

Ultra-low Latency Spiking Neural Networks with Spatio-Temporal Compression and Synaptic Convolutional Block

Changqing Xu^{a,b,*}, Yi Liu^b, Yintang Yang^b

^a*Guangzhou Institute of Technology, Xidian University, Xi'an 710071, China.*

^b*School of Microelectronics, Xidian University, Xi'an 710071, China.*

Abstract

Spiking neural networks (SNNs), as one of the brain-inspired models, has spatio-temporal information processing capability, low power feature, and high biological plausibility. The effective spatio-temporal feature makes it suitable for event streams classification. However, neuromorphic datasets, such as N-MNIST, CIFAR10-DVS, DVS128-gesture, need to aggregate individual events into frames with a new higher temporal resolution for event stream classification, which causes high training and inference latency. In this work, we proposed a spatio-temporal compression method to aggregate individual events into a few time steps of synaptic current to reduce the training and inference latency. To keep the accuracy of SNNs under high compression ratios, we also proposed a synaptic convolutional block to balance the dramatic changes between adjacent time steps. And multi-threshold Leaky Integrate-and-Fire (LIF) models with learnable membrane time constants are introduced to increase its information processing capability. We evaluate the proposed method for event streams classification tasks on neuromorphic N-MNIST, CIFAR10-DVS, DVS128 gesture datasets. The experiment results show that our proposed method outperforms the state-of-the-art accuracy on nearly all datasets, using fewer time steps.

Keywords: Leaky Integrate-and-Fire model, multi-threshold, spatio-temporal

*Corresponding author

Email address: cqxu@xidian.edu.cn (Changqing Xu)

1. Introduction

Inspired by human neurons' working patterns, spiking neural networks (SNNs) are considered as the third generation artificial neural network[1]. With the development of SNNs, a large range of applications have been demonstrated including image classification [2][3], video processing [4] [5], posture and gesture recognition[6][7], voice recognition[8] [9]. Compared with traditional artificial neural networks (ANNs) which consist of static and continuous-valued neuron models, spiking neural networks (SNNs) have a unique event-driven computation characteristic that can respond to the events in a nearly latency-free and power-saving way[10][11], and it is naturally more suitable for processing event stream class.

To take full advantage of the event-driven advantages of spiking neural networks, neuromorphic sensors, such as DVS (Dynamic Vision Sensor)[12][13], Asynchronous Time-based Image Sensor(ATIS)[14], are usually used to transform datasets into neuromorphic datasets by encoding the time, location, and polarity of the brightness change. However, the event streams recorded by neuromorphic sensor based cameras are usually redundant in the temporal dimension, which is caused by high temporal resolution and irregular dynamic scene changes[15]. This characteristic makes event streams almost impossible to be processed directly by deep spiking neural networks, which are based on dense computation.

To make neuromorphic datasets suitable for deep spiking neural networks, variable pre-processing methods are proposed. The event-to-frame integrating method for pre-processing neuromorphic datasets is widely used. In [16], events are split into N slices with nearly the same number of events in each slice and integrate events to frames. The event-to-frame integrating method can convert the event streams into tens of frames at most. Otherwise, the accuracy of the SNNs will drop significantly. In [7], a temporal compression method is proposed

which can reduce the length of event streams by shrinking the duration of the input event trains. However, this method is only applied to the trained SNNs, which limits its potential. In [17], the normalized pixels of the static pictures are taken directly as the input current and multi-threshold neuron models are applied, which makes the SNNs can obtain a good performance in only two time steps. This method is suitable for static image classification, but it is difficult to directly apply it to event stream classification.

In this paper, we proposed a spatio-temporal compression method to aggregate event streams into few time steps of synaptic current to reduce the training and inference latency. To keep the accuracy of SNNs under high compression ratios, we also proposed a synaptic convolutional block in which a synaptic layer is applied to balance the dramatic change between adjacent time steps. To increase the information processing capability of neuron models, parametric multi-threshold Leaky Integrate-and-Fire models is introduced in our SNNs. We evaluate our method on neuromorphic datasets, such as N-MNIST[14], CIFAR10-DVS[13], DVS128-gesture[12] and experimental results show that our method outperforms the state-of-the-art accuracy on nearly all tested datasets using fewer time-steps.

2. Approach

In this section, we first introduce the proposed spatio-temporal compression method and synaptic convolutional block in Sec. 2.1 and 2.2. Then multi-threshold Leaky Integrate-and-Fire models with learnable membrane constants are introduced in Sec. 2.3. At last, we describe the proposed network structure of SNNs.

2.1. Spatio-temporal compression method

The spatio-temporal compression block, which is shown in Fig. 1, is used to pre-processing neuromorphic datasets, such as N-MNIST[14], DVS128 Gesture[12], CIFAR10-DVS[13], etc. neuromorphic datas are usually in the formulation of

$E(x_i, y_i, t_i, p_i)$ that represent the event’s coordinate, time and polarity. We split the event’s time T_{event} into T slices with nearly the same time interval in each slice and integrate these events. Note that T is also the simulating time-step. As Fig.1 shows, in each slice, spiking events are evenly divided into N_r parts in the temporal dimension. N_r is the resolution that user can change based on their requirement. Those spiking events in the same part will be integrated firstly and multiplied by the weights as the formulation (1) is shown. Due to the spiking events having two polarities, two channels frame is used to represent the compressed neuromorphic data.

$$\begin{aligned}
 j_l &= \text{floor}\left(\frac{\text{max}(T_{event})}{T}\right) \cdot j \\
 j_u &= \begin{cases} \text{floor}\left(\frac{\text{max}(T_{event})}{T}\right) \cdot (j + 1), & j < T - 1 \\ \text{max}(T_{event}), & j = T - 1 \end{cases} \\
 F(j, p, x, y) &= \sum_{k=0}^{N_r-1} (2^k \left(\sum_{i=j_l + \text{floor}\left(k \frac{j_u - j_l}{N_r}\right)}^{j_l + \text{floor}\left((k+1) \frac{j_u - j_l}{N_r}\right)} I_{E(x,y,t,p)}(x_i, y_i, t_i, p_i)\right))
 \end{aligned} \tag{1}$$

where j_l and j_u are the lower and upper bounds of the j_{th} slices, respectively, $\text{floor}()$ is the function that rounds the elements to the nearest integers towards minus infinity and $I_{E(x,y,t,p)}(x_i, y_i, t_i, p_i)$ is an indicator function of the event which equals to 1 only when $(x, y, t, p) = (x_i, y_i, t_i, p_i)$

2.2. Synaptic convolutional block

Since the synaptic convolution layer is used to balance the dramatic change between the adjacent time steps of the compressed neuromorphic data, the synaptic convolution block is used to replace the first convolution layer. The synaptic convolution block consists of a convolution layer, a synaptic layer and an optional average pooling layer, which is shown in Fig. 2. Whether there is an average pooling layer depends on the network structure.

The key of the synaptic convolutional block is the synaptic layer. In the

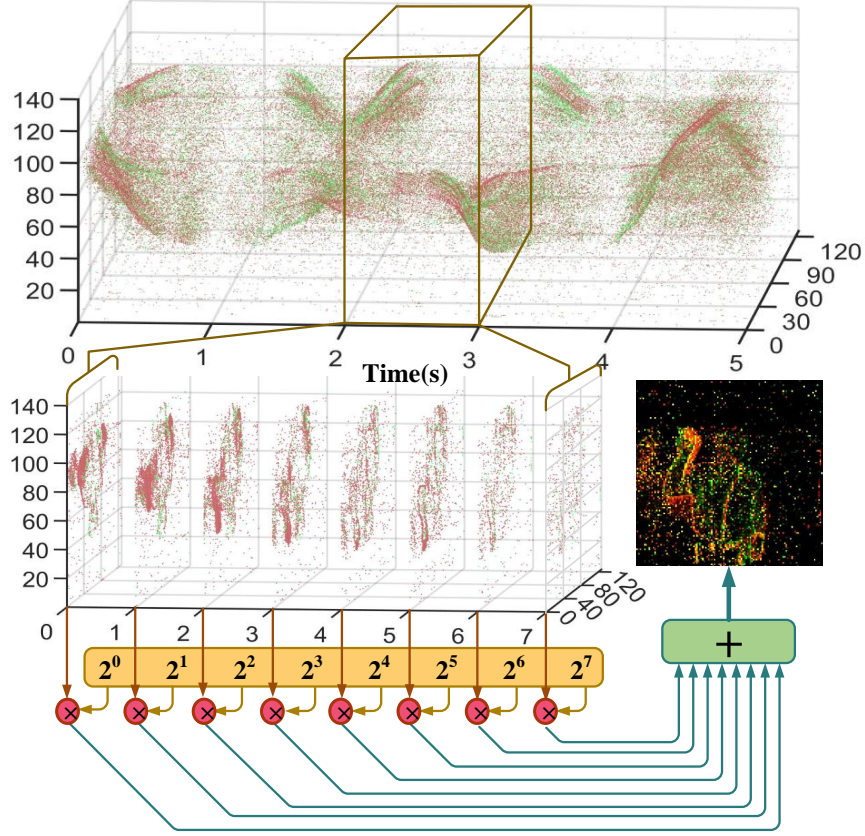


Figure 1: spatio-temporal compression block

synaptic layer, the first-order synaptic model is applied, which is shown below.

$$I_{syn}(t) = e^{-\frac{1}{\tau_{syn}} t} I_{syn}(t-1) + I_{in}(t) \quad (2)$$

where $I_{syn}(t)$ is the output current of the synapse at time t , τ_{syn} is the time constant of synapse and I_{in} is the input current at time t . To facilitate the calculation and simulation, we convert Eq. (2) to

$$I_{syn}[t_k] = \left(1 - \frac{1}{\tau_{syn}}\right) I_{syn}[t_{k-1}] + I_{in}[t_k] \quad (3)$$

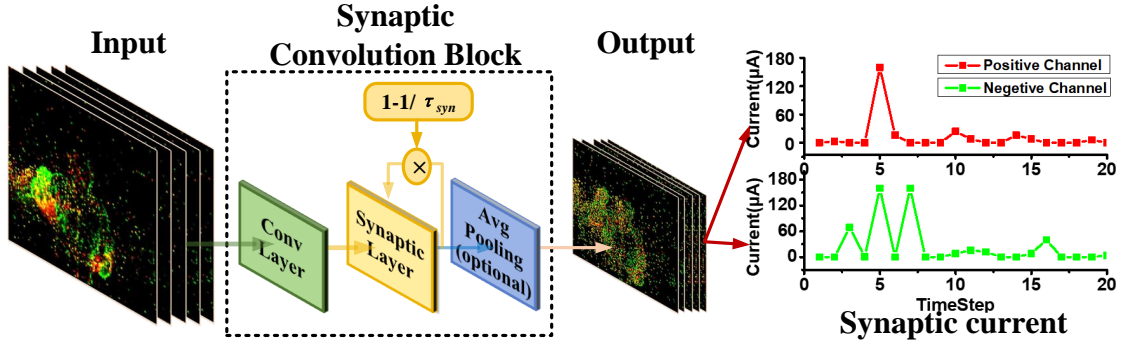


Figure 2: Synaptic convolution block

where $I_{syn}[t_k]$ represents the output current of the synapse at time step $t_k - 1$. τ_{syn} can control how much synaptic current at time step t_{k-1} can be retained at the time step t_k . To make the synaptic layer adjust adaptively according to the characteristics of datasets, τ_{syn} can be calculated based on

$$\tau_{syn}[t_k] = \frac{1}{1 - \frac{C_{valid}[t_k]}{C_{total}}} \quad (4)$$

where $C_{valid}[t_k]$ denotes the number of channels whose synaptic current is not zero at time step t_k . C_{total} is the total number of channels. As Eq. (4) is shown, the higher the value of $C_{valid}[t_k]$ is, the higher the value of $(1 - \frac{1}{\tau_{syn}})$ is, which means more information at time step t_{k-1} will be retained at time t_k .

2.3. Parametric Multi-threshold Leaky Integrate-and-Fire models(PMLIF)

It is known that Leaky Integrate-and Fire (LIF) model is one of the most widely applied models to describe the neuronal dynamics in SNNs. In this paper, we introduce the parametric multi-threshold Leaky Integrate-and-Fire models (PMLIF), whose membrane constants are learnable, to increase the information processing capability of neuron models. The neuronal membrane potential of neuron i at time t is

$$\tau_m \frac{du_i(t)}{dt} = -u_i(t) + I(t) + u_{reset}(t) \quad (5)$$

Where $I(t)$ is the pre-synaptic input current at time t and τ_m is a time constant of membrane voltage. $u_{reset}(t)$ denotes the reset function, which reduces the membrane potential by a certain amount V_{th} after the neuron i fires. The pre-synaptic input $I(t)$ is given by

$$I(t) = \sum_{j=1}^N \omega_{ij} s_j(t) \quad (6)$$

Where $s_j(t)$ is the output spike of pre-synaptic neuron j at time t and ω_{ij} is the presynaptic weight from the neuron j in the pre-synaptic layer to the neuron i in the post-synaptic layer. Due to the discrete time steps in the simulation, we apply the fixed-step first-order Euler method to discretize (7) to

$$u_i[t] = (1 - \frac{1}{\tau_m})u_i[t - 1] + I[t] + u_{reset}[t] \quad (7)$$

Where $u_{reset}[t]$ is equal to $-s_i[t]V_{th}$ and $s_i[t]$ is the output spike of neuron i . In this paper, we extend the LIF model into multi-threshold LIF model, in which the output of the neuron i can be expressed by

$$s_i[t] = \begin{cases} 0, & u_i[t] < V_{th} \\ \text{floor}(\frac{u_i[t]}{V_{th}}), & V_{th} \leq u_i[t] < S_{max}V_{th} \\ s_{max}, & u_i[t] \geq S_{max}V_{th} \end{cases} \quad (8)$$

Where S_{max} is the upper limit of the output spikes and $\text{floor}()$ is the function that rounds the elements to the nearest integers towards minus infinity. Since τ_m is a learnable parameter and its value should be positive, we use a sigmoid function about synaptic time constant weight ω_{m_i} of the neuron i to replace the $(1 - \frac{1}{\tau_m})$ and the (7) becomes

$$u_i[t] = S(\omega_{m_i})u_i[t - 1] + I[t] + u_{reset}[t] \quad (9)$$

where $S(\omega_{m_i})$ is the sigmoid function about ω_{m_i} .

2.4. Error Backpropagation of PMLIF

To present the error backpropagation of PMLIF, we define the loss function $L[t_k]$ in which the mean square error for each output neuron at time step t_k .

$$L[t_k] = \frac{1}{2} \sum_{i=0}^{N_o} (y_i[t_k] - s_i[t_k])^2 \quad (10)$$

Where N_o is the number of neurons in the output layer, $y_i[t_k]$ and $s_i[t_k]$ denotes the desired and the actual firing event of neurons i in the output layer at time step t_k . By combining (5)-(9), it can be seen that loss function $L[t_k]$ is a function of presynaptic weight $\omega_{i,j}$ and synaptic time constant weight ω_m . In this paper, we use $W^{(l)} = [\omega_1^{(l)}; \dots; \omega_{N_l}^{(l)}]$ to represent the presynaptic weight matrix of layer l in which $\omega_{N_l}^{(l)} = [\omega_{N_l,1}^{(l)}, \omega_{N_l,2}^{(l)}, \dots, \omega_{N_l,N_{l-1}}^{(l)}]$. $W_m^{(l)}$ denotes the synaptic time constant weight matrix, which is equal to $[\omega_{m_1}^{(l)}, \dots, \omega_{m_{N_l}}^{(l)}]$. The aim of error backpropagation is to update the presynaptic weight $W^{(l)}$ and synaptic time constant weight $W_m^{(l)}$ using the error gradient $\frac{\partial L[t_k]}{\partial W^{(l)}}$ and $\frac{\partial L[t_k]}{\partial W_m^{(l)}}$. Using the chain rule, the error gradient with the respect to the presynaptic weight $W^{(l)}$ in the layer l is

$$\frac{\partial L[t_k]}{\partial W^{(l)}} = \frac{\partial L[t_k]}{\partial u^{(l)}[t_k]} \frac{\partial u^{(l)}[t_k]}{\partial W^{(l)}} = \delta^{(l)}[t_k] \frac{\partial u^{(l)}[t_k]}{\partial W^{(l)}} \quad (11)$$

Where $\delta^{(l)}[t_k]$ is the back propagated error of layer l at time t_k , which is equal to $\frac{\partial L[t_k]}{\partial u^{(l)}[t_k]}$.

$$\begin{aligned} \delta^{(l)}[t_k] &= \frac{\partial u^{(l+1)}[t_k]}{\partial u^{(l)}[t_k]} \frac{\partial L[t_k]}{\partial u^{(l+1)}[t_k]} \\ &= \frac{\partial u^{(l+1)}[t_k]}{\partial s^{(l)}[t_k]} \frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]} \delta^{(l+1)}[t_k] \\ &= (W^{(l+1)})^T \frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]} \delta^{(l+1)}[t_k] \end{aligned} \quad (12)$$

The key to calculate $\delta^{(l)}[t_k]$ is to obtain $\frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]}$. Theoretically, $s^{(l)}[t]$ is a non-differentiable function and we cannot obtain the value of $\frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]}$ directly.

In this paper, we use an approximate curve in [17] to surrogate the original derivative function. The function of the approximate curve f_1 is shown below.

$$f_1(u) = \sum_{i=1}^{S_{\max}} \alpha_H e^{-(u-iV_{th})^2/\alpha_W} \quad (13)$$

where α_H and α_W determine the curve shape and steep degree. S_{max} is the upper limit of the output spikes. From (11), the second term $\frac{\partial u^{(l)}[t_k]}{\partial W^{(l)}}$ is given by

$$\begin{aligned} \frac{\partial u^{(l)}[t_k]}{\partial W^{(l)}} &= S(W_m^{(l)}) \frac{\partial u^{(l)}[t_{k-1}]}{\partial W^{(l)}} + s^{l-1}[t_k] - \frac{\partial s^{(l)}[t_k]}{\partial W^{(l)}} V_{th} \\ &= S(W_m^{(l)}) \frac{\partial u^{(l)}[t_{k-1}]}{\partial W^{(l)}} + s^{l-1}[t_k] - \frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]} \frac{\partial u^{(l)}[t_k]}{\partial W^{(l)}} V_{th} \end{aligned} \quad (14)$$

From (14), we can get

$$\frac{\partial u^{(l)}[t_k]}{\partial W^{(l)}} = \frac{S(W_m^{(l)}) \frac{\partial u^{(l)}[t_{k-1}]}{\partial W^{(l)}} + s^{l-1}[t_k]}{1 + \frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]} V_{th}} \quad (15)$$

The error gradient with the respect to the synaptic time constant weight $W_m^{(l)}$ can be calculated by

$$\frac{\partial L[t_k]}{\partial W_m^{(l)}} = \delta^{(l)}[t_k] \frac{\partial u^{(l)}[t_k]}{\partial W_m^{(l)}} \quad (16)$$

From (16), the second term $\frac{\partial u^{(l)}[t_k]}{\partial W_m^{(l)}}$ is given by

$$\frac{\partial u^{(l)}[t_k]}{\partial W_m^{(l)}} = \frac{\partial(S(W_m^{(l)})u^{(l)}[t_{k-1}])}{\partial W_m^{(l)}} - \frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]} \frac{\partial u^{(l)}[t_k]}{\partial W_m^{(l)}} V_{th} \quad (17)$$

From (18), we can obtain

$$\frac{\partial u^{(l)}[t_k]}{\partial W_m^{(l)}} = \frac{\frac{\partial(S(W_m^{(l)})u^{(l)}[t_{k-1}])}{\partial W_m^{(l)}}}{1 + \frac{\partial s^{(l)}[t_k]}{\partial u^{(l)}[t_k]} V_{th}} \quad (18)$$

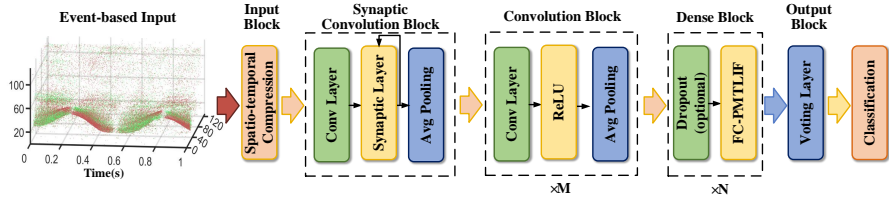


Figure 3: Architecture of the proposed SNN

2.5. Network architecture

We proposed a flexible network structure for event stream classification tasks. The proposed network structure is illustrated in Fig. 3. The network consists of a spatio-temporal compression layer, a synaptic convolution block, M end-to-end connected convolution blocks, N end-to-end connected dense blocks, and a voting layer. Since the output of the synaptic convolution block are synaptic currents which are real values, not binary values, we use the convolution block which consists of sequentially connected layers including a convolution layer, a ReLU layer, and an average pooling layer, to extracted features, directly. In dense blocks, there are a dropout layer and a fully connected layer which consists of PMLIFs. A voting layer after the last dense block is used to boost classifying robustness. The output of the last dense block will be divided into N_{class} groups, randomly, and connected to the voting layer, where N_{class} is the number of the classes. The voting layer is implemented by calculating the average value of each group and selecting the maximum value as the classification result.

3. Experiments and results

We test our proposed SNN model and training method on three neuromorphic datasets N-MNIST[14], DVS128 Gesture[12], CIFAR10-DVS[13] with different sizes and structures of SNNs. And we compared our training method with several previously reported state-of-the-art results with the same or similar networks including different SNNs trained by BP-based methods, converted

Table 1: Parameters setting

Parameters	Description	Value
V_{th}	Threshold	10 mV
α_m	Derivative approximation parameters	1
α_W	Derivative approximation parameters	20
N_r	Resolution of spatio-temporal compression	8
S_{max}	Upper limit of output spikes	15
N_{Batch}	Batch Size (N-MNIST/DVS128/CIFAR10-DVS)	64,16,16
η	Learning rate (N-MNIST/DVS128/CIFAR10-DVS)	0.0002, 0.00004, 0.0001
$\beta_1, \beta_2, \lambda$	Adam parameters	0.9, 0.999, $1 - 10^{-8}$

SNNs, and traditional ANNs.

3.1. Experiment settings

All reported experiments below are conducted on an NVIDIA Tesla V100 GPU. The implementation of our proposed method is on the Pytorch framework[18]. The experimented SNNs are based on the network structure described in Sec. 2.5. Only 2-5 time steps are used to demonstrate the proposed ultra low-latency spiking neural network. No refractory period is used. Adam [19] is applied as the optimizer. If not otherwise specified, the accuracy in this paper refers to the best results obtained by repeating the experiments five times.

The initialization of parameters, such as the weights, threshold, time constant of membrane voltage and synapse, and other parameters, directly affect the convergence speed and stability of the whole network. We should simultaneously make sure enough spikes transmit information between neural network layers and avoid too many spikes that reduce the neuronal selectivity. In this paper, we use a fixed threshold in each neuron for simplification and initialize

Table 2: Network structure

Dataset	Network structure
N-MNIST	128SC3-128C3-AP2-256C3-AP2 -512C3-AP4-DP-512FC-10Voting
DVS128	32SC3-32C3-AP2-64C3-AP2-128C3-AP2 -256C3-AP2-512C3-AP4-DP-512FC-11Voting
Cifar10-DVS	32C3-32C3-AP2-64C3-AP2-128C3-AP2 -256C3-AP2-512C3-AP4-DP-512FC-10Voting

128SC3 represents synaptic convolution block with 128 3×3 filters. 128C3 represents convolution block with 128 3×3 filters. AP2 represents average pooling layer with 2×2 filters. DP denotes dropout layer and 512FC means a fully connected layer that consists of 512 PMLIFs.

the weight $W^{(l)}$ parameters sampling from the normal distribution.

$$W^{(l)} \sim \left[\frac{V_{th}}{N_{l-1}}, 0.5 \right] \quad (19)$$

where V_{th} is the threshold of membrane voltage, N_{l-1} is the number of neurons of pre-layer. The synaptic time constant weights W_m are initialized to zeros. The set of other parameters is presented in Table 1. In addition, we do not apply complex skill, such as error normalization[20], weight regularization[21], warm-up mechanism [22], etc. All testing accuracy is reported after training 50 epochs in our experiments.

3.2. Dataset experiments

3.2.1. N-MNIST

The Neuromorphic-MNIST (N-MNIST) dataset is a spiking version of the original frame-based MNIST dataset, which consists of the same 60 000 training and 10 000 testing samples as the original MNIST dataset[14]. Each N-MNIST example is captured at the same visual scale as the original MNIST dataset (28x28 pixels). The N-MNIST dataset was captured by mounting the Asynchronous Time Based Image (ATIS) sensor on a motorized pan-tilt unit and

Table 3: Comparisons with SNNs on N-MNIST

Models	Method	Time Step	Epoch	ACC(%)
TSSL-BP[24]	SNN	30	100	99.23
LISNN[25]	SNN	20	100	99.45
NeuNorm SNN[26]	SNN	50	200	99.53
BackEISNN[27]	SNN	100	200	99.57
LMCSNN[23]	SNN	10	200	99.61
This work	SNN	2	50	99.63

having the sensor move while it views. For N-MNIST dataset, the network structure we applied is shown in Table 2. We compare our proposed network with several spiking convolutional neural networks which have similar network structures. Table 3 shows that our proposed spiking neural network can achieve 99.63% which outperforms other results. In addition to the performance improvement, our proposed method also has a large reduction of time step count. Compared with LMCSNN[23] which only has 10 time steps, our method still has 5 times improvement.

3.2.2. DVS128-gesture

The IBM DVS128-gesture[12] is an event-based gesture recognition dataset, which has the temporal resolution in μs level and 128×128 spatial resolution. It records 1342 samples of 11 gestures, such as hand clips, arm roll, etc., collected from 29 individuals under three illumination conditions, and each gesture has an average duration of 6 seconds. Compared with N-MNIST, DVS128-gesture is a more challenging dataset which has more temporal information. Table 4 shows that our method obtains 98.96% test accuracy with only 5 time steps in 50 epochs. BPSTA-SNN [28] uses the SNN which is pre-trained by TSSL-BP [24] to obtain the same accuracy with 16 time steps in 300 epochs.

Table 4: Comparisons with DNNs and SNNs on DVS128 Gesture

Models	Method	Time Step	Epoch	Trick	ACC(%)
STBP-TDBN[29]	SNN	40	-		96.87
RG-CNN[30]	DNN	8	150		97.20
LMCSNN[23]	SNN	20	200		97.57
STFilter[31]	DNN	12	-	Spatiotem- poral filters	97.75
BPSTA-SNN[28]	SNN	16	300	Pre-train by TSSL-BP	98.96
Our Method	SNN	5	50		98.96

Table 5: Comparisons with SNNs on Cifar10-DVS

Models	Method	Time Step	Epoch	ACC(%)
NeuNorm SNN[26]	SNN	100	200	60.5
SR-ANN[33]	ANN to SNN	60	-	66.75
STBP-TDBN[29]	SNN	40	-	67.8
LIAF-SNN[34]	SNN	10	60	70.4
TA-SNN[32]	SNN	10	150	72
Our proposed	SNN	5	50	73.8

3.2.3. Cifar10-DVS

To validate our method, we apply a deeper network structure which contains six convolution layers and five average pooling, and a dense layer on the dataset Cifar10-DVS[13] which is an important benchmark for comparison in SNN domains. Cifar10-DVS is a neuromorphic version converted from the Cifar10 dataset in which 10,000 frame-based images are converted into 10,000 event streams with DVS. Our proposed method obtains 73.8% test accuracy with 5 time steps. Compared with TA-SNN [32], our proposed method obtain 1.8% accuracy improvement with only a half of time steps. When we apply the 10 time steps, we can achieve 76.9% test accuracy in 50 epochs.

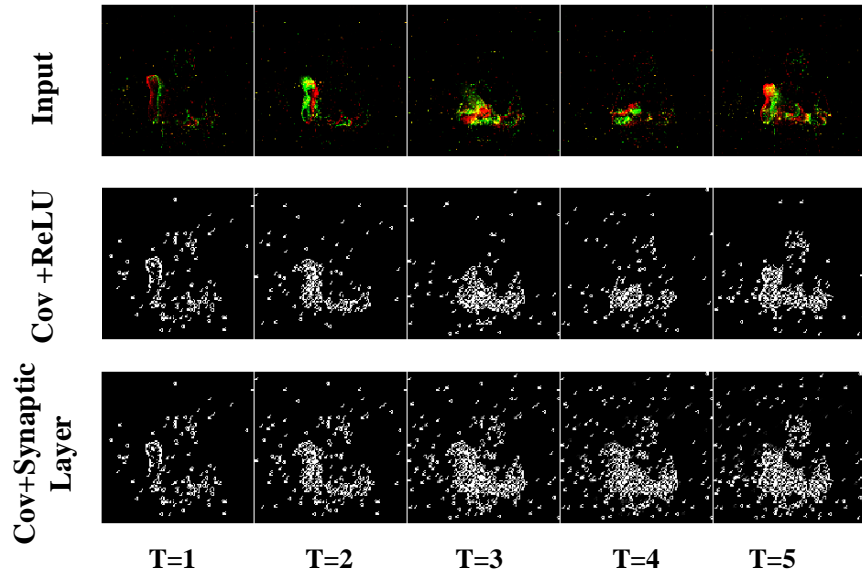


Figure 4: Different between Cov+ReLU layer and Cov+Synaptic Layer

3.3. Ablation Study

We conduct extensive ablation studies on DVS128-gesture to evaluate the validity of the method. The reason that we chose DVS128-gesture is that it is a more challenging dataset, especially for SNNs with few time steps, due to the more complex spatial structure and stronger temporal information. The function of the synaptic convolutional block is to balance the dramatic change between adjacent time steps. The dramatic change may slow down the convergence speed of the spiking neural network or even cause no converge during the training process. We compare the output of a synaptic convolutional block with the output of a convolution layer + ReLU layer. In Fig. 4, the figures in the first row are the input data at different time steps. Compared with figures in the second row, the figures in the third row are more continuous. It is shown that the main difference between Cov+ReLU and Cov+Synaptic layer is that the synaptic layer retains some information of the previous time step which avoids the dramatic changes between adjacent time steps.

Table 6: Ablation study of different strategies based on DVS128-gesture

Methods	T=2			T=5			T=10		
	Mean(%)	Std(%)	Best(%)	Mean(%)	Std(%)	Best(%)	Mean(%)	Std(%)	Best(%)
S0	95.60	0.40	96.52	98.47	0.31	98.96	98.75	0.28	98.96
S1	95.35	0.17	95.49	98.06	0.52	98.96	98.33	0.38	98.61
S2	94.79	0.31	95.14	97.99	0.40	98.61	98.13	0.28	98.61
S3	94.65	0.36	95.14	97.98	0.24	98.26	97.92	0.22	98.26

To evaluate the influence of synaptic convolutional blocks and PMLIF models, we design an ablation study of different strategies, which consist of four: S0, SNNs with synaptic convolutional blocks and PMLIF; S1: SNNs with only PMLIF; S2: SNNs with only synaptic convolutional blocks; S3: SNNs without synaptic convolutional blocks and PMLIF. As shown in table 6, The SNNs with Synaptic convolutional blocks and PMLIF have a higher testing accuracy and the performance improvement is more significant with the reduction of time steps. Compared with the SNN without synaptic convolutional block and PMLIF, the SNN with synaptic convolutional block and PMLIF has a 1.38% improvement on testing accuracy, when time step count is 2.

3.4. Performance analysis

3.4.1. Influence of compression ratio

To reduce the training and inference latency of SNNs, we proposed a spatio-temporal compression method to aggregate event streams into a few time steps. Since neuromorphic datasets will be integrated into tens or hundreds of frames in the reported works, the compression ratio is defined as the ratio of average frames or time steps used in the mentioned work to the time steps we applied. The lines in Fig. 5 denotes the mean of the accuracy obtained by repeating the experiments five times, the shade is the fluctuation range of the data.

As Fig. 5 (a) is shown, the accuracy still can keep 98.96% when compression ratio is 26.04% for DVS128-gesture. When compression ratio is 10.42%, classification accuracy drops from 98.96% to 95.83%. The reason for the significant accuracy degradation is the time steps are only two when the compression

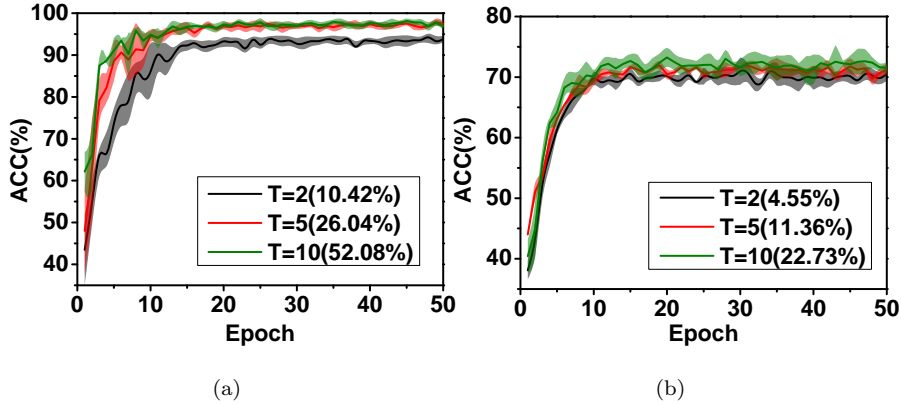


Figure 5: The test accuracy of different compression ratios on (a) DVS128-gesture and (b) Cifar10-DVS dataset

sion ratio is 10.42%, which causes too much temporal information to be lost. For Cifar10-DVS dataset, a higher compression ratio is applied. Compared with DVS128-gesture, the accuracy drop is not significant with the increase of the compression ratio. The reason for that is Cifar10-DVS is a neuromorphic dataset converted from the static dataset, which contains less temporal information, high compression ratio mainly causes temporal information loss rather than spatial information.

3.4.2. Distribution of τ_{mem} after training

In PMLIF model, we introduce the learnable membrane constants to improve the convergence speed and stability of the whole network. The results in section 3.3 have shown the benefit of applying the PMLIF model. In this section, we study the distribution of the membrane constants under different compression ratios. As Fig. 6 is shown, the distribution of membrane constants basically conformed to a Gaussian distribution. We can find that the mean of membrane constants increases with the increase of the compression ratio. The membrane time constants are used to control the retained information at the previous time step. A larger membrane constant means more information at the previous time step will be retained. The increase of compression ratio means that more

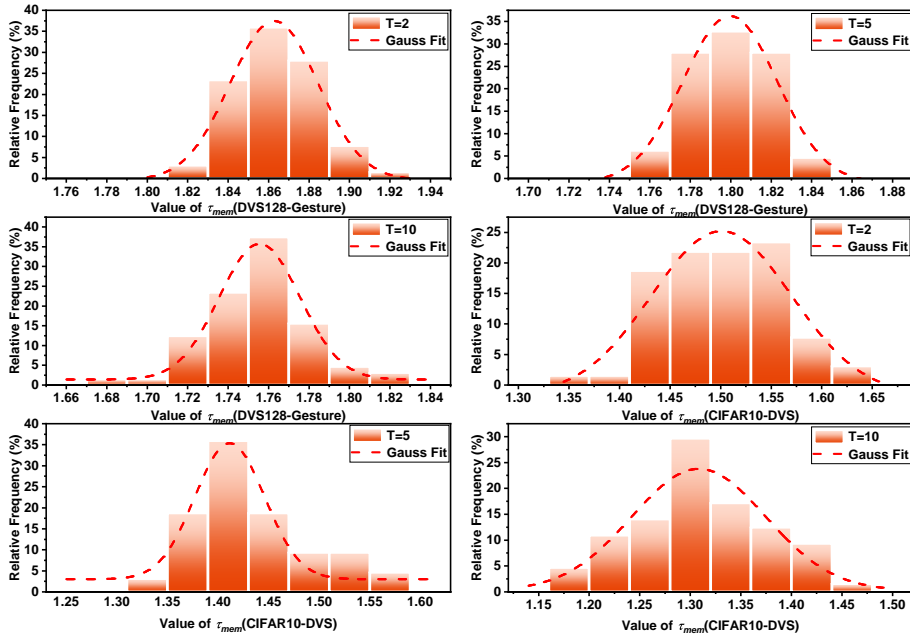


Figure 6: Distribution of τ_{mem} after training on DVS128-Gesture and CIFAR10-DVS

information will be integrated into one time step. The analysis result is that a larger membrane time constant should be applied to keep more information of previous time step retained when the compression ratio increases.

4. Conclusion

In this work, we proposed an ultra-low latency spiking neural network with spatio-temporal compression and a Synaptic convolutional block for event stream classification. The proposed spatio-temporal compression method is used to compress event streams into few time steps to reduce the training and inference latency. We also proposed a synaptic convolutional block as the first layer of the SNN, in which a synaptic layer is applied to balance the dramatic change between adjacent time steps. The parametric multi-threshold Leaky Integrate-and-Fire models, whose membrane constants are learnable, are introduced in our SNNs. We evaluate our proposed method and compare it with state-of-the-art methods. Experiment results show that our proposed method outperforms

state-of-the-art methods with the same or similar network architecture on neuromorphic datasets, such as N-MNIST, CIFAR10-DVS, DVS128-gesture.

Acknowledgment

This work was supported in part by the National Natural Science Foundation of China under Grant 62004146, by the China Postdoctoral Science Foundation funded project under Grant 2021M692498, and by the Fundamental Research Funds for the Central Universities.

References

- [1] W. Maass, Networks of spiking neurons: the third generation of neural network models, *Neural networks* 10 (9) (1997) 1659–1671.
- [2] T. Iakymchuk, A. Rosado-Muñoz, J. F. Guerrero-Martínez, M. Bataller-Mompeán, J. V. Francés-Víllora, Simplified spiking neural network architecture and stdp learning algorithm applied to image classification, *EURASIP Journal on Image and Video Processing* 2015 (1) (2015) 1–11.
- [3] G. Datta, S. Kundu, A. R. Jaiswal, P. A. Beerel, Hyper-snn: Towards energy-efficient quantized deep spiking neural networks for hyperspectral image classification, *arXiv preprint arXiv:2107.11979*.
- [4] G. E. Hinton, N. Srivastava, A. Krizhevsky, I. Sutskever, R. R. Salakhutdinov, Improving neural networks by preventing co-adaptation of feature detectors, *arXiv preprint arXiv:1207.0580*.
- [5] Y. Hu, H. Liu, M. Pfeiffer, T. Delbruck, Dvs benchmark datasets for object tracking, action recognition, and object recognition, *Frontiers in neuroscience* 10 (2016) 405.
- [6] B. Zhao, R. Ding, S. Chen, B. Linares-Barranco, H. Tang, Feedforward categorization on aer motion events using cortex-like features in a spiking

- neural network, *IEEE transactions on neural networks and learning systems* 26 (9) (2014) 1963–1978.
- [7] C. Xu, W. Zhang, Y. Liu, P. Li, Boosting throughput and efficiency of hardware spiking neural accelerators using time compression supporting multiple spike codes, *Frontiers in Neuroscience* 14 (2020) 104.
- [8] W. Zhang, P. Li, Spike-train level backpropagation for training deep recurrent spiking neural networks, *Advances in neural information processing systems* 32.
- [9] Y. Jin, W. Zhang, P. Li, Hybrid macro/micro level backpropagation for training deep spiking neural networks, *Advances in neural information processing systems* 31.
- [10] J. Pei, L. Deng, S. Song, M. Zhao, Y. Zhang, S. Wu, G. Wang, Z. Zou, Z. Wu, W. He, et al., Towards artificial general intelligence with hybrid tianjic chip architecture, *Nature* 572 (7767) (2019) 106–111.
- [11] K. Roy, A. Jaiswal, P. Panda, Towards spike-based machine intelligence with neuromorphic computing, *Nature* 575 (7784) (2019) 607–617.
- [12] A. Amir, B. Taba, D. Berg, T. Melano, J. McKinstry, C. Di Nolfo, T. Nayak, A. Andreopoulos, G. Garreau, M. Mendoza, et al., A low power, fully event-based gesture recognition system, in: *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7243–7252.
- [13] H. Li, H. Liu, X. Ji, G. Li, L. Shi, Cifar10-dvs: an event-stream dataset for object classification, *Frontiers in neuroscience* 11 (2017) 309.
- [14] G. Orchard, A. Jayawant, G. K. Cohen, N. Thakor, Converting static image datasets to spiking neuromorphic datasets using saccades, *Frontiers in neuroscience* 9 (2015) 437.
- [15] M. Yao, H. Gao, G. Zhao, D. Wang, Y. Lin, Z. Yang, G. Li, Temporal-wise attention spiking neural networks for event streams classification, in: *Pro-*

- ceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 10221–10230.
- [16] W. Fang, Y. Chen, J. Ding, D. Chen, Z. Yu, H. Zhou, Y. Tian, other contributors, Spikingjelly, <https://github.com/fangwei123456/spikingjelly>, accessed: 2022-2-24 (2020).
- [17] C. Xu, Y. Liu, Y. Yang, Direct training via backpropagation for ultra-low latency spiking neural networks with multi-threshold (2021). [arXiv: 2112.07426](https://arxiv.org/abs/2112.07426).
- [18] A. Paszke, S. Gross, F. Massa, A. Lerer, S. Chintala, Pytorch: An imperative style, high-performance deep learning library.
- [19] D. Kingma, J. Ba, Adam: A method for stochastic optimization, *Computer Science*.
- [20] P. U. Diehl, D. Neil, J. Binas, M. Cook, S. C. Liu, Fast-classifying, high-accuracy spiking deep networks through weight and threshold balancing, in: *International Joint Conference on Neural Networks*, 2015.
- [21] J. H. Lee, T. Delbruck, M. Pfeiffer, Training deep spiking neural networks using backpropagation, *Frontiers in neuroscience* 10 (2016) 508.
- [22] W. Zhang, P. Li, Temporal spike sequence learning via backpropagation for deep spiking neural networks, *Advances in Neural Information Processing Systems* 33 (2020) 12022–12033.
- [23] W. Fang, Z. Yu, Y. Chen, T. Masquelier, T. Huang, Y. Tian, Incorporating learnable membrane time constant to enhance learning of spiking neural networks, in: *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 2661–2671.
- [24] W. Zhang, P. Li, Spike-train level backpropagation for training deep recurrent spiking neural networks, *arXiv preprint arXiv:1908.06378*.

- [25] X. Cheng, Y. Hao, J. Xu, B. Xu, Lisnn: Improving spiking neural networks with lateral interactions for robust object recognition., in: IJCAI, 2020, pp. 1519–1525.
- [26] Y. Wu, L. Deng, G. Li, J. Zhu, Y. Xie, L. Shi, Direct training for spiking neural networks: Faster, larger, better, in: Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33, 2019, pp. 1311–1318.
- [27] D. Zhao, Y. Zeng, Y. Li, Backeisnn: A deep spiking neural network with adaptive self-feedback and balanced excitatory-inhibitory neurons, arXiv preprint arXiv:2105.13004.
- [28] G. Shen, D. Zhao, Y. Zeng, Backpropagation with biologically plausible spatio-temporal adjustment for training deep spiking neural networks, arXiv preprint arXiv:2110.08858.
- [29] H. Zheng, Y. Wu, L. Deng, Y. Hu, G. Li, Going deeper with directly-trained larger spiking neural networks, arXiv preprint arXiv:2011.05280.
- [30] Y. Bi, A. Chadha, A. Abbas, E. Bourtsoulatze, Y. Andreopoulos, Graph-based spatio-temporal feature learning for neuromorphic vision sensing, IEEE Transactions on Image Processing 29 (2020) 9084–9098.
- [31] R. Ghosh, A. Gupta, A. Nakagawa, A. Soares, N. Thakor, Spatiotemporal filtering for event-based action recognition, arXiv preprint arXiv:1903.07067.
- [32] M. Yao, H. Gao, G. Zhao, D. Wang, Y. Lin, Z. Yang, G. Li, Temporal-wise attention spiking neural networks for event streams classification, in: Proceedings of the IEEE/CVF International Conference on Computer Vision, 2021, pp. 10221–10230.
- [33] A. Kugele, T. Pfeil, M. Pfeiffer, E. Chicca, Efficient processing of spatio-temporal data streams with spiking neural networks, Frontiers in Neuroscience 14 (2020) 439.

- [34] Z. Wu, H. Zhang, Y. Lin, G. Li, M. Wang, Y. Tang, Liaf-net: Leaky integrate and analog fire network for lightweight and efficient spatiotemporal information processing, *IEEE Transactions on Neural Networks and Learning Systems*.