

HARDVS: Revisiting Human Activity Recognition with Dynamic Vision Sensors

Xiao Wang¹, Zongzhen Wu¹, Bo Jiang^{1*}, Zhimin Bao², Lin Zhu³
Guoqi Li⁴, Yaowei Wang⁵, Yonghong Tian^{5,6}

¹School of Computer Science and Technology, Anhui University, Hefei, China

²Tencent YouTu Lab, Hefei, China ³Beijing Institute of Technology, Beijing, China

⁴Institute of Automation, Chinese Academy of Sciences, Beijing, China

⁵Peng Cheng Laboratory, Shenzhen, China ⁶Peking University, Beijing, China

Abstract

The main streams of human activity recognition (HAR) algorithms are developed based on RGB cameras which are suffered from illumination, fast motion, privacy-preserving, and large energy consumption. Meanwhile, the biologically inspired event cameras attracted great interest due to their unique features, such as high dynamic range, dense temporal but sparse spatial resolution, low latency, low power, etc. As it is a newly arising sensor, even there is no realistic large-scale dataset for HAR. Considering its great practical value, in this paper, we propose a large-scale benchmark dataset to bridge this gap, termed HARDVS, which contains 300 categories and more than 100K event sequences. We evaluate and report the performance of multiple popular HAR algorithms, which provide extensive baselines for future works to compare. More importantly, we propose a novel spatial-temporal feature learning and fusion framework, termed ESTF, for event stream based human activity recognition. It first projects the event streams into spatial and temporal embeddings using StemNet, then, encodes and fuses the dual-view representations using Transformer networks. Finally, the dual features are concatenated and fed into a classification head for activity prediction. Extensive experiments on multiple datasets fully validated the effectiveness of our model. Both the dataset and source code will be released on <https://github.com/Event-AHU/HARDVS>.

1. Introduction

With the rapid development of the smart city, recognizing human behavior (i.e., Human Activity Recognition, HAR) accurately and efficiently is becoming an extremely urgent task. Most researchers develop the HAR algorithms [1, 29] based on RGB cameras which are widely de-

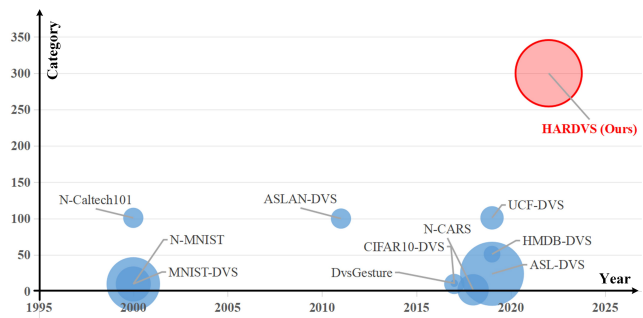


Figure 1. Comparison between existing datasets and our proposed HARDVS dataset for event based video classification.

ployed and easy to collect the data. With the help of large-scale benchmark datasets [7, 24, 26, 30, 43, 50, 53] and deep learning, HAR in regular scenarios has been studied to some extent. However, the storage, transmission, and analysis of surveillance videos set limits the demands for the practical systems due to the usage of RGB sensors. More in detail, the standard RGB cameras have a limited frame rate (e.g., 30 FPS) which makes it hard to capture the fast-moving objects and is easily influenced by motion blur. The low dynamic range (60 dB) makes the RGB sensors work poorly in low illumination. It also suffers from the high redundancy between nearby frames which needs more storage and energy consumption. Privacy protection also greatly limits its development, therefore, a natural question is *do we have to recognize human activities using the RGB sensors?*

Recently, the biologically inspired sensors (called event cameras), such as DAVIS [6], CeleX [13], ATIS [47], and PROPHESSEE¹, drawing more and more attention of researchers. Different from RGB cameras which record light in a synchronous way (i.e., the video frame), the event cameras output events (or spikes) asynchronously which corresponds to the illumination variation. In another word, each pixel of event cameras independently records a binary value

*Corresponding author: Bo Jiang

¹<https://www.prophesee.ai>

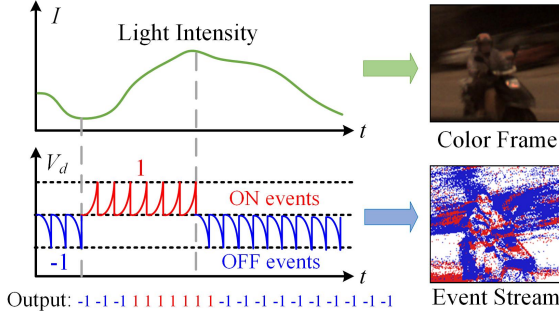


Figure 2. Comparison between the imaging principles of the color frame and event stream.

only when the light changes exceed a threshold. Events for the increase and decrease of illumination are called ON and OFF events respectively. Due to the unique sampling mechanism, the asynchronous events are spatially sparse but temporally dense. It is less affected by motion blur, therefore, is suitable for capturing fast-moving human actions, such as the magician’s fast-moving palm, and movement recognition of sports players. It has a higher dynamic range (120 dB) and lower latency, which enables it to work well even in low illumination compared with standard RGB cameras. In addition, the storage and energy consumption are also significantly reduced [21, 34, 58, 67, 68]. Event streams highlight the contour information and protect personal privacy to a large extent. According to the aforementioned observation and thinking, we are inspired to address human activity recognition in the wild using event cameras. A comparison of the imaging principles of the color frame and event camera is illustrated in Fig. 2.

Although there are already several benchmark datasets proposed for classification [2, 5, 28, 30, 33, 49, 53]. However, most of them are simulated/synthetic datasets that are transformed from RGB videos with the simulator. Some researchers attain the event data by recording the screen while displaying RGB videos. Obviously, these datasets are hard to reflect the features of event cameras in real-world scenarios, especially fast-motion and low-light scenarios. ASL-DVS [5] is proposed by Bi et al. which is consisted of 100, 800 samples but can only be used for hand gesture recognition with 24 classes. DvsGesture [2] is also limited by its scale and categories in the deep learning era. In addition, some datasets have become saturated in performance, for example, Wang et al. [57] already achieved 97.08% on the DvsGesture [2] dataset. Therefore, the research community still has insistent demands for a large-scale HAR benchmark dataset recorded in the wild.

In this paper, we propose a large-scale benchmark dataset, termed HARDVS, to address the problem of the lack of real event data. Specifically, the HARDVS dataset contains more than 100K video clips recorded with a

DAVIS346 camera, each of them lasting for about 5-10 seconds. It contains 300 categories of human activities in daily life, such as *drinking, riding a bike, sitting down, washing hands*. The following factors are taken into account to make our data more diverse, including *multi-views, illuminations, motion speed, dynamic background, occlusion, flashing light, photographic distance*. To the best of our knowledge, our proposed HARDVS is the first large-scale and challenging benchmark dataset for human activity recognition in the wild. A comparison between existing recognition datasets and our HARDVS is illustrated in Fig. 1.

Based on our newly proposed HARDVS dataset, we construct a novel event-based human action recognition framework, termed ESTF (Event-based Spatial-Temporal Transformer). As shown in Fig. 4, the ESTF transforms the event streams into spatial and temporal tokens and learns the dual features by employing SpatialFormer (SF) and TemporalFormer (TF) respectively. Further, we propose a FusionFormer to realize the message passing between the spatial and temporal features. The aggregated representation is added with features of dual branches as the input for subsequent learning blocks, respectively. The outputs will be concatenated and fed into two MLP layers for the final action prediction.

To sum up, the contributions of this paper can be concluded as the following three aspects:

- We propose a large-scale neuromorphic dataset for human activity recognition, termed HARDVS. It contains more than 100K samples with 300 classes, and fully reflects the challenging factors in the real world. To the best of our knowledge, it is the first large-scale realistic neuromorphic dataset for HAR.
- We propose a novel Event-based Spatial-Temporal Transformer (ESTF) approach for human action recognition by exploiting spatial and temporal feature learning and fusing them with Transformer networks. It is the first Transformer based spatial-temporal representation learning framework for event stream-based HAR.
- We re-train and report the performance of multiple popular HAR algorithms, which provide extensive baselines for future works to compare on the HARDVS dataset. Extensive experiments on multiple event-based classification datasets fully demonstrate the effectiveness of our proposed ESTF approach.

2. Related Work

HAR with Event Sensors. Compared with RGB cameras, few researchers focus on event camera-based HAR [2, 3, 10, 14]. Arnon et al. [2] propose the first gesture recognition system based on TrueNorth neurosynaptic processor. Xavier et al. [14] propose an event-based luminance-free feature for local corner detection and global gesture recog-

nition. Chen et al. [10] propose a hand gesture recognition system based on DVS and also design a wearable glove with a high-frequency active LED marker that fully exploits its properties. A retinomorph event-driven representation (EDR) is proposed by Chen et al. [11], which can realize three important functions of the biological retina, i.e., the logarithmic transformation, ON/OFF pathways, and integration of multiple timescales. The authors of [31] represent the recent temporal activity within a local spatial neighborhood, and utilize the rich temporal information provided by events to create contexts in the form of time-surfaces, termed HOTS, for the recognition task. Wu et al. first transform the event flow into images, then, predict and combine the human pose with event images for HAR [62]. Graph neural networks (GNN) and SNNs are also exploited for event-based recognition [9, 12, 23, 32, 37, 41, 45, 48, 60, 64]. Specifically, Chen et al. [12] treat the event flow as a 3D point cloud and use dynamic GNNs to learn the spatial-temporal features for gesture recognition. Wang et al. [60] adopt GNNs and CNNs for gait recognition. Xing et al. design a spiking convolutional recurrent neural network (SCRNN) architecture for event-based sequence [64]. According to our observations, these works are evaluated only on simple HAR datasets or simulated datasets. It is necessary and urgent to introduce a large-scale HAR dataset for current evaluation.

Event Benchmark Datasets for HAR. As shown in Table 1, most of the existing event camera-based datasets for recognition are artificial datasets. Usually, the researchers display the RGB HAR datasets on a large screen and record the activity with neuromorphic sensors. For example, the N-Caltech101 [44] and N-MNIST [44] are recorded with an ATIS camera which contains 101 and 10 classes, respectively. Bi et al. [5] also transform popular HAR datasets into simulated event flow, including HMDB-DVS [5, 30], UCF-DVS [5, 53], and ASLAN-DVS [28], which further expands the number of datasets available for HAR. However, these simulated event datasets hardly reflect the advantages of event cameras, such as low light, fast motion, etc. There are three realistic event datasets for classification, i.e., the DvsGesture [2], N-CARS [51] and ASL-DVS [5], but these benchmarks are limited by their scale, categories, and scenes. Specifically, these datasets contain 11, 2, and 24 classes only, and also rarely take challenging factors like multi-view, motion, and glitter into consideration. Compared with existing datasets, our proposed HARDVS dataset is large-scale (100K samples) and category-wide (300 classes) for deep neural networks. Our sequences are recorded in the wild and fully reflect the features of the aforementioned attributes. We believe our proposed benchmark dataset greatly promotes the development of event-based HAR.

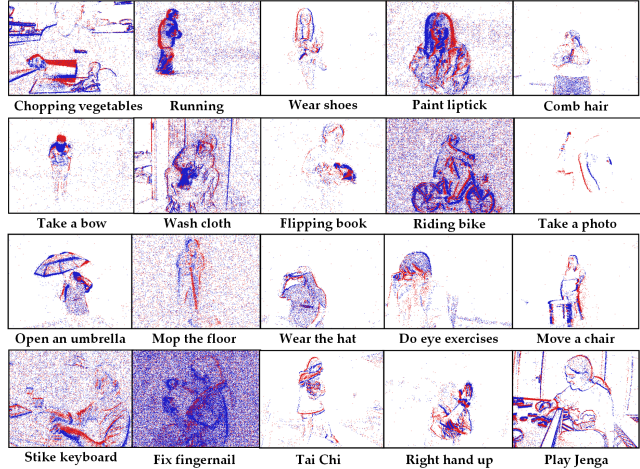


Figure 3. Illustration of some representative samples of our proposed HARDVS dataset.

3. HARDVS Benchmark Dataset

3.1. Protocols

We aim to provide a good platform for the training and evaluation of DVS-based human activity recognition. When constructing the HARDVS benchmark dataset, we obey the following protocols:

- 1). Large-scale:** As we all know, large-scale datasets play a very important role in the deep learning era. In this work, we collect more than 100k DVS event sequences to meet the needs for large-scale training and evaluation of HAR.
- 2). Wide varieties:** Thousands of human activities can exist in the real world, but existing DVS-based HAR datasets only contain limited categories. Therefore, it is hard to fully reflect the classification and recognition ability of HAR algorithms. Our newly proposed HARDVS contains 300 classes which are several times larger than other DVS datasets.
- 3). Various challenges:** Our dataset considers multiple challenging factors which may influence the results of HAR with the DVS sensor. The detailed introductions can be found below:
 - (a). Multi-view:** We collect different views of the same behavior to mimic practical applications, including front-, side-, horizontal-, top-down-, and bottom-up-views.
 - (b). Multi-illumination:** High dynamic range is one of the most important features of DVS sensors, therefore, we collect the videos under scenarios with strong-, middle-, and low-light (60% of each category). Our dataset also contains many videos with *glitter*, because we find that the DVS sensor is sensitive to flashing lights, especially at the night.
 - (c). Multi-motion:** We also highlight the features of DVS sensors by recording many actions with various motion speeds, such as slow-, moderate-, and high-speed.
 - (d). Dynamic background:** As it is relatively easy to recognize actions without background objects, i.e., sta-

Table 1. **Comparison of event datasets for human activity recognition.** M-VW, M-ILL, M-MO, DYB, OCC, and DR denotes multi-view, multi-illumination, multi-motion, dynamic background, occlusion, and duration of the action, respectively. Note that we only report these attributes of realistic DVS datasets for HAR.

Dataset	Year	Sensors	Scale	Class	Resolution	Real	M-VW	M-ILL	M-MO	DYB	OCC	DR	Link
ASLAN-DVS [5, 28]	2011	DAVIS240c	3, 697	432	240 × 180	✗	-	-	-	-	-	-	URL
MNIST-DVS [49]	2013	DAVIS128	30, 000	10	128 × 128	✗	-	-	-	-	-	-	URL
N-Caltech101 [44]	2015	ATIS	8, 709	101	302 × 245	✗	-	-	-	-	-	-	URL
N-MNIST [44]	2015	ATIS	70, 000	10	28 × 28	✗	-	-	-	-	-	-	URL
CIFAR10-DVS [33]	2017	DAVIS128	10, 000	10	128 × 128	✗	-	-	-	-	-	-	URL
HMDB-DVS [5, 30]	2019	DAVIS240c	6, 766	51	240 × 180	✗	-	-	-	-	-	-	URL
UCF-DVS [5, 53]	2019	DAVIS240c	13, 320	101	240 × 180	✗	-	-	-	-	-	-	URL
N-ImageNet [27]	2021	Samsung-Gen3	1, 781, 167	1000	480 × 640	✗	-	-	-	-	-	-	URL
ES-ImageNet [36]	2021	-	1, 306, 916	1000	224 × 224	✗	-	-	-	-	-	-	URL
DvsGesture [2]	2017	DAVIS128	1, 342	11	128 × 128	✓	✗	✓	✗	✗	✗	-	URL
N-CARS [51]	2018	ATIS	24, 029	2	304 × 240	✓	✗	✗	✗	✗	✗	-	URL
ASL-DVS [5]	2019	DAVIS240	100, 800	24	240 × 180	✓	✗	✗	✗	✗	✗	0.1s	URL
PAF [42]	2019	DAVIS346	450	10	346 × 260	✓	✗	✗	✗	✗	✗	5s	URL
DailyAction [38]	2021	DAVIS346	1, 440	12	346 × 260	✓	✓	✓	✗	✗	✗	5s	URL
HARDVS (Ours)	2022	DAVIS346	107, 646	300	346 × 260	✓	✓	✓	✓	✓	✓	5s	URL

tionary DVS camera, we also collect many actions with a dynamic moving camera to make our dataset challenging enough. (e). *Occlusion*: In the real world, human action can be occluded commonly. Thus, we also add occlusion issues into the HARDVS dataset with hand or other things. **4). Different capture distance**: The HARDVS dataset is collected under various distances, i.e., 1-2, 3-4, and more than 5 meters. **5). Long-term**: Most of the existing DVS-based HAR datasets are microsecond-level, in contrast, each video in our HARDVS dataset lasts for about 5 seconds. **6). Dual-modality**: The DAVIS346 camera can output both RGB frames and event flow, therefore, our dataset can also be used for HAR by fusing video frames and events. In this work, we focus on HAR with DVS only, but the RGB frames will also be released to support the research on dual-modality fusing based HAR.

3.2. Data Collection and Statistic Analysis

The HARDVS dataset is collected with a DAVIS346 camera whose resolution is 346×260 . We take the aforementioned protocols in mind when recording videos. Therefore, our dataset fully reflects the unique features of DVS sensors in challenging scenarios, such as low-illumination, high-speed, clutter background, etc. The main characters are also diverse, generally speaking, there is a total of five persons involved in the data collection stage.

From a statistical perspective, our dataset contains a total of 107, 646 video sequences and 300 classes of common human activities. We split 60%, 10%, and 30% of each category for training, validating, and testing, respectively. Totally, the number of videos in the training, validating, and testing subset is 64526|10734|32386, respectively. A direct comparison with existing classification benchmark datasets can be found in Table 1 and Fig. 1. With the aforementioned characteristics, we believe our HARDVS dataset will be a better evaluation platform for the neuromorphic classifica-

tion problem, especially for the human activity recognition task.

4. Methodology

4.1. Overview

In this section, we devise a new Event-based Spatial-Temporal Transformer (ESTF) approach for event-stream data learning. As shown in Fig. 4, the proposed ESTF architecture contains three main learning modules, i.e., i) Initial Spatial and Temporal Embedding, ii) Spatial and Temporal Enhancement Learning, and iii) Spatial-Temporal Fusion Transformer. Specifically, given the input event-stream data, we first extract the initial spatial and temporal embeddings respectively. Then, a Spatial and Temporal Feature Enhancement Learning module is devised to further enrich the event-stream data representations by deeply capturing both **spatial correlation** and **temporal dependence** of event stream. Finally, an effective Fusion Transformer (FusionFormer) block is designed to integrate the spatial and temporal cues together for the final feature representation. The details of these modules are introduced below.

4.2. Initial Spatial and Temporal Embedding

Different from visible sensors which capture a global image at each time, the event cameras asynchronously capture the intensity variations in the log-scale. That is, each pixel outputs a discrete event (or spike) independently when the visual changing exceeds a pre-defined threshold. Usually, we use a 4-tuple $\{x, y, t, p\}$ to represent the discrete event of a pixel captured with DVS, where x, y are spatial coordinates, t is timestamp, and $p \in \{1, -1\}$ is the polarity of brightness variation. Following previous works [18, 59, 65, 66], we first transform the asynchronous event flow into the synchronous *event images* by stacking the events in a time interval based on the exposure time. Let

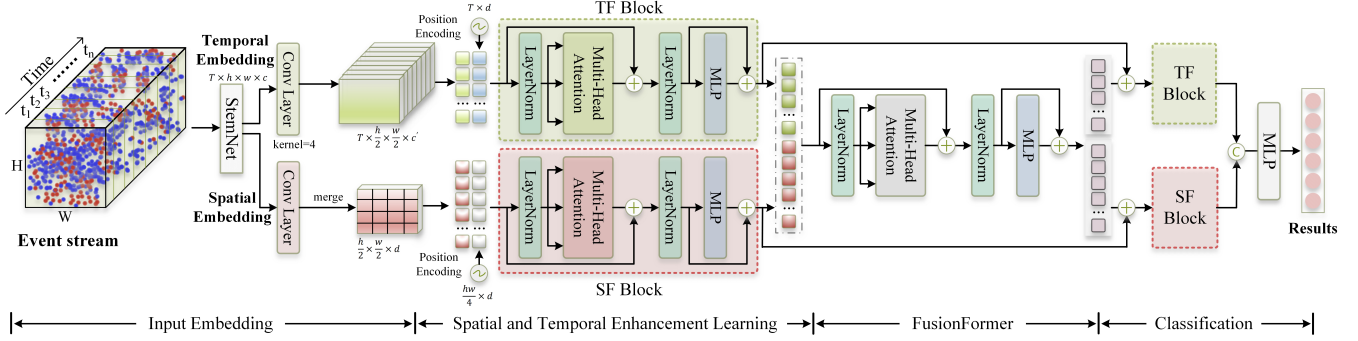


Figure 4. An overview of our proposed ESTF framework for event-based human action recognition. It transforms the event streams into spatial and temporal tokens and learns the dual features using multi-head self-attention layers. Further, a FusionFormer is proposed to realize message passing between the spatial and temporal features. The aggregated features are added with dual features as the input for subsequent TF and SF blocks, respectively. The outputs will be concatenated and fed into MLP layers for action prediction.

$\mathcal{E} = \{E_1, E_2 \dots E_T\} \in \mathbb{R}^{H \times W \times T}$ be the collection of the sampled input event frames. In our experiments, we set $T = 8$, as used in works [55]. For each event frame E_t , we adopt StemNet [25] to extract an initial CNN feature descriptor for it and denote $\mathcal{X} = \{\mathcal{X}_1, \mathcal{X}_2 \dots \mathcal{X}_T\} \in \mathbb{R}^{T \times h \times w \times c}$ as the collection of T event frames. Based on it, we respectively extract spatial and temporal embeddings. To be specific, for temporal branch, we adopt a convolution layer to reduce the feature size to obtain $\mathcal{X}^t \in \mathbb{R}^{T \times \frac{h}{2} \times \frac{w}{2} \times c'}$ and reshape it to the matrix form as $X^t \in \mathbb{R}^{T \times d}$ where $d = \frac{h}{2} \times \frac{w}{2} \times c'$. For spatial branch, we first adopt a convolution layer to resize the features \mathcal{X} to $\mathcal{X}^s \in \mathbb{R}^{T \times \frac{h}{2} \times \frac{w}{2} \times d}$. Then, we conduct the merging/summation operation on the time dimension and reshape it to the matrix form $X^s \in \mathbb{R}^{N \times d}$ where $N = \frac{hw}{4}$. Hence, both spatial and temporal embeddings have the same d -dim feature descriptors.

4.3. Spatial and Temporal Enhancement Learning

Based on the above initial spatial embeddings $X^s \in \mathbb{R}^{N \times d}$ and temporal embeddings $X^t \in \mathbb{R}^{T \times d}$, we then devise our Spatial and Temporal Enhancement Learning (STEL) module to further enrich their representations. The proposed STEL module involves two blocks, i.e., Spatial Transformer (SF) block, and Temporal Transformer (TF) block, which respectively capture the spatial correlations and temporal dependences of event data to learn context enriched representations. The SF block includes multi-head self-attention (MSA) and MLP module with a LayerNorm (LN) used between two modules. A residual connection is also employed, as shown in Fig. 4. To be specific, given spatial embeddings $X^s \in \mathbb{R}^{N \times d}$, we first incorporate the position encoding [17] to obtain $\tilde{X}^s \in \mathbb{R}^{N \times d}$ which represents N number of the input tokens with d -dim feature descriptor. Then, the outputs of SF block are summarized

as follows,

$$Y^s = LN(\tilde{X}^s + MSA(LN(\tilde{X}^s))) \quad (1)$$

$$\tilde{X}^s = Y^s + MLP(Y^s) \quad (2)$$

In contrast to input \tilde{X}^s , the output \tilde{X}^s provides the spatial-aware enhanced representations by employing the MSA mechanism to model the spatial relationships of different event patches. Similarly, given $\tilde{X}^t \in \mathbb{R}^{T \times d}$ representing T temporal tokens with position encoding, the outputs of TF block are summarized as follows,

$$Y^t = LN(\tilde{X}^t + MSA(LN(\tilde{X}^t))) \quad (3)$$

$$\tilde{X}^t = Y^t + MLP(Y^t) \quad (4)$$

Compared with the input \tilde{X}^t , the outputs $\tilde{X}^t \in \mathbb{R}^{T \times d}$ provide a temporal-context enhanced representations for T number of frame tokens thanks to the MSA mechanism to model the dependencies of different event frames.

4.4. Fusion Transformer

In order to conduct the interaction between the above ST and TF blocks and learn a unified spatio-temporal contextual data representations, we also design a Fusion Transformer (FusionF) module. To be specific, let \tilde{X}^s and \tilde{X}^t denote the outputs of previous SF and TF blocks respectively. We first collect the N spatial and T temporal tokens together and feed them to a unified Transformer block which includes multi-head self-attention (MSA) and MLP submodule, i.e.,

$$Z = [\tilde{X}^t, \tilde{X}^s] \in \mathbb{R}^{(T+N) \times c} \quad (5)$$

$$Y = LN(Z + MSA(LN(Z))) \quad (6)$$

$$\tilde{Z} = Z + Y + MLP(Y) \quad (7)$$

Afterward, we split \tilde{Z} into $\{\tilde{Z}^s, \tilde{Z}^t\}$ where $\tilde{Z}^s \in \mathbb{R}^{N \times d}$ and $\tilde{Z}^t \in \mathbb{R}^{T \times d}$ and further employ the above SF (Eqs.(1,2)) and TF (Eqs.(3,4)) block to respectively enhance their representations as follows,

$$F^s = SF(\tilde{Z}^s), F^t = TF(\tilde{Z}^t) \quad (8)$$

Finally, we concatenate both F^s and F^t together and reshape the concatenated features to the vector form. After that, we utilize a two-layer MLP to output the final class label prediction, as shown in Fig. 4.

4.5. Loss Function

Our proposed ESTF framework can be optimized in an end-to-end way. The standard cross-entropy loss function is adopted to measure the distance between our model prediction and ground truth:

$$Loss = -\frac{1}{B} \sum_{b=1}^B \sum_{n=1}^N Y_{bn} \log P_{bn} \quad (9)$$

where B denotes the batch size, N denotes the number of event classes. Y and P represent the ground truth and predicted class labels of the event sample, respectively.

5. Experiments

5.1. Dataset and Evaluation Metrics

In this work, three datasets are adopted for the evaluation of our proposed model, including **N-Caltech101** [44], **ASL-DVS** [5], and our newly proposed **HARDVS**. More details about these datasets can be found in Table 1. The widely used **top-1** and **top-5 accuracy** are adopted as evaluation metrics.

5.2. Implementation Details

Given the event streams, we stack them into image-like representations to make full use of CNN. More in detail, the time window is set based on the exposure time of color frames, when generating the event images. The batch size is 60, and the initial learning rate is 0.01, which is reduced to 10% of the original every 15 epochs. The stochastic gradient descent (SGD) [54] is selected as the optimizer to train our network. Our code is implemented based on Python 3.8, PyTorch 1.10.2+cu113 [46], on a server with RTX3090. The source code and pre-trained models will be released to help other researchers reproduce our experimental results.

5.3. Comparison with SOTA Algorithms

Results on N-Caltech101 [44]. As shown in Table 2, our proposed method achieves 0.832 on the top-1 accuracy metric which is significantly better than the compared models by a large margin. For example, the VMV-GCN

Table 2. Results on N-Caltech101 [44] Dataset.

EventNet	Gabor-SNN	RG-CNNs	VMV-GCN	EV-VGCNN	EST
0.425	0.196	0.657	0.778	0.748	0.753
ResNet-50	MVF-Net	M-LSTM	AMAE	HATS	Ours
0.637	0.687	0.738	0.694	0.642	0.832

Table 3. Results on the ASL-DVS [5] dataset.

EST	AMAE	M-LSTM	MVF-Net	ResNet-50
0.979	0.984	0.980	0.971	0.886
EventNet	RG-CNNs	EV-VGCNN	VMV-GCN	Ours
0.833	0.901	0.983	0.989	0.999

achieves 0.778 on this benchmark dataset which ranks second place, meanwhile, our model outperforms it by up to 5.4%. The M-LSTM is an adaptive event representation learning model which obtained 0.738 only on this dataset. EV-VGCNN is a graph neural network based model which obtains 0.748 and is also worse than ours. These experimental results fully demonstrate the effectiveness of our proposed spatial-temporal feature learning for event-based pattern recognition.

Results on ASL-DVS [5]. As shown in Table 3, the performance on this dataset is already close to saturation and most of the compared models achieve more than 0.95+ on the top-1 accuracy, including EST [22] (0.979), AMAE [16] (0.984), M-LSTM [8] (0.980), MVF-Net [15] (0.971). Note that, the VMV-GCN [63] achieves 0.989 on this benchmark dataset which ranks the second place. It is very hard to beat these models. Thanks to our proposed spatial-temporal feature learning and fusion modules, we set new state-of-the-art performance on this dataset, i.e., 0.999 on the top-1 accuracy. Therefore, we can draw the conclusion that our method almost completely solves the simple gesture recognition problem defined in the ASL-DVS [5].

Results on HARDVS. From the experimental results reported in the ASL-DVS [5] and N-Caltech101 [44], we can find that existing event based recognition datasets are almost saturated. The newly proposed HARDVS dataset can bridge this gap and further boost the development of event based human action recognition. As shown in Table 4, we re-training and testing multiple state-of-the-art models for future works to compare on the HARDVS benchmark dataset, including C3D [55], R2Plus1D [56], TSM [52], ACTION-Net [61], TAM [40], Video-SwinTrans [39], TimeSformer [4], SlowFast [20]. It is easy to find that these popular and strong recognition models still perform poorly on our newly proposed HARDVS dataset. To be specific, the R2Plus1D [56], ACTION-Net [61], and SlowFast [20] only achieves 49.06|56.43, 46.85|56.19, and 46.54|54.76 on the top-1 and top-5 accuracy respectively. The recently proposed TAM [40] (ICCV-2021), Video-SwinTrans [39] (CVPR-2022), TimeSformer [4] (ICML 2021) also obtains

Table 4. Results on the newly proposed HARDVS dataset.

No.	Algorithm	Publish	Backbone	Top1	Top5
01	ResNet18 [25]	CVPR-2016	ResNet18	49.20	56.09
02	C3D [55]	ICCV-2015	CNN	50.52	56.14
03	R2Plus1D [56]	CVPR-2018	ResNet-34	49.06	56.43
04	TSM [35]	ICCV-2019	ResNet-50	52.63	60.56
05	ACTION-Net [61]	CVPR-2021	ResNet-50	46.85	56.19
06	TAM [40]	ICCV-2021	ResNet-50	50.41	57.99
07	Video-SwinTrans [39]	CVPR-2022	Swin Transformer	51.91	59.11
08	TimeSformer [4]	ICML-2021	VIT	50.77	58.70
09	SlowFast [20]	ICCV-2019	ResNet-50	46.54	54.76
10	X3D [19]	CVPR-2020	ResNet	45.82	52.33
11	ESTF (Ours)	-	ResNet18	51.22	57.53

50.41|57.99, 51.91|59.11, and 50.77|58.70 on the two metrics respectively. Compared with these models, our proposed spatial-temporal feature learning and fusion modules perform comparable or even better than these SOTA models, i.e., 51.22|57.53. All in all, our proposed model is effective for event based human action recognition task and may be a good baseline for future works to compare.

5.4. Ablation Study

To help researchers better understand our proposed module, in this subsection, we conduct extensive experiments to analyze the contributions of each key component and the influence of different settings for our model.

Component Analysis. As shown in Table 5, three main modules are analyzed on the N-Caltech101 dataset, including SpatialFormer (SF), TemporalFormer (TF), and FusionFormer. We can find that our baseline method ResNet18 [25] achieves 72.14 on the top-1 accuracy metric. When introducing the TemporalFormer (TF) into the recognition framework, the overall performance can be significantly improved by +9.4, and achieves 81.54. When the SpatialFormer (SF) is adopted for long-range global feature relation mining, the recognition results can be enhanced to 80.47, and the improvement is up to +8.33. When both modules are all utilized for joint spatial-temporal feature learning, a better result can be obtained, i.e., 82.89. If the FusionFormer is adopted to achieve interactive feature learning and information propagation between the spatial and temporal Transformer branches, the best results can be achieved, i.e., 83.17 on the top-1 accuracy. Based on the experimental analysis for Table 5 and Table 2, we can draw the conclusion that our proposed modules all contribute to final recognition results.

Analysis on Number of Input Frames. In this paper, we transform the event streams into an image-like representation for classification. In our experiments, 8 frames are adopted for the evaluation of our model. Actually, various event frames can be obtained with different intervals of the time windows. In this part, we test our model with 4, 6, 8, 10, 12, and 16 frames on the N-Caltech101 dataset and report the results in Fig. 5. It is easy to find that the mean

Table 5. Component Analysis on the N-Caltech101 Dataset.

No.	ResNet	TF	SF	FusionFormer	Accuracy
1	✓				72.14
2	✓	✓			81.54
3	✓		✓		80.47
4	✓	✓	✓		82.89
5	✓	✓	✓	✓	83.17

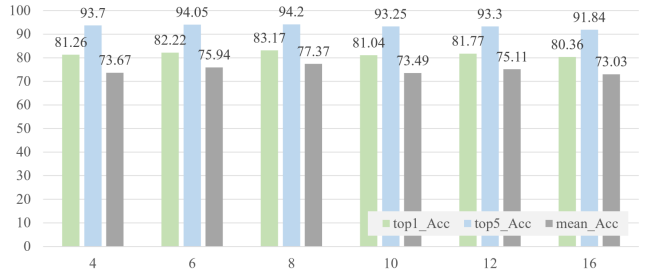


Figure 5. Experimental results of different input frames.

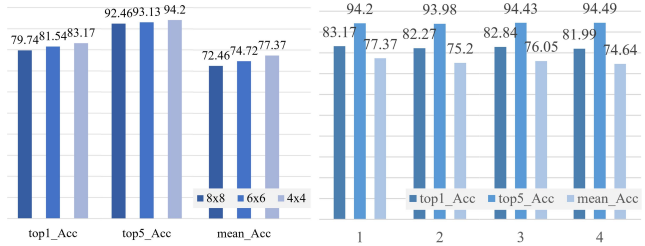


Figure 6. Experimental results of different partition patches (left) and Transformer layers (right).

accuracy is 73.67, 75.94, 77.37, 73.49, 75.11, and 73.03, correspondingly, and the highest mean accuracy can be obtained when 8 frames are adopted. For the decrease in accuracy when the frames are larger than 8, we think this may be caused by the fact that the event streams are partitioned into more frames and each frame will be more sparse. Therefore, this will lead to sparse edge information which is very important for recognition.

Analysis on Split Patches of Spatial Data. In this paper, the spatial features are partitioned into non-overlapped patches. We test multiple scales in this subsection, including 8×8 , 6×6 , and 4×4 . As illustrated in Fig. 6 (left), the best performance can be obtained when 4×4 is adopted, i.e., 83.17, 94.20, and 77.37 on the top-1, top-5, and mean accuracy respectively.

Analysis on Layers of Transformer Layers. As we all know, the self-attention or Transformer layers can be stacked multiple times for more accurate recognition, as validated in many works. In this experiment, we also test different Transformer layers to check their influence on our model. As shown in Fig. 6 (right), four different settings are

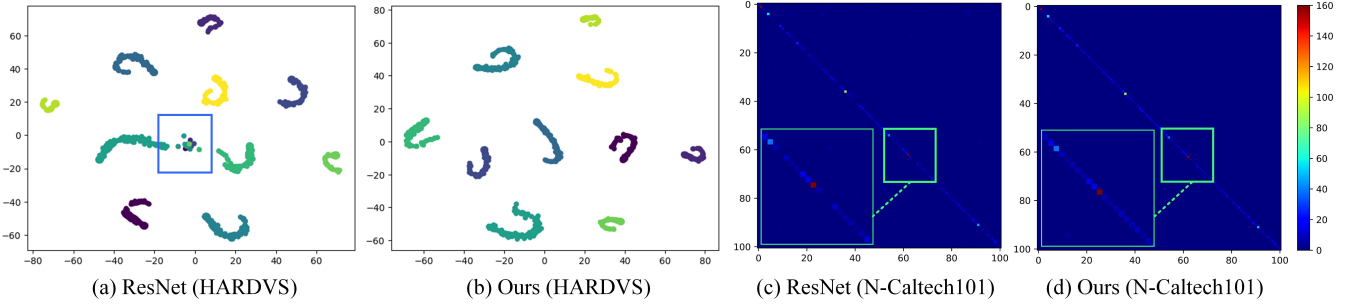


Figure 7. Visualization of feature distribution of our baseline and newly proposed ESTF on HARDVS dataset (a, b) and confusion matrix of baseline ResNet and our model on N-Caltech101 dataset (c, d). Best viewed by zooming in.

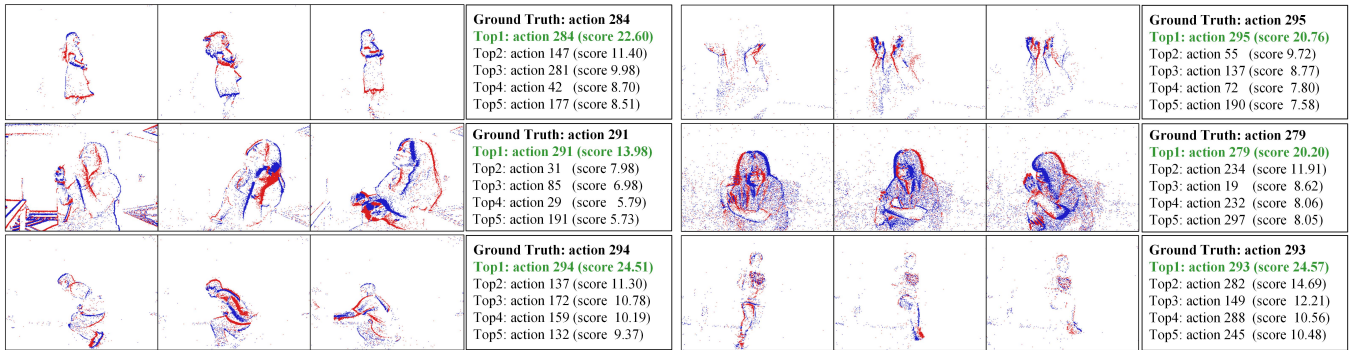


Figure 8. Visualization of the top-5 predicted actions using our model.

tested, i.e., 1, 2, 3, and 4 layers, and the corresponding mean accuracy is 77.37, 75.20, 76.05, and 74.64. We can find that higher recognition results can be obtained when the Transformer is set as 1 to 3 layers. Maybe a larger dataset is needed to train deeper Transformer layers.

Model Parameters and Running Efficiency. The storage space occupied by our checkpoint is 377.34 MB and the number of parameters is 46.71 M. The MAC score is 17.62 G tested using toolkit *ptflops*². Our model spends 25 ms for each video (8 frames used) in our proposed HARDVS dataset.

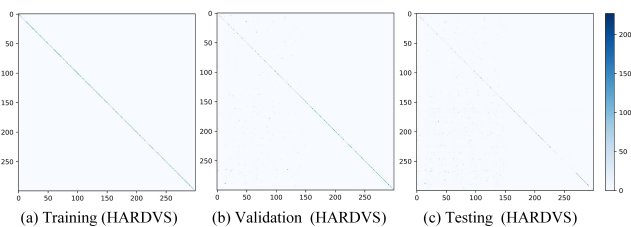


Figure 9. Visualization of confusion matrix on the HARDVS dataset.

5.5. Visualization

In the previous subsections, we conduct extensive experiments to validate the effectiveness of our model from a quantitative point of view. In this part, we resort to the visualization to help the readers better understand our proposed model.

Feature Visualization & Confusion Matrix. As shown in Fig. 7 (a, b), we select 10 classes of actions defined in the HARDVS dataset and visualize the features by projecting them into 2D plane using *tSNE* toolkit³. It is easy to find that partial data samples are not discriminated well using the baseline ResNet18, such as the regions highlighted in blue bounding box. In contrast, our proposed ESTF model achieves a better feature representation learning and more of the categories are classified well. For the confusion matrix on N-Caltech101 dataset, as shown in Fig. 7 (c, d), we can find that our proposed ESTF achieves significant improvement compared with our baseline ResNet18. All in all, we can draw the conclusion that our proposed spatial-temporal feature learning module works well for event based action recognition.

Confusion Matrix. As shown in Fig. 9, we visualize the confusion matrix of our model based on the results pre-

²<https://pypi.org/project/ptflops/>

³<https://github.com/mx11990/tsne-pytorch>

dicted in the training, validation, and testing phase, respectively. One can note that our model achieves better results in the training phase, but the overall performance in the testing phase is still weak. This demonstrates that our proposed HARDVS dataset is challenging and there is still plenty of room for further improvement.

Recognition Results. As shown in Fig. 8, we provide the top-5 predicted actions and corresponding confidence scores. The ground truth and top-1 results are highlighted in **black** and **green**. It is easy to find that our model can predict the human activities accurately.

6. Conclusion

In this paper, we propose a large-scale benchmark dataset for event-based human action recognition, termed HARDVS. It contains 300 categories of human activities and more than 100K event sequences captured from DAVIS346 camera. These videos reflect various views, illuminations, motions, dynamic backgrounds, occlusion, etc. More than 10 popular and recent classification models are evaluated for future works to compare. In addition, we also propose a novel Event-based Spatial-Temporal Transformer (short for ESTF) that conducts spatial-temporal enhanced learning and fusion for accurate action recognition. Extensive experiments on multiple benchmark datasets validated the effectiveness of our proposed framework. It sets the new SOTA performances on N-Caltech101 and ALS-DVS datasets. We hope the proposed dataset and baseline approach will boost the further development of event camera based human action recognition. In our future works, we will consider combining the color frames and event streams together for high-performance action recognition.

References

- [1] Tasweer Ahmad, Lianwen Jin, Xin Zhang, LuoJun Lin, and Guozhi Tang. Graph convolutional neural network for action recognition: A comprehensive survey. *IEEE Transactions on Artificial Intelligence*, 2021.
- [2] 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.
- [3] Stefanie Anna Baby, Bimal Vinod, Chaitanya Chinni, and Kaushik Mitra. Dynamic vision sensors for human activity recognition. In *2017 4th IAPR Asian Conference on Pattern Recognition (ACPR)*, pages 316–321. IEEE, 2017.
- [4] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, volume 2, page 4, 2021.
- [5] Yin Bi, Aaron Chadha, Alhabib Abbas, Eirina Bourtsoulatze, and Yiannis Andreopoulos. Graph-based spatio-temporal feature learning for neuromorphic vision sensing. *IEEE Transactions on Image Processing*, 29:9084–9098, 2020.
- [6] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A 240×180 130 db 3 μ s latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014.
- [7] Fabian Caba Heilbron, Victor Escorcia, Bernard Ghanem, and Juan Carlos Niebles. ActivityNet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 961–970, 2015.
- [8] Marco Cannici, Marco Ciccone, Andrea Romanoni, and Matteo Matteucci. A differentiable recurrent surface for asynchronous event-based data. In *European Conference on Computer Vision*, pages 136–152. Springer, 2020.
- [9] Enea Ceolini, Charlotte Frenkel, Sumit Bam Shrestha, Gemma Taverni, Lyes Khacef, Melika Payvand, and Elisa Donati. Hand-gesture recognition based on emg and event-based camera sensor fusion: A benchmark in neuromorphic computing. *Frontiers in Neuroscience*, 14:637, 2020.
- [10] Guang Chen, Zhongcong Xu, Zhijun Li, Huajin Tang, Sanqing Qu, Kejia Ren, and Alois Knoll. A novel illumination-robust hand gesture recognition system with event-based neuromorphic vision sensor. *IEEE Transactions on Automation Science and Engineering*, 18(2):508–520, 2021.
- [11] Huaijin Chen, WanJia Liu, Rishab Goel, Rhonald C Lua, Siddharth Mittal, Yuzhong Huang, Ashok Veeraraghavan, and Ankit B Patel. Fast retinomorphic event-driven representations for video gameplay and action recognition. *IEEE Transactions on Computational Imaging*, 6:276–290, 2019.
- [12] Junming Chen, Jingjing Meng, Xinchao Wang, and Junsong Yuan. Dynamic graph cnn for event-camera based gesture recognition. In *2020 IEEE International Symposium on Circuits and Systems (ISCAS)*, pages 1–5. IEEE, 2020.
- [13] Shoushun Chen and Menghan Guo. Live demonstration: Celex-v: a 1m pixel multi-mode event-based sensor. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pages 1682–1683. IEEE, 2019.
- [14] Xavier Clady, Jean-Mathieu Maro, Sébastien Barré, and Ryad B Benosman. A motion-based feature for event-based pattern recognition. *Frontiers in neuroscience*, 10:594, 2017.
- [15] Yongjian Deng, Hao Chen, and Youfu Li. Mvf-net: A multi-view fusion network for event-based object classification. *IEEE Transactions on Circuits and Systems for Video Technology*, 2021.
- [16] Yongjian Deng, Youfu Li, and Hao Chen. Amae: Adaptive motion-agnostic encoder for event-based object classification. *IEEE Robotics and Automation Letters*, 5(3):4596–4603, 2020.
- [17] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.

- [18] Wei Fang, Zhaofei Yu, Yanqi Chen, Tiejun Huang, Timothée Masquelier, and Yonghong Tian. Deep residual learning in spiking neural networks. *NeurIPS*, 2021.
- [19] Christoph Feichtenhofer. X3d: Expanding architectures for efficient video recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 203–213, 2020.
- [20] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019.
- [21] Guillermo Gallego, Tobi Delbruck, Garrick Michael Orchard, Chiara Bartolozzi, Brian Tabar, Andrea Censi, Stefan Leutenegger, Andrew Davison, Jorg Conradt, Kostas Daniilidis, et al. Event-based vision: A survey. *IEEE transactions on pattern analysis and machine intelligence*.
- [22] Daniel Gehrig, Antonio Loquercio, Konstantinos G Derpanis, and Davide Scaramuzza. End-to-end learning of representations for asynchronous event-based data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5633–5643, 2019.
- [23] Arun M George, Dighanchal Banerjee, Sounak Dey, Arijit Mukherjee, and P Balamurali. A reservoir-based convolutional spiking neural network for gesture recognition from dvs input. In *2020 International Joint Conference on Neural Networks (IJCNN)*, pages 1–9. IEEE, 2020.
- [24] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6047–6056, 2018.
- [25] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [26] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [27] Junho Kim, Jaehyeok Bae, Gangin Park, Dongsu Zhang, and Young Min Kim. N-imagenet: Towards robust, fine-grained object recognition with event cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2146–2156, 2021.
- [28] Orit Kliper-Gross, Tal Hassner, and Lior Wolf. The action similarity labeling challenge. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(3):615–621, 2011.
- [29] Yu Kong and Yun Fu. Human action recognition and prediction: A survey. *arXiv preprint arXiv:1806.11230*, 2018.
- [30] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *2011 International conference on computer vision*, pages 2556–2563. IEEE, 2011.
- [31] Xavier Lagorce, Garrick Orchard, Francesco Galluppi, Bertram E Shi, and Ryad B Benosman. Hots: a hierarchy of event-based time-surfaces for pattern recognition. *IEEE transactions on pattern analysis and machine intelligence*, 39(7):1346–1359, 2016.
- [32] Hongmin Li, Guoqi Li, Xiangyang Ji, and Luping Shi. Deep representation via convolutional neural network for classification of spatiotemporal event streams. *Neurocomputing*, 299:1–9, 2018.
- [33] Hongmin Li, Hanchao Liu, Xiangyang Ji, Guoqi Li, and Luping Shi. Cifar10-dvs: an event-stream dataset for object classification. *Frontiers in neuroscience*, 11:309, 2017.
- [34] Jianing Li, Xiao Wang, Lin Zhu, Jia Li, Tiejun Huang, and Yonghong Tian. Retinomorphic object detection in asynchronous visual streams. In *Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, February 22-March 1, 2022*, pages 1332–1340. AAAI Press, 2022.
- [35] Ji Lin, Chuang Gan, and Song Han. Tsm: Temporal shift module for efficient video understanding. In *Proceedings of the IEEE International Conference on Computer Vision*, 2019.
- [36] Yihan Lin, Wei Ding, Shaohua Qiang, Lei Deng, and Guoqi Li. Es-imagenet: A million event-stream classification dataset for spiking neural networks. *Frontiers in neuroscience*, page 1546, 2021.
- [37] Qianhui Liu, Gang Pan, Haibo Ruan, Dong Xing, Qi Xu, and Huajin Tang. Unsupervised aer object recognition based on multiscale spatio-temporal features and spiking neurons. *IEEE transactions on neural networks and learning systems*, 31(12):5300–5311, 2020.
- [38] Qianhui Liu, Dong Xing, Huajin Tang, De Ma, and Gang Pan. Event-based action recognition using motion information and spiking neural networks. In *IJCAI 2021, Virtual Event / Montreal, Canada, 19-27 August 2021*, pages 1743–1749. ijcai.org, 2021.
- [39] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021.
- [40] Zhaoyang Liu, Limin Wang, Wayne Wu, Chen Qian, and Tong Lu. Tam: Temporal adaptive module for video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 13708–13718, 2021.
- [41] Aref Moqadam Mehr, Saeed Reza Kheradpisheh, and Hadi Farahani. Action recognition using supervised spiking neural networks. *arXiv preprint arXiv:1911.03630*, 2019.
- [42] Shu Miao, Guang Chen, Xiangyu Ning, Yang Zi, Kejia Ren, Zhenshan Bing, and Alois Knoll. Neuromorphic vision datasets for pedestrian detection, action recognition, and fall detection. *Frontiers in neurorobotics*, 13:38, 2019.
- [43] Mathew Monfort, Alex Andonian, Bolei Zhou, Kandan Ramakrishnan, Sarah Adel Bargal, Tom Yan, Lisa Brown, Quanfu Fan, Dan Gutfreund, Carl Vondrick, et al. Moments in time dataset: one million videos for event understanding. *IEEE transactions on pattern analysis and machine intelligence*, 42(2):502–508, 2019.
- [44] Garrick Orchard, Ajinkya Jayawant, Gregory K Cohen, and Nitish Thakor. Converting static image datasets to spiking

- neuromorphic datasets using saccades. *Frontiers in neuroscience*, 9:437, 2015.
- [45] Priyadarshini Panda and Narayan Srinivasa. Learning to recognize actions from limited training examples using a recurrent spiking neural model. *Frontiers in neuroscience*, 12:126, 2018.
- [46] 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, 2019.
- [47] Christoph Posch, Daniel Matolin, and Rainer 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*, 46(1):259–275, 2010.
- [48] Ali Samadzadeh, Fatemeh Sadat Tabatabaei Far, Ali Javadi, Ahmad Nickabadi, and Morteza Haghiri Chehreghani. Convolutional spiking neural networks for spatio-temporal feature extraction. *arXiv preprint arXiv:2003.12346*, 2020.
- [49] Teresa Serrano-Gotarredona and Bernabé Linares-Barranco. Poker-dvs and mnist-dvs. their history, how they were made, and other details. *Frontiers in neuroscience*, 9:481, 2015.
- [50] Gunnar A Sigurdsson, Gül Varol, Xiaolong Wang, Ali Farhadi, Ivan Laptev, and Abhinav Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *European Conference on Computer Vision*, pages 510–526. Springer, 2016.
- [51] Amos Sironi, Manuele Brambilla, Nicolas Bourdis, Xavier Lagorce, and Ryad Benosman. Hats: Histograms of averaged time surfaces for robust event-based object classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1731–1740, 2018.
- [52] Xiaolin Song, Cuiling Lan, Wenjun Zeng, Junliang Xing, Xiaoyan Sun, and Jingyu Yang. Temporal–spatial mapping for action recognition. *IEEE Transactions on Circuits and Systems for Video Technology*, 30(3):748–759, 2019.
- [53] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- [54] Ilya Sutskever, James Martens, George Dahl, and Geoffrey Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147. PMLR, 2013.
- [55] Du Tran, Lubomir Bourdev, Rob Fergus, Lorenzo Torresani, and Manohar Paluri. Learning spatiotemporal features with 3d convolutional networks. In *Proceedings of the IEEE international conference on computer vision*, pages 4489–4497, 2015.
- [56] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018.
- [57] Qinyi Wang, Yexin Zhang, Junsong Yuan, and Yilong Lu. Space-time event clouds for gesture recognition: From rgb cameras to event cameras. In *2019 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1826–1835. IEEE, 2019.
- [58] 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. *arXiv preprint arXiv:2108.05015*, 2021.
- [59] Yanxiang Wang, Bowen Du, Yiran Shen, Kai Wu, Guangrong Zhao, Jianguo Sun, and Hongkai Wen. Ev-gait: Event-based robust gait recognition using dynamic vision sensors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6358–6367, 2019.
- [60] Yanxiang Wang, Xian Zhang, Yiran Shen, Bowen Du, Guangrong Zhao, Lizhen Cui Cui Lizhen, and Hongkai Wen. Event-stream representation for human gaits identification using deep neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [61] Zhengwei Wang, Qi She, and Aljosa Smolic. Action-net: Multipath excitation for action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13214–13223, 2021.
- [62] Xiao Wu and Junsong Yuan. Multipath event-based network for low-power human action recognition. In *2020 IEEE 6th World Forum on Internet of Things (WF-IoT)*, pages 1–5. IEEE, 2020.
- [63] Bochen Xie, Yongjian Deng, Zhanpeng Shao, Hai Liu, and Youfu Li. Vmv-gcn: Volumetric multi-view based graph cnn for event stream classification. *IEEE Robotics and Automation Letters*, 7(2):1976–1983, 2022.
- [64] 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.
- [65] Man Yao, Huanhuan Gao, Guangshe Zhao, Dingheng Wang, Yihan Lin, Zhaoxu Yang, and Guoqi Li. Temporal-wise attention spiking neural networks for event streams classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10221–10230, 2021.
- [66] Alex Zihao Zhu and Liangzhe Yuan. Ev-flownet: Self-supervised optical flow estimation for event-based cameras. In *Robotics: Science and Systems*, 2018.
- [67] Lin Zhu, Jianing Li, Xiao Wang, Tiejun Huang, and Yonghong Tian. Neuspikes-net: High speed video reconstruction via bio-inspired neuromorphic cameras. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2400–2409, 2021.
- [68] Lin Zhu, Xiao Wang, Yi Chang, Jianing Li, Tiejun Huang, and Yonghong Tian. Event-based video reconstruction via potential-assisted spiking neural network. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3594–3604, 2022.