged time surfaces for robust event-based object classification," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 1731–1740.
- [11] M. Cannici, M. Ciccone, A. Romanoni, and M. Matteucci, "Matrix-LSTM: A differentiable recurrent surface for asynchronous event-based data," 2020, pp. 136–152.
- [12] H. Rebecq, R. Ranftl, V. Koltun, and D. Scaramuzza, "Events-to-video: Bringing modern computer vision to event cameras," in *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 3857–3866.
- [13] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, "EV-flowNet: Selfsupervised optical flow estimation for event-based cameras," in *Proc. Robotics: Sci. Syst.*, 2018, p. 62.
- [14] F. Gu, W. Sng, T. Taunyazov, and H. Soh, "TactileSGNet: A spiking graph neural network for event-based tactile object recognition," in *Proc. IEEE/RSJ Int. Conf. Intell. Robots Syst.*, *IROS*, 2020, pp. 9876–9882.
- [15] A. Amir et al., "A low power, fully event-based gesture recognition system," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., 2017, pp. 7243–7252.
- [16] Q. Wang, Y. Zhang, J. Yuan, and Y. Lu, "Space-time event clouds for gesture recognition: From RGB cameras to event cameras," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, WACV, 2019, pp. 1826–1835.
- [17] Q. Wu, M. McGinnity, L. Maguire, A. Belatreche, and B. Glackin, "Edge detection based on spiking neural network model," in *Proc. Int. Conf. Intell. Comput.*, Springer, 2007, pp. 26–34.
- [18] J. H. Lee, T. Delbruck, and M. Pfeiffer, "Training deep spiking neural networks using backpropagation," *Front. Neurosci.*, vol. 10, 2016, Art. no. 508.

- [19] J. Botzheim, T. Obo, and N. Kubota, "Human gesture recognition for robot partners by spiking neural network and classification learning," in *Proc. IEEE 6th Int. Conf. Soft Comput. Intell. Syst.*, 13th Int. Symp. Adv. Intell. Syst., 2012, pp. 1954–1958.
- [20] P. A. Merolla *et al.*, "A million spiking-neuron integrated circuit with a scalable communication network and interface," *Science*, vol. 345, no. 6197, pp. 668–673, 2014.
- [21] M. Davies *et al.*, "Loihi: A neuromorphic manycore processor with on-chip learning," *IEEE Micro*, vol. 38, no. 1, pp. 82–99, Jan./Feb. 2018.
- [22] J. A. Pérez-Carrasco, B. Zhao, C. Serrano, B. Acha, T. S- Gotarredona, and S. Chen, "Mapping from frame-driven to frame-free event-driven vision systems by low-rate rate coding and coincidence processing–application to feedforward convNets," *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 35, no. 11, pp. 2706–2719, Nov. 2013.
- [23] S. B. Shrestha and G. Orchard, "Slayer: Spike layer error reassignment in time," in *Proc. Adv. Neural Inf. Process. Syst.*, 2018, pp. 1412–1421.
- [24] Y. Wu, L. Deng, G. Li, J. Zhu, Y. Xie, and L. Shi, "Direct training for spiking neural networks: Faster, larger, better," in *Proc. AAAI Conf. Artif. Intell.*, 2019, vol. 33, pp. 1311–1318.
- [25] D. Neil, M. Pfeiffer, and S.-C. Liu, "Phased LSTM: Accelerating recurrent network training for long or event-based sequences," in *Proc. Adv. neural Inf. Process. Syst.*, 2016, pp. 3882–3890.
- [26] M. Cannici, M. Ciccone, A. Romanoni, and M. Matteucci, "Attention mechanisms for object recognition with event-based cameras," in *Proc. IEEE Winter Conf. Appl. Comput. Vis.*, WACV, 2019, pp. 1127–1136.
- [27] N. Messikommer, D. Gehrig, A. Loquercio, and D. Scaramuzza, "Eventbased asynchronous sparse convolutional networks," in *Proc. Eur. Conf. Comput. Vis.*, Springer, 2020, pp. 415–431.
- [28] K. Roy, A. Jaiswal, and P. Panda, "Towards spike-based machine intelligence with neuromorphic computing," *Nature*, vol. 575, no. 7784, pp. 607–617, 2019.
- [29] G. Orchard, A. Jayawant, G. K. Cohen, and N. Thakor, "Converting static image datasets to spiking neuromorphic datasets using saccades," *Frontiers Neurosci.*, vol. 9, 2015, Art. no. 437.
- [30] L. Fei-Fei, R. Fergus, and P. Perona, "Learning generative visual models from few training examples: An incremental Bayesian approach tested on 101 object categories," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshop*, 2004, pp. 178–178.
- [31] W. W. Lee *et al.*, "A neuro-inspired artificial peripheral nervous system for scalable electronic skins," *Sci. Robot.*, vol. 4, no. 32, 2019, Art. no eaax2198.
- [32] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [33] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Int. Conf. Learn. Representations*, 2015.
- [34] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "MobileNetV2: Inverted residuals and linear bottlenecks," in *Proc. IEEE Conf. Comput. Vis. pattern Recognit.*, 2018, pp. 4510–4520.
- [35] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, "Rethinking the inception architecture for computer vision," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 2818–2826.
- [36] O. Russakovsky et al., "Imagenet large scale visual recognition challenge," Int. J. Comput. Vis., vol. 115, no. 3, pp. 211–252, 2015.
- [37] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," 2014, arXiv:1412.6980.
- [38] L. Prechelt, "Early stopping-but when?," in *Neural Networks: Tricks of the Trade*. Berlin, Germany: Springer, 1998, pp. 55–69.