1



A New Spiking Convolutional Recurrent Neural Network (SCRNN) with applications to Event-based Hand Gesture Recognition

Yannan Xing $^{1,\ast},$ Gaetano Di Caterina 1 and John Soraghan 1

¹Neuromorphic Sensor Signal Processing Laboratory, Centre for Signal and Image Processing(CeSIP), Department of Electronic and Electrical Engineering, University of Strathclyde, Glasgow, UK

Correspondence*: Yannan Xing yannan.xing@strath.ac.uk

2 ABSTRACT

The combination of neuromorphic visual sensors and spiking neural network offers a high effici-З ent bio-inspired solution to real-world applications. However, processing event- based sequences 4 still remain challenging because of the nature of their asynchronism and sparsity behaviour. In 5 this paper, a novel spiking convolutional recurrent neural network (SCRNN) architecture that 6 takes advantage of both convolution operation and recurrent connectivity to maintain the spatial 7 and temporal relations from event-based sequence data are presented. The use of recurrent 8 architecture enables the network to have arbitrary length of sampling window allowing the netw-9 ork to exploit temporal correlations between event collections. Rather than standard ANN to 10 SNN conversion techniques, the network utilizes supervised Spike Layer Error Reassignment 11 (SLAYER) training mechanism that allows the network to adapt to neuromorphic (event-based) 12 data directly. The network structure is validated on the DVS gesture dataset and it has achieved 13 a 10 class gesture recognition accuracy of 96.59% and 11 class gesture recognition accuracy of 14 15 90.28%.

1 INTRODUCTION

During the past couple of decades, computer vision applications have become increasingly important 16 in many industrial domains such as security systems, robotics, medical devices. Many Deep Neural 17 Network(DNN) based algorithms have outperformed human performance in different image recognition 18 tasks such as the success of Convolutional Neural Networks (CNN) (Krizhevsky et al., 2012) in the 2012 19 ILSVRC image classification challenge. However, it remains a challenge to extend the achievements 20 21 in static image recognition to dynamic scene recognition, which has strong both temporal and spatial 22 correlations. Human hand gesture recognition is one such problem that is significant for human-computer 23 interaction (Rautaray and Agrawal, 2012; Haria et al., 2017; Mitra and Acharya, 2007). The hand's 24 movement conveys certain information that can be used as a tool to communicate with computers. The hand 25 gesture recognition has been shown a significant value in applications such as virtual reality (Wickeroth 26 et al., 2009; Frati and Prattichizzo, 2011), robot control (Droeschel et al., 2011; Liu and Wang, 2018) and sign language recognition (Pigou et al., 2015; Liang and Ouhyoung, 1998; Yang et al., 2010). The 27

importance of developing intelligent models for complex Spatio-temporal processing is widely recognized 28 29 for solving dynamic scene based recognition problems. In recent years, recurrent neural network (RNN) structures such as the long-short-term-memory (LSTM) (Hochreiter and Schmidhuber, 1997) have been 30 shown to be effective for time-based sequence to sequence classification and prediction tasks. However, the 31 LSTM is still inherently inefficient for the dynamic scene recognition since it does not deal with any spatial 32 information. Research has shown the effectiveness of combining the recurrent structure and convolution 33 operation in the dynamic scene recognition such as CNN-LSTM structure (Wang et al., 2017a; Donahue 34 et al., 2017) and convLSTM structure (Shi et al., 2015; Song et al., 2018; Zhou et al., 2018). Such a 35 mechanism allows feature extraction to use both temporal and spatial information. 36

37 Concerning the data acquisition side, the traditional vision sensor is a digital camera that repeatedly 38 refreshes its entire array of pixel values at a predefined frame rate. However, using the digital camera has three drawbacks for dynamic motion recognition. First, a digital camera normally operates with a 39 40 predefined frame sampling rate (typically range 25-50 frames per second), which limits the temporal 41 resolution of activities observed. Secondly, consecutive frames and redundant pixels in each frame waste 42 significant storage resources and computation. Thirdly, the dynamic range of traditional image sensors is 43 limited by its exposure time and integration capacity. Most cameras suffer from saturating linear response 44 with dynamic range limited to 60-70dB where light from natural scenes can reach approximately 140dB of dynamic range (Posch et al., 2011a). The dynamic vision sensor (DVS) (Lichtsteiner et al., 2008a; Posch 45 46 et al., 2011b; Brandli et al., 2014) provides a solution to these problems. The DVS using address event 47 representation (AER) is an event-driven technology based on the human visual system. The benefit of the event-based sensor on dynamic scene recognition task is that it offers very high temporal resolution when 48 a large fraction of scene changes, which can only be matched by a high-speed digital camera with the 49 50 requirement of high power and significant resources.

In DVS, information is coded and transmitted as electric pulses (or spikes), which is similar to the 51 processing mechanism in biological sensory systems. The output of DVS is generated asynchronously 52 by comparing each activity of a retina pixel with a certain threshold. The emergence of dynamic vision 53 sensor (DVS)(Lichtsteiner et al., 2008b) demonstrated significant potential in applications of ultra-fast 54 power efficient computing. Compared to traditional vision sensors, DVS returns unsynchronized events 55 rather than sampled time-based frame series. For a given real-world input, DVS records only changes in 56 pixel intensity values and outputs a stream of ON/OFF discrete events regarding the changing polarity. 57 Such an event- based acquisition mechanism offers many advantages such as low power consumption, less 58 59 redundant information, low latency and high dynamic range. Despite the advantages of DVS, it is still challenging to apply the traditional computer vision algorithms to unsynchronized DVS output data. 60

The spiking neural network (SNN) provides an efficient solution to event-based data processing. As the 61 DVS mimics the biological retina, spiking neural network (SNN) mimics the human brain's functionality 62 by utilizing bio-inspired neuron and synapse models. The major difference between SNN and traditional 63 ANNs is the information carrier between their fundamental processing units. The SNN propagates only 64 individual spikes rather than floating-point numbers. Such characteristic provides an effective and low 65 power computing strategy for event-driven inputs. Previous work has demonstrated application examples 66 of combining SNN and event-based visual sensor such as extracting car trajectories on a freeway[10], 67 recognition of human postures (Pérez-Carrasco et al., 2010; Jiang et al., 2019), object tracking(Hinz 68 et al., 2017) and human gesture recognition (Amir et al., 2017). However, to our knowledge to date, 69 the convolutional recurrent network structure that is particularly designed for gesture recognition has 70 not been widely investigated in the SNN domain. Wang et al. (Wang et al., 2019b) presented a spiking 71

recurrent neural network that used for action recognition, but the term "spiking" in their work does not 72 73 represent the event-based processing but a spiking signal that was used to help a traditional RNN correct its contaminated memory. Vyacheslav et al. (Demin and Nekhaev, 2018) proposed a bio-inspired learning rule 74 75 FEELING with an attempt on the recurrent structure, which is applied to the handwritten digit recognition. 76 The FEELING algorithm was further implemented by (Nekhaev and Demin, 2020) with an convolutional recurrent structure that is proved to be more energy efficient on hand digit recognition. However, this work 77 78 had not considered the research line that the combination of convolutional and recurrent structure is more 79 significant in dynamic scene based recognition(i.e., hand gesture recognition). Besides, this work ignored 80 the adaptability of SNN with neuromorphic hardware and sensors.

In this paper, we present a novel spiking neural network structure that can adapt to neuromorphic vision 81 data-based recognition problem especially for those data that contains strong spatiotemporal correlations 82 83 such as human hand gesture recognition. The convolutional operation and recurrent neural network connections are combined in an SNN that uses a supervised learning based spiking convolutional recurrent 84 neural network (SCRNN). By adjusting the integration period of the input data sequence and convolution 85 kernel, SCRNN can achieve arbitrary Spatio-temporal resolution related to the recognition demand. 86 Besides, The Spike Layer Error Reassignment (SLAYER) training algorithm (Shrestha and Orchard, 2018) 87 is successfully deployed to the SCRNN for the purpose of generalization and training stability. It utilizes 88 89 both temporal error and axonal delay credit assignment to minimize the computational complexity. The use of SLAYER effectively prevents the common gradient vanishing and explosion problem associated 90 91 with recurrent neural networks. Since the recurrent propagation between the SCRNN cells relies on the information fusion from inputs of current timestamps and output from previous timestamps. Particularly for 92 SCRNN, a spiking feature map integration method is developed in the SCRNN cell to maintain information 93 continuity in the temporal domain. Furthermore, The SCRNN is validated by a series of experiments on 94 the DVS gesture dataset (Amir et al., 2017) to prove its robustness for the motion-based neuromorphic 95 action recognition problem. 96

97 The remainder of this paper is organized as follows. Section 2 introduces the related work in the spiking 98 recurrent neural network and SLAYER training algorithm. In Section 3, detailed descriptions are provided 99 in terms of individual SCRNN cell and overall SCRNN topology. The experiment results on the DVS 100 gesture dataset is presented and discussed in Section 4. The experiment result is analyzed and compared 101 with previous work. Finally, the conclusions are provided in Section V.

2 PRELIMINARIES

This section gives an explanation of the background of SNN, the SLAYER training algorithms (Shresthaand Orchard, 2018) as well as relevant previous works on convolutional recurrent neural networks.

104 2.1 Spiking Neural Network

In recent years, deep learning technologies have rapidly revolutionized the field of machine learning. Traditional deep neural networks are trained using supervised learning algorithms, which are usually based on gradient descent backpropagation. A neural network comprises several fundamental computing units (neurons) containing weighted and biased continuous activation function. The typical example of these activation functions are sigmoid, hyperbolic tangent and ReLU (Nair and Hinton, 2010). With the feed-forward and recurrent structure, this computation strategy allows them to be able to approximate any analog function universally (Vreeken, 2002).

Although DNNs were initially brain-inspired, their structure, neural information processing and learning 112 method are still fundamentally different from the brain. One of the most distinctive difference is the means 113 in which information is carried between neurons. That is one of the main reasons for the increased interest in 114 spiking neural networks (SNNs). SNN raises the level of biological realism of ANNs by utilizing individual 115 spikes as information carriers. This allows the network computation and communication to incorporate 116 spatial-temporal information. The spikes used in SNN, however, are sparse in time with uniform amplitude, 117 but rich in their information content when they occur in time. The information in SNNs is presented by 118 spike timing e.g. latency, frequency or the population of the neuron that are emitted spikes (Gerstner et al., 119 120 2014).

121 The SNN is an ideal universal spike generation model that mimics the actual biophysical mechanisms 122 describes by Hodgkin and Huxley (Hodgkin and Huxley, 1990a). The spikes are only identified at the time instant when they arrive at the post-synaptic neuron. Non-linear differential equations are commonly used 123 in SNN neuron modeling to generated the membrane potential through the time (Abbott, 1999; Hodgkin 124 and Huxley, 1990b; Teka et al., 2014; Gerstner, 2009). Figure 1 illustrates the basic operating mechanism 125 of a spiking neuron. This illustrates a single spiking neuron that receives incoming spike trains from s_1 , 126 s_2 and s_3 and generates an output spike as shown in Figure 1(a). The incoming spikes to a neuron are 127 integrated and transferred to the membrane potential dynamics u(t) as is shown in Figure 1(b). Whenever 128 the membrane potential reaches a certain threshold value ϑ , the spiking neuron will emit a spike and reset 129 the membrane potential to its resting value u_{rest} . After a spike activity, the neuron enters the refractory 130 131 period and cannot fire any further spikes until its membrane potential resets to its resting value.





A typical spiking neuron model can contain additional parameters that approximate the membrane potential dynamics in the neural cortex. Commonly used spiking neuron model in SNNs include: Integrate and fire neurons(IF) (Feng, 2001; Feng and Brown, 2000), Leaky integrated and fire neurons(LIF) (Liu and Wang, 2001), Hodgkin-Huxley model (Bower et al., 1995) and Spike Response Model(SRM) (Gerstner, 2008) etc.

Recent research has successfully demonstrated examples of SNN based applications including object 137 recognition (Kheradpisheh et al., 2018; Diehl and Cook, 2015), speech processing (Loiselle et al., 2006; 138 Tavanaei and Maida, 2017; Wysoski et al., 2010), pattern recognition (Kasabov et al., 2013; Mohemmed 139 et al., 2012; Han and Taha, 2010; Dhoble et al., 2012). Furthermore, many developed neuromorphic 140 computing platforms have demonstrated tremendous potential in real-world power limited applications. 141 The IBM TrueNorth systems consist of 5.4 billion transistors with only 70mW power density consumption, 142 which accounts for only 1/10000 of traditional computing units (Akopyan et al., 2015). The SpiNNaker 143 144 platform (Furber et al., 2014, 2013) developed by Researchers at Manchester provides ASIC solutions to hardware implementations of SNNs. It utilized multiple ARM cores and FPGAs to configure the hardware 145 and PyNN (Davison et al., 2009) software API to enable the scalability of the platform. The Loihi NM 146 chip (Davies et al., 2018) is a digital NM computing platform that was recently announced by Intel. One of 147 the most attractive features of Loihi is the potential of online-learning. Loihi has a special programmable 148 microcode engine for SNN training on the fly. The emergence of these hardware technologies demonstrates 149 strong suitability of applying power efficient neuromorphic computing into real-world mobile units. 150

151 2.2 Spike Layer Error Reassignment in Time(SLAYER)

Currently, the training procedure of most ANNs relies on the combination of continuously differentiable activation function and gradient descent convergence algorithm. Spiking Neural Networks are similar to traditional neural networks in topology but differ in the way of information carrier and the choice of neuron models. The non-differentiable nature of biological-plausible spiking neurons is the main challenge of the development of SNN training algorithms. Spike Layer Error Reassignment in Time (SLAYER) alternatively approximates the derivative of the spike function based on the neuron state changes and assigns the error to previous layers. A description of SLAYER training algorithm is provided in the next subsection.

The neuron model used for the SLAYER is the Spike Response Model (SRM). The membrane potential generation process of a SRM neuron is achieved by convolving a spike response kernel $\sigma(t)$ with the incoming spike train $s_i(t)$ to this neuron to form a spike response signal as $a(t) = (\sigma(t) * s_i(t))$. Here the index *i* represents the i_{th} input channel. The spike response signal is further weighted by the synaptic weight *w*. Similarly, the refractory response signal can be obtained via convolving a refractory kernel $\nu(t)$ with the neuron output spike train $s_o(t)$ as $r(t) = (\nu(t) * s_o(t))$. The overall neuron membrane potential u(t) can be obtained by summing all the spike response signal and refractory response signal as:

$$u(t) = \sum w_i(\sigma(t) * s_i(t)) + (\nu(t) * s_o(t))$$

= $\mathbf{W}^\top \mathbf{a}(t) + r(t)$ (1)

166 The generated membrane potential u(t) is then compared with a predefined threshold ϑ and output spike 167 when $u(t) > \vartheta$ like is shown in Figure 1. In a multilayer feedforward spiking neural network architecture, 168 instead of directly managing the non-differentiable spike neuron equations, SLAYER approximates the 169 derivative of the spike function as a probability density function (PDF) of spike state changes. Further 170 details of the model and its use in training the SNN can be found in (Shrestha and Orchard, 2018). With

171 a good estimation PDF as the derivative term of spike change state, the SLAYER can easily derive the

172 gradient of weights and delays in each layer from a feedforward SNN. This allows the network to adapt

173 developed gradient descent method for optimization purpose such as ADAM (Kingma and Ba, 2015),

174 RmsProp (Hinton et al., 2012).

175 2.3 Convolutional Recurrent Neural Network

The convolutional recurrent neural network (CRNN) structure has been well studied in the second generation of ANNs. The convolution operation in the ANNs usually acts as a spatial visual feature extractor that assumes features are in different levels of hierarchy. The recurrent structure introduces memory to the network and an ability to deal with sequential data dependently.

A significant design of the CRNN structure is the ConvLSTM structure (Shi et al., 2015) that was initially designed for forecasting precipitation. By replacing the general gate activation by the convolutional operation, the network is able to exploit an extracted 3D tensor as the cell state. The ConvLSTM was also evaluated on the moving MNIST (Srivastava et al., 2015) dataset and was shown to successfully separate the overlapping digits and predicted the overall motion with a high level of accuracy.

Another CRNN structure CNN-LSTM concatenates a CNN and an LSTM to formulate a collaborative network. The LSTM in the structure is placed behind a pretrained CNN that directly takes the output feature vector from the CNN as the input sequence. The implementation of this structure however is highly dependent on a well pre-trained CNN that was designed for the interest as the feature extractor. The CNN-LSTM is proved powerful in many application domains such as acoustic scene classification (Bae et al., 2016), emotion recognition (Fan et al., 2016), action recognition Wang et al. (2017b) etc.

Over the past few years, researchers have successfully applied CRNN in medical applications (Wang, 191 Lebo and Li, Kaiming and Chen, Xu and Hu, 2019), speech processing (Tan and Wang, 2018; Cakir et al., 192 2017), music classification (Choi et al., 2017). Adopting a recurrent structure enables the neural network to 193 194 encapsulate the global information while local features are extracted by the convolution layers. Yang et al. (Haodong Yang, Jun Zhang, Shuohao Li and Chen, 2018) demonstrated a Convolutional LSTM network 195 that was successfully evaluated on various action recognition datasets. The importance of using CRNN 196 structure in the application of human action recognition is that unlike action recognition in images, the 197 same task in videos relies on motion dynamics in addition to visual appearance. Although CNNs and its 198 variants like 3D convolution (Ji et al., 2013; Karpathy et al., 2014) achieves good performance, they still do 199 not make sufficient use of temporal relations between frames. More recently, Maj et al. (Maid, Mahshid 200 and Safabakhsh, 2019) designed a motion-ware ConvLSTM for the action recognition task which is an 201 LSTM unit that considers the correlation of consecutive video frames in addition to the Spatio-temporal 202 information. 203

However, in the SNN domain, the CRNN structure has not been widely investigated especially for the action recognition problem. One of the main challenges in developing a spiking CRNN is how to manage the training process of spiking neurons. Besides, the consecutive information recurrency is difficult to achieve in the SNN since the traditional probabilistic based functions do not comply with spikes. In this paper, the SLAYER algorithm is used as an efficient, general supervised training mechanism for SNNs. Based on the spiking model of SLAYER, we design a network structure that can achieve both forward and recurrent information propagation.



Figure 2. The 3D spiking convolution operation, Red: represents the spiking convolution through a defined 3D volume

3 SPIKING CONVOLUTIONAL RECURRENT NEURAL NETWORK(SCRNN)

In this section, the novel system using SCRNN for action recognition is described. The fundamentals of3D spiking convolution and the related SCRNN model are described in the following subsections.

213 3.1 Spiking Convolution Operation

Consider an input sequence S(n), n = 0, 1, 2, ...N as is illustrated in Figure 2. At each time step, S(n) is a 3D tensor with shape $\{u, v, t\}$ where u and v denote the width and height of each frame and t correspond to the pre-defined time resolution. For a given event-based video stream, it can be arbitrarily segmented into several tensors according to the desired temporal frequency. For example, for a 1.5sec 128x128 resolution events data stream with 30ms temporal resolution and 1ms sampling time can form a input sequence S(n), n = 0, 1, 2, ...50. For each segments, the tensor shape is $\{128, 128, 30\}$.

220 The sampled input tensor S(n) with a shape of $\{u, v, t\}$ is convolved with a 3D convolutional kernel to generate a spiking neuronal feature map. The spikes within an arbitrary kernel can be regarded as a bunch 221 of spike trains $s_{u,v}(t)$ where each spike train corresponds to the spikes at a specific coordinate (u, v) within 222 the temporal resolution window t. Each neuron in the feature map receives the spikes from the neurons in 223 224 the 3D convolutional kernel. The spikes in the region of the kernel are integrated to generate membrane 225 potential for a single neuron in the feature map. The neurons in a map detect the Spatio-temporal dynamic patterns in different 3D volumes. Unlike the standard feature map generated by CNN, the information at 226 each coordinate in a spiking feature map is expressed by spike trains which can be considered as a spiking 227 representation of detected patterns. 228

The convolutional kernel is highly overlapped to make sure the proper detection of features. The SRM neuron model is used to describe the 3D spiking convolution operation, which gathers all the input spikes from pre-synaptic neurons and outputs spike when the membrane potential reaches the pre-defined threshold. In the SLAYER, this is done by convolving the spike trains in the kernel with a spike response kernel and followed by the threshold function. Each spike train will be transferred to the spike response signal then further to the membrane potential of the postsynaptic neuron. The process can be expressed as:

$$a_{u,v}(t) = s_{u,v}(t) * \sigma(t) \tag{2}$$

$$u_{j,k}(t) = \sum_{m=1}^{K} \sum_{n=1}^{K} \mathbf{W}_{m,n} a_{j+m-1,k+n-1}(t) + (s_{j,k}(t) * \nu(t))$$
(3)

$$s_{j,k}(t) = 1 \& u_{j,k}(t) = 0 \text{ when } u_{j,k}(t) \ge V_{thr}$$
(4)

where **W** denotes to the synaptic weights. u and v are the vertical and horizontal coordinate index of the input tensor. j and k represents the vertical and horizontal coordinate in the feature map. K represents the convolution kernel width and height.

The 3D spiking convolution can decompose the input event based data into several spatio-temporal pattern feature maps, where each spike in the map corresponds to a specific pattern. When multiple spiking convolution layers are used, the feature in a layer is a combination of several low level features extracted from the previous layer.



Figure 3. The proposed single SCRNN cell. The state spiking feature map and input feature map are combined in the cell with an output feature map recurrently connected to the cell

242 3.1.1 SCRNN Cell

243 The SCRNN cell is designed as the fundamental unit of the SCRNN system. The idea was inspired by the structure of the ConvLSTM cell (Shi et al., 2015). A graphical illustration of a single SCRNN cell is 244 shown in Figure 3. The inputs to the cell comprise two parts: First is the spiking feature map generated by 245 the outside events(e.g., a fragment from an event-based action data). The second part is the hidden spiking 246 states which represent the fused feature map of previous states and the feature map generated by the current 247 input. To ensure the state feature map has the same shape as the input, a padding technique is needed 248 before the actual convolution operation, which means padding empty events(zeros) on the boundary of 249 state maps. This can be viewed as the current state having no prior knowledge in terms of the region outside 250 the current receptive field. At zero time index, the internal state needs to be initialized randomly or set 251

empty which represents no prior knowledge at the beginning from the temporal perspective. Consequently, the 3D spiking convolution operation is applied to both input-to-internal state transitions and state-to-state transitions in an SCRNN cell. The future state to state transition is achieved by utilizing another 3D convolution layer that contains a pre-defined number of hidden neurons. Two feature maps are concatenated to form a single map. Then the spikes in the same kernel of the fusion map are accumulated and activated to generate the membrane potential signal for future states. Consider an input segment X_i . The entire computation process within an SCRNN cell can be written as:

$$s_i(t) = \theta\{\sum W_{ih}(X_i * \sigma)\}\tag{5}$$

$$s_h(t) = \theta\{\sum W_{hi}(s_h(t-1)*\sigma)\}\tag{6}$$

$$s_h(t+1) = \theta\{\sum W_{hh}(s_i(t) * \sigma + s_h(t) * \sigma)\}$$
(7)

$$s_o(t) = \theta\{\sum W_{ho}(s_i(t) * \sigma + s_h(t) * \sigma)\}$$
(8)

where θ represents the thresholding operation. W_{ih} , W_{hi} , W_{hh} and W_{ho} denotes the weight input to state, state to input, state to state and state to output respectively. It can be seen from equation (7) and (8) that the output of an SCRNN cell comprises two terms: $s_h(t+1)$ is the spiking states that can be used for future cells and the $s_o(t)$ represents the output spike train. The output from the cell represents the 3D feature map extracted from the current cell that allows the network to go deeper by using the $s_o(t)$ as the input of the next layer.



Temporal domain

Figure 4. The proposed SCRNN structure which is comprised by prior defined individual SCRNN cells. The information going through the vertical direction in the Figure 4 is the spiking convolutional operation in the spatial domain. The information processing along with the horizontal direction in the Figure 4 is the recurrent process between the SCRNN cells which is in the temporal domain. h_1 , h_2 and h_3 is the initial feature map assumption prior to the zero index. X_n and Y_n represents the n_{th} input or the output sequences

265 3.2 Spiking Convolutional Recurrent Neural Network

The overall SCRNN architecture shown in Figure 4 comprises a combination of single cells that are 266 stacked in both temporal and spatial processing domain. From a temporal point of view, the cells can 267 process the input sequence separately using the internal state correlations. Furthermore, the input can be 268 further decomposed by adding additional cells at each time step, thus allowing the network to form greater 269 computational complexity and processing higher level spatial features. In other words, at a specific time 270 step, the concatenated SCRNN cells (layers) can be treated as a standard spiking convolutional neural 271 network wherein each input of an SCRNN cell is the output signal of the previous cell. It should be noted 272 that additional initial states are needed for every added layer. 273

274 Similarly to the conventional recurrent neural network, the SCRNN can also be unrolled to form a short-term feed-forward structure that increases the network parameter capacity. Unrolling a recurrent 275 structure represents a trade-off between the network performance and the computational cost. Although 276 theoretically the cells can be unrolled up to the length of the input sequence, the computation cost in the 277 training process increases dramatically along with the number of cells. Moreover, to guarantee the network 278 performance in terms of temporal information, the backpropagation through time (BPTT) (Werbos, 1990) 279 is used which is another factor that affects the training speed. BPTT calculates and accumulates errors 280 across each time step, which can be computationally expensive as the number of time step increases. 281



Figure 5. The demonstration of DVS gesture dataset with integral time of 0.5s. The gesture showing in the example is hand waving. The green and red edges in each Figure 5 represents the ON/OFF polarities of spikes

4 EXPERIMENT RESULTS

In this section, the experimental result of action recognition using SCRNN will be presented. To validate the robustness of the SCRNN, we evaluated the network structure by performing the recognition task on the IBM DVS gesture dataset (Amir et al., 2017). The DVS gesture dataset comprises recordings of 29 different actors carrying out 10 different hand gesture actions. All recordings are captured by an Inilabs 128 x 128 dynamic vision sensor under three different lighting conditions. Each gesture sample has a duration of approximately 6 second. Figure 5 shows an example of hand waving gesture with 0.5s integral time interval
in nature light condition. The goal is to classify the gesture event video data into a corresponding label.
The DVS gesture dataset is split as 1176 samples for training and 288 samples for testing as annotated. We
construct a three layer SCRNN to solve this problem as is shown in Figure 4. The SRM response neuron
parameters are shown in Table 1.

Table 1. The neuron parameter setting for the SCRNN simulation.

ϑ_{neuron}	τ_{neuron}	$ au_{ref}$	C_{ref}	tau_f	C_f
5	10	1	2	1	1

The parameters define the standard neuron dynamics behavior which is used in all SCRNN networks. Where ϑ_{neuron} is the neuron firing threshold. τ_{neuron} is the neuron time constant, τ_{ref} is the neuron refractory time constant, C_{ref} is the refractory response scaling coefficient, tau_f is the neuron spike function derivative time constant, and the C_f is the neuron spike function derivative scaling coefficient.

As the gesture recognition is a many-to-one problem, only the output from the last layer and last time step SCRNN cell are taken into account for the loss calculation. The loss function used in this method is defined as the square error based on the number of spikes between the target and actual output in a time window according to Shrestha and Orchard (2018). With the S_o denotes to the output spike train of the last layer of SCRNN and \hat{S} indicates to the target spike train, the loss function L can be expressed as follows.

$$L = \frac{1}{2} \sum_{1}^{N} \left(\int S_o(\tau) d\tau - \int \hat{S}(\tau) d\tau \right)^2 \tag{9}$$

where N is the number of output neurons of the last layer. At each time step, the error signal is 301 calculated according to the current output spike count and target spike count. It should be noted that 302 the backpropagation pipeline covers both spatial and temporal propagating routes through the recurrent 303 connection. To save on computation resources, only 1.5s out of 6s of each gesture samples were used for the 304 experiment. The input event sequence is integrated into several frames based on pre-defined segmentation 305 length l_s . The segmentation length significantly affects the sparsity and the number of integrated frames. A 306 small l_s will results in a large number of sparse frames, on the contrary a chosen of large l_s will reduce the 307 number of frames but increase the number of events in each frame. 308

309 To evaluate the performance of SCRNN, we carried out different combinations of network parameters to perform the action recognition task. The following hyper-parameters were used in the experiments: 310 Number of filters in the convolutional layer, the segmentation length(time resolution) l_s , the target true 311 spike count Tg_{True} and target false spike count Tg_{False} . Figure 6 illustrates the output spike activities 312 before and after the training of the last layer of the SCRNN. The vertical dash line in the figures simulates 313 the time window that spikes will be counted for an input sample. In other words, the spikes between two 314 dash lines are the output from a single input instance. The output neuron index from 1 to 10 represents 10 315 different gesture classes. The red bars are target spike(labels) and the black bars are actual network output 316 317 spikes. It should be noted that the loss for the SLAYER training algorithms is calculated from the error



Figure 6. The last layer SCRNN output: (a)Before Training (b)After Training

signal that was generated according to the difference between the number of actual output spikes from the 318 network and the target spikes (Tg_{True} and Tg_{False}). If the actual spikes count of output neuron match that 319 from the target spike count then a correct prediction is implied. As shown in Figure 6(a), the SCRNN has 320 zero output before training and gradually learns to generate spikes that match the target spike in terms of 321 the target spike quantity. Figure 6(b) demonstrates the output spike monitoring after-training the SCRNN. 322 It can be clearly seen from Figure 6(b) that the actual spikes(shown in black) now have similar spike counts 323 as target spikes(shown in red) for the input samples. It should be noted that, the target spikes and actual 324 spikes have different spike timings but similar spike counts in each window. 325

Conv1	Conv2	Conv3	FC1	FC2	Tg_{Ture}	Tg_{False}	$l_s(ms)$	Trainacc	Testacc
5x5x16	3x3x32	3x3x64	1024	512	30	5	25	90.73%	85.23%
3x3x16	3x3x32	3x3x64	512	128	30	5	25	87.92%	84.64%
5x5x32	3x3x64	3x3x128	1024	512	30	5	25	93.54%	89.15%
5x5x16	3x3x32	3x3x64	1024	512	60	10	50	95.45%	91.67%
3x3x16	3x3x32	3x3x64	512	128	60	10	50	95.08%	89.39%
5x5x32	3x3x64	3x3x128	1024	512	60	10	50	98.48%	96.59%
5x5x16	3x3x32	3x3x64	1024	512	80	15	75	95.45%	88.64%
3x3x16	3x3x32	3x3x64	512	128	80	15	75	93.18%	93.56%
5x5x32	3x3x64	3x3x128	1024	512	80	15	75	96.59%	90.90%

Table 2. Comparisons of SCRNNs performance on DVS gesture dataset with different hyper-parameters. Tg_{Ture} : The preliminary setting of target True spikes count; Tg_{Ture} : The preliminary setting of target False spike count; $l_s(ms)$: The segmentation length(time resolution)

The experiment results are shown in Table 2, where each listed architecture is simulated for 100 epoch 326 over the full dataset. For each structure listed in the table, the accuracy is obtained by averaging the best 327 testing accuracy among 5 repeated experiments with different random initialized weights. Among these 328 experiments, the best testing accuracy of 10 class gesture is 96.59% with the 3 layer SCRNN structure 329 with the first convolutional layer consisted of 32 5x5 convolutional filters, second and third convolution 330 layer has 64 and 128 3x3 convolutional kernels respectively. The l_s is 50ms which represents there are total 331 1000/50 = 20 time steps. The loss and training curve for the best network structure is shown in Figure 7(a) 332 and Figure 7(b). This structure also was used to train the 11 class gesture (plus a random other gesture 333 action) and obtained a testing accuracy of 90.28%. 334



Figure 7. (a):The training and testing loss changes for 3 layer SCRNN with conv1: 5x5x32; conv2:3x3x64; conv3:3x3x128; $l_s=50ms(b)$:The training and testing accuracy changes for 3 layer SCRNN with conv1: 5x5x32; conv2:3x3x64; conv3:3x3x128; $l_s=50ms$. c: The confusion matrix for 3 layer SCRNN with conv1: 5x5x32; conv2:3x3x64; conv3:3x3x128; $l_s=50ms$; The 0-9 represents the 10 categories of gestures. 0: hand clapping; 1:right hand wave; 2: left hand wave; 3: right arm clockwise; 4: right arm counter clockwise; 5: left arm clockwise; 6: left arm counter clockwise; 7: arm roll; 8: air drums; 9: air guitar



Figure 8. The example of 3 layer SCRNN misclassification case. The 4 figures demonstrate a similarity of event dynamics between the hand clapping gesture and air drum gesture. Top left: the 3D view of a hand clapping sample with duration of approximately 1s. Top right: the 2D view of a hand clapping gesture that integrated all spikes within 1s. Bottom left: the 3D view of a air drum gesture sample with duration of approximately 1s. Top right: the integrated all spikes within 1s. Bottom left: the 2D view of a air drum gesture sample with duration of approximately 1s. Bottom right: the 2D view of a air drum gesture that integrated all spikes within 1s.

Thus, the loss can be very large at the start compared with normal loss value since the network can have 335 an empty output with untrained weights and delays. It was found that setting the $l_s = 50ms$ produces the 336 best result for SCRNN structure which can be explained as follows. First, the time resolution is matched 337 with the frame continuity for this dataset, which means the individual segmented frame can either contain 338 limited or redundant information with $l_s = 25ms$ or $l_s = 75ms$. This can possibly weaken the connection 339 between the frames from the perspective of recurrent convolutional operation. Secondly, the spike emitting 340 of neurons in each layer is important to the training process. A proper selection of l_s can make sure the 341 sparsity of frames which guaranteed the stability of the training process. 342

The confusion matrix in Figure 7(c) shows a detailed performance of the SCRNN for the 10 gesture 343 recognition tasks. Note that the amount of samples of arm roll is twice than other gestures in the original 344 dataset. It can be seen that the SCRNN achieved an overall good performance except that the confusion 345 between the hand-clapping and air drums gesture where there are totally 3 + 4 = 7 instances that SCRNN 346 misclassified the hand clapping or air drum as each other. This is due to the dynamic similarity of these 347 two gestures for some instances. Figure 8 demonstrates an example of misclassification which shows both 348 3D and 2D view of dynamics of these two gesture. From our observations, some of hand-clapping and 349 air drum gestures exhibit a strong similar spike change pattern which is a potential reason that leads to 350

Method	Type of processing	10 class	11 class
IBM TrueNorth Eedn (Amir et al., 2017)	spiking	96.49%	$\begin{array}{c} 94.59\%\\ 93.64\%\pm0.49\%\\ 95.32\%\\ 92.01\%\end{array}$
SLAYER CNN (Shrestha and Orchard, 2018)	spiking	unknown	
PointNet++ (Wang et al., 2019a)	Non-spiking	97.08%	
SCRNN	spiking	96.59%	

Table 3. Comparison of SCRNN gesture recognition results with previous work

misclassification. This further matches our initial design purpose of SCRNN, which is an action dynamics
 sensitive, event stream pattern based recognition network.

353 For comparison purpose, results from previously published work (Amir et al., 2017; Shrestha and Orchard, 2018; Wang et al., 2019a) on the IBM DVS gesture dataset is carried out which is shown in Table 3. It can 354 be seen that the SCRNN approaches the state of the art recognition accuracy and surpassing the benchmark 355 accuracy of IBM's work in 10 categories gesture classification tasks. The original work from IBM that 356 357 running on TrueNorth was trained with Eedn (Amir et al., 2017) and required extra filters and preprocessing 358 before the CNN. On the other hand, the SCRNN takes the neuromorphic data directly from the sensor 359 and the training process does not require any additional processing to the data. The SLAYER algorithms 360 (Shrestha and Orchard, 2018) using CNN with a feedforward structure achieved an accuracy of 93.64%on average for the 11 class recognition . Although the SCRNN does not outperform the SLAYER based 361 CNN network in 11 class classification, the SCRNN is still competitive at 90.28%. We conclude this 362 363 accuracy drop for the 11 class recognition task is due to the introduction of the additional class of random 364 gesture. The "other" class in the DVS gesture dataset consists of random samples and each of those is neither same as other samples nor falls into the first ten categories. The SCRNN with designed recurrent 365 366 convolution operation is found to be less effective to such type of training data. Although the SCRNN 367 although does not outperform the SLAYER based CNN network in 11 class classification, the SCRNN is still competitive at 92.01%. The pointnet++ (Wang et al., 2019a) processed individual event data by 368 369 an MLP based feedforward neural network which achieved the best accuracy in both 10 and 11 category 370 gesture recognition tasks. However, the pointnet++ is not a spiking based training algorithm that has less potential to be applied to neuromorphic hardware and the DVS data in their method needs to be modeled as 371 multiple points cloud with each spike $\{x,y,z\}$ is fed into an MLP. 372

5 EFFECT OF RECURRENT CONNECTION

To further demonstrate the effectiveness of SCRNN for the category-limited dynamic scene recognition. 373 A mini-experiment is designed to directly compare the effect of the recurrence for the 10 class gesture 374 375 recognition. A feedforward spiking convolutional neural network and an SCRNN is designed following a "same learning capacity rule" as is shown in Figure 9. The spike pooling operation was applied to reduce 376 the computational cost. The pooling was done by reducing all the spikes in a pooling kernel into one over 377 the spike presentation time. The two structures are exactly the same in neuron parameters, the number of 378 379 neurons and number of layers except the SCRNN has a recurrent connection in each convolution layer. 380 For both structure, with the segmentation length of l_s , the first layer is a pooling layer with a kernel size of $4x4xl_s$, which reduced the dimension of data from $128x128xl_s$ to $32x32xl_s$. The second layer is a 381 382 convolutional layer that has a kernel size of $3x3xl_s$ with 16 hidden neurons. The third layer is a pooling 383 layer using 2x2 kernels to further reduce the dimension of each feature map to $16x16xl_s$. The fourth layer is a convolutional layer with 32 hidden neurons with the kernel size of $3x3xl_s$, which the output is 384

flattened and fed into a fully connected layer with 5256 neurons followed by the output layer to performthe classification.

The feedforward CNN is different from the SCRNN in the training phase. For CNN, the first 1s event data of each sample with a temporal resolution of $1\text{ms}(l_s = 1000)$ is used as the input data which only needs to be fed to the network once per sample. The SCRNN takes the same length of input data in total for each sample but a segmentation length of $l_s = 50$ is selected to partition the input into 20 subsets. This represents that the SCRNN need to iteratively take the data to perform the recurrent processing.



Figure 9. The network structure for the experiments of comparison between the feedforward Spiking Convolutional Neural Network and SCRNN

Both of the designed structures are trained 100 epochs for 5 trials with different weight initializations, 392 the averaged testing accuracy dynamics of these two experiments are plotted in Figure 10. The SCRNN 393 compared to standard feedforward spiking CNN with a similar learning condition can provide a faster 394 convergence speed. As is shown in Figure 10, the averaged testing accuracy of SCRNN is stabilized after 395 approximately 40 epochs while the CNN requires about additional 25 epochs to fully converge with the 396 data. Besides, the SCRNN without the inference of the unknown class can provide a recognition accuracy 397 of 88.64% on the 10 class gesture recognition in this particular structure, while the feedforward CNN only 398 achieves 84.09%. 399

6 CONCLUSION

In this paper we presented a novel spiking convolutional recurrent neural network that was designed for
efficient human hand gesture recognition. The individual cell is able to extract the spatial features by 3D
spiking convolution operation and transferring the information recurrently.



Figure 10. The testing accuracy curve for the designed experiments

The SCRNN is successfully deployed to the DVS 128 gesture dataset. The SCRNN tested on the IBM DVS gesture dataset achieving an averaged recognition accuracy of 96.59% for 10 category classification and 90.28% for 11 category classification. We have shown that the designed SCRNN compared to standard feedforward CNN structure performs less competitive for the 'unknown' class but has the advantages in terms of convergence speech and accuracy for the fixed amount of categories.

However, we believe that the usage of SCRNN is not only limited to action recognition but can be extended to various dynamic scene recognition and prediction tasks. A further extension of this work could be a spiking-flownet-like network that used for optical flow estimation (Dosovitskiy et al., 2015). Additionally, using new neuromorphic hardware with low SWaP(Size, Weight and Power) profile, the SCRNN has the potential to be implemented as an efficient training algorithm for neuromorphic action recognition based applications. The SCRNN also has a strong potential to be implemented on Loihi chip due to the use of SLAYER algorithm.

REFERENCES

- 415 [Dataset] Abbott, L. F. (1999). Lapicque's introduction of the integrate-and-fire model neuron (1907).
 416 doi:10.1016/S0361-9230(99)00161-6
- Akopyan, F., Sawada, J., Cassidy, A., Alvarez-Icaza, R., Arthur, J., Merolla, P., et al. (2015). TrueNorth:
 Design and Tool Flow of a 65 mW 1 Million Neuron Programmable Neurosynaptic Chip. *IEEE*
- Transactions on Computer-Aided Design of Integrated Circuits and Systems doi:10.1109/TCAD.2015.
 2474396
- 421 Amir, A., Taba, B., Berg, D., Melano, T., Mckinstry, J., Di Nolfo, C., et al. (2017). A low power, fully
- event-based gesture recognition system. In *Proceedings 30th IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017.* doi:10.1109/CVPR.2017.781
- 424 Bae, S. H., Choi, I., and Kim, N. S. (2016). Acoustic Scene Classification Using Parallel Combination of
- 425 LSTM and CNN. Proceedings of the Detection and Classification of Acoustic Scenes and Events 2016

- 426 *Workshop (DCASE2016)*
- Bower, J. M., Beeman, D., Nelson, M., and Rinzel, J. (1995). The Hodgkin-Huxley Model. In *The Book of GENESIS*. doi:10.1007/978-1-4684-0189-9_4
- Brandli, C., Berner, R., Yang, M., Liu, S. C., and Delbruck, T. (2014). A 240 × 180 130 dB 3 μs latency
 global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits* doi:10.1109/JSSC.
 2014.2342715
- Cakir, E., Parascandolo, G., Heittola, T., Huttunen, H., and Virtanen, T. (2017). Convolutional Recurrent
 Neural Networks for Polyphonic Sound Event Detection. *IEEE/ACM Transactions on Audio Speech and Language Processing* doi:10.1109/TASLP.2017.2690575
- Choi, K., Fazekas, G., Sandler, M., and Cho, K. (2017). Convolutional recurrent neural networks for music
 classification. In *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings*. doi:10.1109/ICASSP.2017.7952585
- 438 Davies, M., Srinivasa, N., Lin, T. H., Chinya, G., Cao, Y., Choday, S. H., et al. (2018). Loihi: A
 439 Neuromorphic Manycore Processor with On-Chip Learning. *IEEE Micro* doi:10.1109/MM.2018.
- 440 112130359
- 441 Davison, A. P., Brüderle, D., Eppler, J., Kremkow, J., Muller, E., Pecevski, D., et al. (2009). PyNN: A
 442 common interface for neuronal network simulators. *Frontiers in Neuroinformatics* doi:10.3389/neuro.11.
 443 011.2008
- Demin, V. and Nekhaev, D. (2018). Recurrent spiking neural network learning based on a competitive
 maximization of neuronal activity. *Frontiers in Neuroinformatics* doi:10.3389/fninf.2018.00079
- 446 Dhoble, K., Nuntalid, N., Indiveri, G., and Kasabov, N. (2012). Online spatio-temporal pattern recognition
 447 with evolving spiking neural networks utilising address event representation, rank order, and temporal
- spike learning. In *Proceedings of the International Joint Conference on Neural Networks*. doi:10.1109/
 IJCNN.2012.6252439
- Diehl, P. U. and Cook, M. (2015). Unsupervised learning of digit recognition using spike-timing-dependent
 plasticity. *Frontiers in Computational Neuroscience* doi:10.3389/fncom.2015.00099
- 452 Donahue, J., Hendricks, L. A., Rohrbach, M., Venugopalan, S., Guadarrama, S., Saenko, K., et al.
 453 (2017). Long-Term Recurrent Convolutional Networks for Visual Recognition and Description. *IEEE*454 *Transactions on Pattern Analysis and Machine Intelligence* doi:10.1109/TPAMI.2016.2599174
- 455 Dosovitskiy, A., Fischery, P., Ilg, E., Hausser, P., Hazirbas, C., Golkov, V., et al. (2015). FlowNet: Learning
 456 optical flow with convolutional networks. In *Proceedings of the IEEE International Conference on*457 *Computer Vision*. doi:10.1109/ICCV.2015.316
- 458 Droeschel, D., Stückler, J., and Behnke, S. (2011). Learning to interpret pointing gestures with a
- time-of-flight camera. In *HRI 2011 Proceedings of the 6th ACM/IEEE International Conference on Human-Robot Interaction*. doi:10.1145/1957656.1957822
- 461 Fan, Y., Lu, X., Li, D., and Liu, Y. (2016). Video-Based emotion recognition using CNN-RNN and C3D
 462 hybrid networks. In *ICMI 2016 Proceedings of the 18th ACM International Conference on Multimodal*463 *Interaction.* doi:10.1145/2993148.2997632
- 464 [Dataset] Feng, J. (2001). Is the integrate-and-fire model good enough? A review. doi:10.1016/
 465 S0893-6080(01)00074-0
- Feng, J. and Brown, D. (2000). Integrate-and-fire models with nonlinear leakage. *Bulletin of Mathematical Biology* doi:10.1006/bulm.1999.0162
- 468 Frati, V. and Prattichizzo, D. (2011). Using Kinect for hand tracking and rendering in wearable haptics. In
- 469 2011 IEEE World Haptics Conference, WHC 2011. doi:10.1109/WHC.2011.5945505

- Furber, S. B., Galluppi, F., Temple, S., and Plana, L. A. (2014). The SpiNNaker project. *Proceedings of the IEEE* doi:10.1109/JPROC.2014.2304638
- [Dataset] Furber, S. B., Lester, D. R., Plana, L. A., Garside, J. D., Painkras, E., Temple, S., et al. (2013).
 Overview of the SpiNNaker system architecture. doi:10.1109/TC.2012.142
- 474 Gerstner, W. (2008). Spike-response model. Scholarpedia doi:10.4249/scholarpedia.1343
- 475 Gerstner, W. (2009). Spiking Neuron Models. In *Encyclopedia of Neuroscience*. doi:10.1016/
 476 B978-008045046-9.01405-4
- Gerstner, W., Kistler, W. M., Naud, R., and Paninski, L. (2014). *Neuronal dynamics: From single neurons to networks and models of cognition*. doi:10.1017/CBO9781107447615
- Han, B. and Taha, T. M. (2010). Acceleration of spiking neural network based pattern recognition on
 NVIDIA graphics processors. *Applied Optics* doi:10.1364/AO.49.000B83
- Haodong Yang, Jun Zhang, Shuohao Li, J. L. and Chen, S. (2018). Attend it again: Recurrent attention
 convolutional neural network for action recognition. *Applied Sciences* 8, 383
- Haria, A., Subramanian, A., Asokkumar, N., Poddar, S., and Nayak, J. S. (2017). Hand Gesture Recognition
 for Human Computer Interaction. In *Procedia Computer Science*. doi:10.1016/j.procs.2017.09.092
- Hinton, G. E., Srivastava, N., and Swersky, K. (2012). Lecture 6a- overview of mini-batch gradient descent.
 COURSERA: Neural Networks for Machine Learning
- Hinz, G., Chen, G., Aafaque, M., Röhrbein, F., Conradt, J., Bing, Z., et al. (2017). Online Multi-object
 Tracking-by-Clustering for Intelligent Transportation System with Neuromorphic Vision Sensor. In
- 489 Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and
- 490 *Lecture Notes in Bioinformatics*). doi:10.1007/978-3-319-67190-1_11
- Hochreiter, S. and Schmidhuber, J. (1997). Long Short-Term Memory. *Neural Computation* doi:10.1162/
 neco.1997.9.8.1735
- Hodgkin, A. L. and Huxley, A. F. (1990a). A quantitative description of membrane current and its
 application to conduction and excitation in nerve. *Bulletin of Mathematical Biology* doi:10.1007/
 BF02459568
- Hodgkin, A. L. and Huxley, A. F. (1990b). A quantitative description of membrane current and its
 application to conduction and excitation in nerve. *Bulletin of Mathematical Biology* doi:10.1007/
 BF02459568
- Ji, S., Xu, W., Yang, M., and Yu, K. (2013). 3D Convolutional neural networks for human action recognition.
 IEEE Transactions on Pattern Analysis and Machine Intelligence doi:10.1109/TPAMI.2012.59
- Jiang, Z., Xia, P., Huang, K., Stechele, W., Chen, G., Bing, Z., et al. (2019). Mixed frame-/event-driven
 fast pedestrian detection. In *Proceedings IEEE International Conference on Robotics and Automation*.
 doi:10.1109/ICRA.2019.8793924
- Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., and Li, F. F. (2014). Large-scale
 video classification with convolutional neural networks. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. doi:10.1109/CVPR.2014.223
- Kasabov, N., Dhoble, K., Nuntalid, N., and Indiveri, G. (2013). Dynamic evolving spiking neural networks
 for on-line spatio- and spectro-temporal pattern recognition. *Neural Networks* doi:10.1016/j.neunet.
 2012.11.014
- 510 Kheradpisheh, S. R., Ganjtabesh, M., Thorpe, S. J., and Masquelier, T. (2018). STDP-based spiking deep
- convolutional neural networks for object recognition. *Neural Networks* doi:10.1016/j.neunet.2017.12.
 005
- 513 Kingma, D. P. and Ba, J. L. (2015). Adam: A method for stochastic optimization. In 3rd International
- 514 Conference on Learning Representations, ICLR 2015 Conference Track Proceedings

515 516	Krizhevsky, A., Sutskever, I., and Hinton, G. E. (2012). ImageNet classification with deep convolutional neural networks. In <i>Advances in Neural Information Processing Systems</i>
517	Liang, R. H. and Ouhvoung, M. (1998). A real-time continuous gesture recognition system for sign langu-
518	age. In Proceedings - 3rd IEEE International Conference on Automatic Face and Gesture Recognition.
519	<i>FG</i> 1998 doi:10 1109/AFGR 1998 671007
520	Lichtsteiner, P., Posch, C., and Delbruck, T. (2008a). A 128 × 128 120 dB 15 µs latency asynchronous
521	temporal contrast vision sensor. <i>IEEE Journal of Solid-State Circuits</i> doi:10.1109/JSSC.2007.914337
522	Lichtsteiner, P., Posch, C., and Delbruck, T. (2008b). A 128 × 128 120 dB 15 µs latency asynchronous
523	temporal contrast vision sensor. <i>IEEE Journal of Solid-State Circuits</i> 43, 566–576. doi:10.1109/JSSC.
524	2007.914337
525	Liu, H. and Wang, L. (2018). Gesture recognition for human-robot collaboration: A review. International
526	Journal of Industrial Ergonomics doi:10.1016/j.ergon.2017.02.004
527	Liu, Y. H. and Wang, X. J. (2001). Spike-frequency adaptation of a generalized leaky integrate-and-fire
528	model neuron. Journal of Computational Neuroscience doi:10.1023/A:1008916026143
529	Loiselle, S., Rouat, J., Pressnitzer, D., and Thorpe, S. (2006). Exploration of rank order coding with
530	spiking neural networks for speech recognition. doi:10.1109/ijcnn.2005.1556220
531	Majd, Mahshid and Safabakhsh, R. (2019). A motion-aware ConvLSTM network for action recognition.
532	Applied Intelligence, 1—-7
533	Mitra, S. and Acharya, T. (2007). Gesture recognition: A survey. IEEE Transactions on Systems, Man and
534	Cybernetics Part C: Applications and Reviews doi:10.1109/TSMCC.2007.893280
535	Mohemmed, A., Schliebs, S., Matsuda, S., and Kasabov, N. (2012). Span: Spike pattern association
536	neuron for learning spatio-temporal spike patterns. International Journal of Neural Systems doi:10.1142/
537	S0129065712500128
538	Nair, V. and Hinton, G. E. (2010). Rectified linear units improve Restricted Boltzmann machines. In ICML
539	2010 - Proceedings, 27th International Conference on Machine Learning
540	Nekhaev, D. and Demin, V. (2020). Competitive maximization of neuronal activity in convoluti-
541	onal recurrent spiking neural networks. In Studies in Computational Intelligence. doi:10.1007/
542	978-3-030-30425-6_30
543	Pérez-Carrasco, J. A., Serrano, C., Acha, B., Serrano-Gotarredona, T., and Linares-Barranco, B. (2010).
544	Spike-based convolutional network for real-time processing. In Proceedings - International Conference
545	on Pattern Recognition. doi:10.1109/ICPR.2010.756
546	Pigou, L., Dieleman, S., Kindermans, P. J., and Schrauwen, B. (2015). Sign language recognition using
547	convolutional neural networks. In <i>Lecture Notes in Computer Science (including subseries Lecture Notes</i>
548	in Artificial Intelligence and Lecture Notes in Bioinformatics). doi:10.1007/978-3-319-16178-5_40
549	Posch, C., Matolin, D., and Wohlgenannt, R. (2011a). A QVGA 143 dB dynamic range frame-free PWM
550	image sensor with lossless pixel-level video compression and time-domain CDS. In <i>IEEE Journal of</i>
551	Solid-State Circuits. doi:10.1109/JSSC.2010.2085952
552	Posch, C., Matolin, D., and Wohlgenannt, R. (2011b). A QVGA 143 dB dynamic range frame-free PWM
553	image sensor with lossless pixel-level video compression and time-domain CDS. In <i>IEEE Journal of</i> Solid State Circuita, doi:10.1100/ISSC.2010.2025052
554 555	Sour-State Circuits. doi:10.1109/JSSC.2010.2063952 Poutoroy S. S. and Agrowal A. (2012). Vision based hand gesture responsition for human commuter
555	interaction: a survey Artificial Intelligence Review doi:10.1007/s10462-012-0356-0
000	$MU_1U_1U_1U_1$, a survey, <i>Internetwork of the New Works</i> $MU_1U_1U_1U_1U_1U_2U_12U_2U_12U_2U_2U_2U_2U_2U_2U_2U_2U_2U_2U_2U_2U_2$

- interaction: a survey. *Artificial Intelligence Review* doi:10.1007/s10462-012-9356-9
 Shi, X., Chen, Z., Wang, H., Yeung, D. Y., Wong, W. K., and Woo, W. C. (2015). Convolutional LSTM
- network: A machine learning approach for precipitation nowcasting. In Advances in Neural Information
- 559 Processing Systems

- Shrestha, S. B. and Orchard, G. (2018). Slayer: Spike layer error reassignment in time. In *Advances in Neural Information Processing Systems*
- Song, H., Wang, W., Zhao, S., Shen, J., and Lam, K. M. (2018). Pyramid Dilated Deeper ConvLSTM for
 Video Salient Object Detection. In *Lecture Notes in Computer Science (including subseries Lecture Notes*
- *in Artificial Intelligence and Lecture Notes in Bioinformatics*). doi:10.1007/978-3-030-01252-6_44

Srivastava, N., Mansimov, E., and Salakhutdinov, R. (2015). Unsupervised learning of video representations
 using LSTMs. In *32nd International Conference on Machine Learning, ICML 2015*

- Tan, K. and Wang, D. L. (2018). A convolutional recurrent neural network for real-time speech enhancement.
 In *Proceedings of the Annual Conference of the International Speech Communication Association, INTERSPEECH*. doi:10.21437/Interspeech.2018-1405
- Tavanaei, A. and Maida, A. (2017). Bio-inspired multi-layer spiking neural network extracts discriminative
 features from speech signals. In *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics*). doi:10.1007/978-3-319-70136-3-95
- 573 Teka, W., Marinov, T. M., and Santamaria, F. (2014). Neuronal Spike Timing Adaptation Described with
- a Fractional Leaky Integrate-and-Fire Model. *PLoS Computational Biology* doi:10.1371/journal.pcbi.
 1003526
- 576 Vreeken, J. (2002). Spiking neural networks, an introduction. Computing
- Wang, Q., Zhang, Y., Yuan, J., and Lu, Y. (2019a). Space-time event clouds for gesture recognition: From
 RGB cameras to event cameras. In *Proceedings 2019 IEEE Winter Conference on Applications of Computer Vision, WACV 2019.* 1826–1835. doi:10.1109/WACV.2019.00199
- Wang, W., Hao, S., Wei, Y., Xiao, S., Feng, J., and Sebe, N. (2019b). Temporal Spiking Recurrent Neural
 Network for Action Recognition. *IEEE Access* doi:10.1109/access.2019.2936604
- Wang, X., Gao, L., Song, J., and Shen, H. (2017a). Beyond Frame-level CNN: Saliency-Aware 3-D CNN
 with LSTM for Video Action Recognition. *IEEE Signal Processing Letters* doi:10.1109/LSP.2016.
 2611485
- Wang, X., Gao, L., Song, J., and Shen, H. (2017b). Beyond Frame-level CNN: Saliency-Aware 3-D
 CNN with LSTM for Video Action Recognition. *IEEE Signal Processing Letters* doi:10.1109/LSP.2016.
 2611485
- Wang, Lebo and Li, Kaiming and Chen, Xu and Hu, X. P. (2019). Application of Convolutional Recurrent
 Neural Network for Individual Recognition Based on Resting State fMRI Data. *Frontiers in Neuroscience* 13, 434. doi:10.3389/fnins.2019.00434
- Werbos, P. J. (1990). Backpropagation Through Time: What It Does and How to Do It. *Proceedings of the IEEE* doi:10.1109/5.58337
- Wickeroth, D., Benölken, P., and Lang, U. (2009). Markerless gesture based interaction for design
 review scenarios. In 2nd International Conference on the Applications of Digital Information and Web
 Technologies, ICADIWT 2009. doi:10.1109/ICADIWT.2009.5273873
- Wysoski, S. G., Benuskova, L., and Kasabov, N. (2010). Evolving spiking neural networks for audiovisual
 information processing. *Neural Networks* doi:10.1016/j.neunet.2010.04.009
- Yang, R., Sarkar, S., and Loeding, B. (2010). Handling movement epenthesis and hand segmentation ambi guities in continuous sign language recognition using nested dynamic programming. *IEEE Transactions on Pattern Analysis and Machine Intelligence* doi:10.1109/TPAMI.2009.26
- 601 Zhou, K., Zhu, Y., and Zhao, Y. (2018). A spatio-temporal deep architecture for surveillance event
- detection based on ConvLSTM. In 2017 IEEE Visual Communications and Image Processing, VCIP
 2017. doi:10.1109/VCIP.2017.8305063