

Cascade Spiking Neuron Network For Event-based Image Classification In Noisy Environment

Yuntao Han, Tao Yu, Silu Cheng, and Jiangtao Xu, *Member, IEEE*

Abstract—Spiking Neuron Network (SNN) has shown advantages in processing event-based data for image classification. However, the classification accuracy of SNNs decreases in noisy environment. The cascade spiking neuron network (cascade-SNN) was proposed in this letter, which was a combination of spiking convolutional neuron network (SCNN) for features extraction and liquid state machine (LSM) for read-out. Spatio-temporal back propagation was used for training this network. Compared with early works on ANNs, this network achieved the state-of-the-art classification accuracy 76.70% in DVS-CIFAR10 dataset and 97.57% in DVS-Gesture dataset, which are both challenging dataset because of noisy environment. The results showed that the proposed structure could improve the classification accuracy of SNN in noisy environments. Ablation experiments were used to examine the cascade structure, showing that this structure could achieve higher classification accuracy but have lower convergence speed. Finally, Hyper-parameter analysis was applied to further explore the proposed network.

Index Terms—spiking neuron network, liquid state machine, image classification, neuromorphic data.

I. INTRODUCTION

WITH the development of the brain-inspired computing [1], event-based images captured by bio-inspired cameras [2] appeared. The event-based images store spatio-temporal information in spikes with address and time label, thus showing advantages in low time latency and high dynamic range [3]. Many approaches based on traditional ANNs have been proposed for processing event-based images, e.g. RNN [4], DART [5], HATS [6]. Spiking neural network (SNN) is the third generation of neuron network. It has the ability in processing event-based images with low power consumption [7] and high classification accuracy [4].

The basic calculation block in SNNs is a mathematical model of neurons. Leaky integrate-and-fire (LIF), as a neuron model [8], achieves a balance between implementation cost and biological credibility [9], and the hardware circuits for LIF neuron model has been studied by many works [10]. So far, LIF is the most frequently used neuron model in SNNs.

The structure of SNNs is usually inspired by the structure of classical convolutional neuron networks (CNNs). For using classical CNN structure (e.g. VGG [11], yolo-v3 [12]) in SNN, many conversion methods have been proposed for converting CNNs into spiking convolutional neuron networks (SCNNs).

Compared with the original models, the converted VGG architecture [13] [14] and Spiking-YOLO [15] achieved similar accuracy with low precision loss.

Liquid state machine (LSM) [16] is also a computation model widely used in SNNs. The structure and function of LSM are based on the neocortex in the central nervous system of mammals. LSM is composed of a large number of LIF neurons with cyclic structure. The existence of cyclic connections makes the network have the ability of memory, because the firing state contains information about current and historical inputs. LSM could contain the temporal information in both the membrane voltages and output spikes, thus utilizing temporal information better than other SNNs.

Like traditional ANNs, SNNs are facing the problem of classification accuracy decreasing [17] in complex noisy scenes. Many approaches has been proposed to promote the performance of SCNN and LSM like modifying ANNs before converting to SNN [18], temporal credit assignment policy [19] and threshold-dependent batch normalization (tdBN) [20].

To address this problem, the cascade spiking neuron network (cascade-SNN) was proposed in this letter. We used SCNN for features extraction and LSM for read-out. The structure of the proposed network was based on VGG model. Because SNNs have time-domain characteristic and the spike activity is non-differentiable, the effect of traditional back-propagation (BP) algorithm is limited. Spatio-temporal back-propagation (STBP) [21] [22] was applied to solving this problem. As far as we know, this was the first time to train LSM with STBP.

To evaluate the classification accuracy of the proposed architecture, this network was trained in neuromorphic datasets, which are challenging datasets because of noisy environment. The results showed that the proposed network achieved state-of-the-art classification accuracy on these datasets. Ablation experiments were implemented for examining the cascade structure. For further exploring the influence of hyper-parameters of the proposed network, the relationship between hyper-parameters and classification accuracy was also analyzed.

II. MATERIALS AND METHODS

The proposed cascade-SNN was shown in Fig. 1. The event-based image was storing in form of spikes. This network accepted input spikes and extracted higher level features via SCNN layer and MaxPool layer. Then the spikes were unfolded from two-dimensional into one-dimensional in spatial domain and sent to LSM layer to read out classification results.

This work was supported by National Natural Science Foundation of China under Grant 61774110.

Yuntao Han, Silu Cheng, and Jiangtao Xu are with the School of Microelectronics, Tianjin University, Tianjin 300072, China (e-mail: hyt_@tju.edu.cn; silu.cheng@tju.edu.cn; xujiantao@tju.edu.cn).

Tao Yu is with the School of Electrical and Information Engineering, Tianjin University, Tianjin 300072, China (e-mail: yutao_1999@tju.edu.cn).

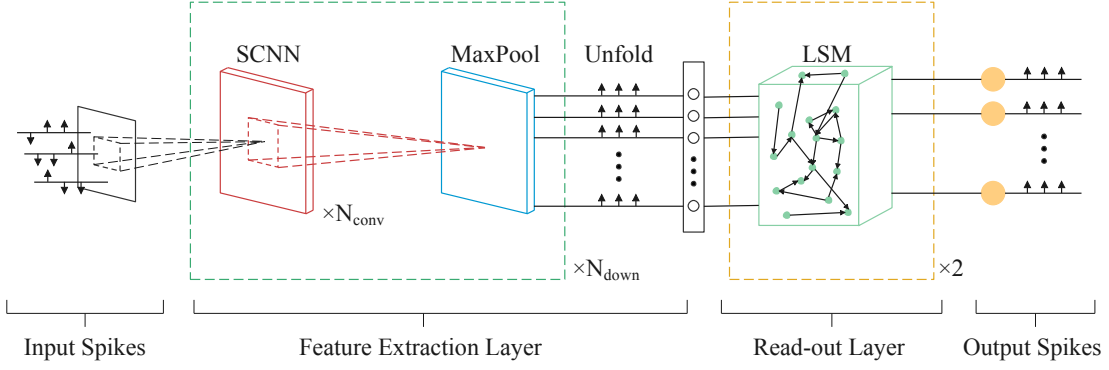


Fig. 1: Network Overview

A. LIF Neuron Model

LIF neuron model [8] is commonly used in SNN considering the biological credibility and implementation cost [9]. This model is the combination of leaky current, integrating pre-synapse spikes and firing spikes. Specifically, it is governed by:

$$\begin{cases} \tau \frac{du}{dt} = -u + I & u < V_{th}, \\ \text{fire a spike \& } u = u_{reset} & u \geq V_{th}. \end{cases} \quad (1)$$

where u is the membrane potential voltage, V_{th} is the threshold voltage, u_{reset} is the reset voltage, τ_{decay} is the time constant, I is the input current.

The potential voltage is determined by pre-synapse spikes received and decay with time. When the potential voltage reaches the threshold voltage, the neuron will fire a spike and its potential voltage will be reset to the u_{reset} .

Considering the implication of the conventional LIF model for the implementations on mainstream machine learning frameworks, [22] proposed an iterative LIF model, which converted differential equations into iterative equations.

$$\begin{cases} u_{n+1}^t = \tau_{decay} u_{n+1}^{t-1} (1 - o_{n+1}^{t-1}) + x_n^t \\ o_{n+1}^t = \varepsilon(u_{n+1}^t - V_{th}) \end{cases} \quad (2)$$

$$(3)$$

where u was the potential voltage, x was the input spikes, o was the output spikes, τ was the decay factor, ε was step function and n stood for the n th neuron and t stood for the t th timestep.

However, $\varepsilon(x)$ is not differentiable. [21] provides four curves to approximate the derivative of this function denoted by $h(u)$. The used function in this letter was shown below.

$$\begin{cases} o = \frac{1}{1 + e^{\frac{V_{th}-u}{\alpha}}} \\ h(u) = \frac{do}{du} = \frac{1}{\alpha} \frac{e^{\frac{V_{th}-u}{\alpha}}}{(1 + e^{\frac{V_{th}-u}{\alpha}})^2} \end{cases} \quad (4)$$

$$(5)$$

This model enables the forward calculation and the gradient back-propagation to be implemented on both spatial and temporal dimensions, which makes it friendly to general machine-learning programming frameworks.

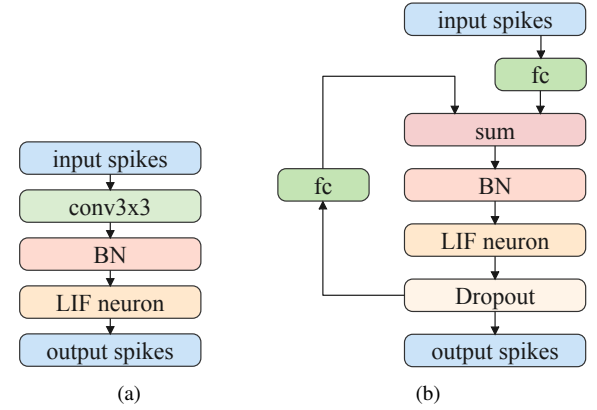


Fig. 2: SCNN and LSM Layer Structure

B. Feature Extraction Layer

SCNN was used for features extraction in SNNs. The structure used in letter was inspired by VGG19 [11], whose structure was shown in Fig. 1. N_{conv} stood for the number of convolutional layer and N_{down} stood for the number of downsampling layer. The feature extraction layer contained SCNN layer and MaxPool layer, and the scale of this layer could be modified by the number of extracted channels.

The basic structure of this layer was shown in Fig. 2(a), which consisted of 3x3 convolution kernel followed by a batch normalization (BN) layer and LIF neurons. On this basis, we proposed the basic SCNN block of the feature extraction layer.

C. LSM Read-out Layer

The LSM layer was used to read out the results of the network. It can contain previous information in this layer with feedback connections, instead of only forward propagation in traditional SNN. Fig. 2(b) showed the basic structure of this layer, which consisted of two fully connected (fc) layers and LIF neurons.

Equations for describing the LSM is shown below:

$$\begin{cases} \frac{du}{dt} = -\tau u + w_{fw} I + w_{fb} o \\ o = \varepsilon(u - V_{th}) \end{cases} \quad (6)$$

$$(7)$$

where u was the potential voltage of neurons in LSM layer, I was the discrete spike signal from the input layer, and o was the output of this layer. w_{fw} and w_{fb} were weight matrixes, which should be optimized by the learning algorithm.

However, this differential equation is not software friendly. The iterative formation was proposed to solve this problem, and it is convenient for the subsequent processing.

$$\begin{cases} \mathbf{u}_t = \tau \mathbf{u}_{t-1} \odot (1 - \mathbf{o}_{t-1}) + \mathbf{X} & (8) \\ \mathbf{X} = w_{fw} \mathbf{I}_t + w_{fb} \mathbf{o}_{t-1} & (9) \\ \mathbf{o}_t = \varepsilon(\mathbf{u}_t - V_{th}) & (10) \end{cases}$$

where \odot stands for the element-wise product.

For regulating our network, BN and dropout were applied. The complete algorithm of LSM was shown below.

Algorithm 1 LSM Layer Processing

Input: input spikes from last layer

$x \in [Timestep, Batch_size, Neuron]$.

Output: output spikes to next layer

$y \in [Timestep, Batch_size, Neuron]$.

Parameters:

forward weights $w_{fw} \in [Neuron_x, Neuron_y]$,

feedback weights $w_{fb} \in [Neuron_y, Neuron_x]$.

Initialization:

voltage of neurons $U \leftarrow 0$,

spike of neurons $O \leftarrow 0$.

1: **for** t **in** timesteps **do**

2: $X = w_{fw}x[t, :, :] + w_{fb}S$

3: $U = U + BatchNorm(X)$

4: $O = \text{if}(V \geq U_{threshold})$

5: **if** training **then**

6: $O = Dropout(O)$

7: **end if**

8: $O_{return}[t] = O$

9: **end for**

10: **return** O_{return}

The Read-out layer contained two LSM layer as shown in Fig. 1. The first layer had 1,000 neurons and the second layer had the same number of neurons as the number of classes.

D. STBP for LSM

STBP was a supervised learning algorithm, which was proposed in [21]. In this letter, we used STBP for training our network. The STBP training methodology for basic SNNs has already been proved in the original paper. The mathematical deduction of STBP for LSM was presented in this section.

The loss function of LSM layer could be written

$$L = \frac{1}{2} \|\mathbf{y} - \frac{1}{T} \sum_{t=1}^T \mathbf{o}^t\|^2 \quad (11)$$

where L was loss function, \mathbf{y} was the label vector in one-hot coding, T was the total simulation timesteps, \mathbf{o} was the output spikes of last LSM layer.

We acquired the required derivative $\frac{\partial L}{\partial \mathbf{o}^t}$ and $\frac{\partial L}{\partial \mathbf{u}^t}$, which were discussed in two cases:

1) $t = T$:

$$\begin{cases} \frac{\partial L}{\partial \mathbf{o}_i^T} = -\frac{1}{T} (y_i - \frac{1}{T} \sum_{t=1}^T o_i^t) & (12) \end{cases}$$

$$\begin{cases} \frac{\partial L}{\partial \mathbf{u}_i^T} = \frac{\partial L}{\partial \mathbf{o}_i^T} \frac{\partial \mathbf{o}_i^T}{\partial \mathbf{u}_i^T} = \frac{\partial L}{\partial \mathbf{o}_i^T} h(\mathbf{u}_i^T) & (13) \end{cases}$$

where $h(x)$ was the approximate derivative function of $\varepsilon(x)$.

2) $t < T$:

$$\begin{cases} \frac{\partial L}{\partial \mathbf{o}_i^t} = \sum_{k=1}^N (\frac{\partial L}{\partial \mathbf{o}_k^{t+1}} \frac{\partial \mathbf{o}_k^{t+1}}{\partial \mathbf{o}_i^t}) + \frac{\partial L}{\partial \mathbf{o}_i^T} & (14) \end{cases}$$

$$\begin{cases} \frac{\partial L}{\partial \mathbf{u}_i^t} = \frac{\partial L}{\partial \mathbf{u}_i^{t+1}} \frac{\partial \mathbf{u}_i^{t+1}}{\partial \mathbf{u}_i^t} = \frac{\partial L}{\partial \mathbf{o}_i^{t+1}} h(\mathbf{u}_i^{t+1}) \tau (1 - o_i^t) & (15) \end{cases}$$

where

$$\frac{\partial \mathbf{o}_k^{t+1}}{\partial \mathbf{o}_i^t} = \begin{cases} \frac{\partial \mathbf{o}_i^{t+1}}{\partial \mathbf{u}_i^{t+1}} (-\tau \mathbf{u}_i^T + w_{fb,ii}), & \text{if } k = i, \\ \frac{\partial \mathbf{o}_k^{t+1}}{\partial \mathbf{u}_k^{t+1}} w_{fb,ki}, & \text{if } k \neq i. \end{cases} \quad (16)$$

where $w_{fb,ii}$ and $w_{fb,ki}$ stood for the element of matrix w_{fb} in position (i, i) and in position (k, i) .

Finally, we obtained the derivatives with respect to w_{fw} and w_{fb} as follows:

$$\begin{cases} \frac{\partial L}{\partial w_{fw}} = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{u}^t} \frac{\partial \mathbf{u}^t}{\partial w_{fw}} = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{u}^t} (\mathbf{I}^t)^T & (17) \end{cases}$$

$$\begin{cases} \frac{\partial L}{\partial w_{fb}} = \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{u}^t} (\mathbf{o}^t)^T & (18) \end{cases}$$

III. EXPERIMENT

The experiment was running in SpikingJelly [23], which is an open-source deep learning framework for SNN based on pytorch. The learning algorithm was STBP and the optimizer was SGD with initial learning rate $r = 0.1$ and momentum 0.9. The learning rate decreased to half every 30 epochs. All the hyper-parameters used in this experiment were listed in TABLE I.

A. Experiments on DVS datasets

To evaluate the proposed SNN architecture and the influence of hyper-parameters, we used the DVS-CIFAR10 dataset [24]. It contains 10,000 event-stream recordings. In this experiment, the dataset was split into training set with 9000 images and testing set with 1000 images. The recordings were downsampled from the original 128×128 image size to 64×64 and the temporal resolution was reduced by accumulating spikes within frame data which was integrated from events [25]. With VGG-19 structure, our methods achieved the best performance with 76.70% accuracy in 20 timesteps. TABLE II compared our results with other models.

DVS-Gesture is a collection of moving gestures performed by 29 different individuals. Each instance in this dataset is a stream of events with size of 128×128 . The recordings were downsampled using the same method in DVS-CIFAR10. We compared our results with other related works on DVS-Gesture, as shown in TABLE III.

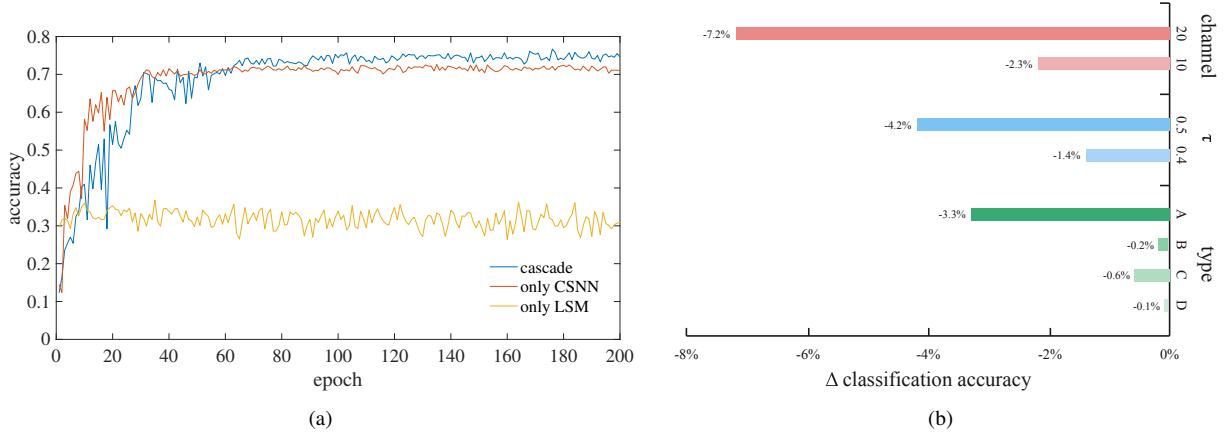


Fig. 3: Hyper-parameter Analysis (a)Ablation Experiment Result (b)Hyper-parameter Analysis

TABLE I: Hyper-parameter

hyper-parameter	value
V_{th}	1
v_{reset}	0
τ	0.7
channel	40
number of convolutional layers	16
dropout rate	0.5
batch size	32

TABLE II: Accuracy Comparisons on DVS-CIFAR1010

Model	Methods	Accuracy
[6]	HATS	52.40%
[18]	streaming rollout ANN	66.75%
[5]	DART	65.78%
[22]	STBP	60.50%
[20]	STBP-tdBN	67.80%
our model	cascade-SNN	76.70%

TABLE III: Accuracy Comparisons on DVS-Gesture

Model	Methods	Accuracy
[4]	STBP	93.40%
[19]	SLAYER	93.64%
[18]	streaming rollout ANN	95.68%
[26]	BPTT	94.59%
[20]	STBP-tdBN	96.87%
our model	cascade-SNN	97.57%

B. Ablation Experiment

To investigate the behavior of cascade-SNN as a proposal method, we conducted several ablation studies. We simulated the network which contained only LSM layer and only CSNN layer on DVS-CIFAR10 dataset and the results were shown in Fig. 3(a). Compared with the single layer structure, the convergence speed of cascade structure is slow, and more epoch needs to be trained to obtain stable classification accuracy, but accuracy is higher than that of single layer structure.

C. Hyper-parameters analysis

Like traditional ANNs, the performance of SNNs is influenced by the hyper-parameters of the network like scale, decay factor and depth. For finding the relationship between each

hyper-parameters and classification accuracy, we simulated the proposed network on DVS-CIFAR10 dataset. The accuracy descending with the different hyper-parameter was shown in Fig. 3(b). Except for the analyzed hyper-parameter, all hyper-parameters remained the same as shown in Table I.

1) *scale of feature extraction layer*: The number of channel represents the number of features extracted by the network. With the increase of the number of channels, the classification accuracy can be improved to a certain extent.

2) *decay factor of neuron model*: The parameter τ affects the influence of the previous time information on the current network output. As τ approaches 1, the attenuation of membrane potential in neuron model slows down, which is more conducive to the accumulation of membrane potential and the release of pulses. It could be seen that the classification accuracy increases with the increase of τ . This shows that SNN has the advantage of memory ability for time-domain signals in spatio-temporal signal processing compared with traditional neural networks.

3) *depth of feature extraction layer*: To measure the improvement brought by the increased depth of SCNN, the structure of the proposed neuron network was switched to different configuration of VGG. We chose configuration A, B, C, D with number of convolutional layers 8, 10, 13, 13. Compared to other configurations, configuration C has a unique convolutional layer whose kernel size is 1. Except configuration C, the classification accuracy increased with the increase of the number of convolutional layers. The network can extract higher-level features from input spikes with the structure going deeper.

IV. CONCLUSION

In this letter we proposed the cascade-SNN architecture which is a combination of SCNN for features extraction and LSM for read-out classification results. The experiment showed that it could achieve SOTA accuracy on complex noisy datasets. Ablation experiments were implemented to examine the cascade structure. We also examined different hyper-parameters of the proposed network. The influences of hyper-parameters in classification accuracy was analyzed.

REFERENCES

- [1] B. Zhang, L. Shi, and S. Song, "Creating more intelligent robots through brain-inspired computing," in *Science Robotics*, p. 3, 2016.
- [2] C. Brandli, R. Berner, M. Yang, S. Liu, and T. Delbruck, "A 240×180 130 db 3 μ s latency global shutter spatiotemporal vision sensor," *IEEE Journal of Solid-State Circuits*, vol. 49, no. 10, pp. 2333–2341, 2014.
- [3] G. Gallego, T. Delbruck, G. M. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. Davison, J. Conradt, K. Daniilidis, and D. Scaramuzza, "Event-based vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.
- [4] W. He, Y. Wu, L. Deng, G. Li, H. Wang, Y. Tian, W. Ding, W. Wang, and Y. Xie, "Comparing snns and rnns on neuromorphic vision datasets: similarities and differences," *Neural Networks*, vol. 132, pp. 108–120, 2020.
- [5] B. Ramesh, H. Yang, G. Orchard, N. A. Le Thi, S. Zhang, and C. Xiang, "Dart: distribution aware retinal transform for event-based cameras," *IEEE transactions on pattern analysis and machine intelligence*, vol. 42, no. 11, pp. 2767–2780, 2019.
- [6] A. Sironi, M. Brambilla, N. Bourdis, X. Lagorce, and R. Benosman, "Hats: Histograms of averaged time surfaces for robust event-based object classification," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1731–1740.
- [7] K. Roy, A. Jaiswal, and P. Panda, "Towards spike-based machine intelligence with neuromorphic computing," *Nature*, vol. 575, no. 7784, pp. 607–617, 2019.
- [8] P. Dayan and L. F. Abbott, *Theoretical neuroscience: computational and mathematical modeling of neural systems*. Computational Neuroscience Series, 2001.
- [9] E. M. Izhikevich, "Which model to use for cortical spiking neurons?" *IEEE Transactions on Neural Networks*, vol. 15, no. 5, pp. 1063–1070, 2004.
- [10] V. Kornijcuk, H. Lim, J. Y. Seok, G. Kim, S. K. Kim, I. Kim, B. J. Choi, and D. S. Jeong, "Leaky integrate-and-fire neuron circuit based on floating-gate integrator," *Frontiers in neuroscience*, vol. 10, p. 212, 2016.
- [11] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [12] J. Redmon and A. Farhadi, "Yolov3: An incremental improvement," *arXiv preprint arXiv:1804.02767*, 2018.
- [13] B. Rueckauer, I.-A. Lungu, Y. Hu, M. Pfeiffer, and S.-C. Liu, "Conversion of continuous-valued deep networks to efficient event-driven networks for image classification," *Frontiers in neuroscience*, vol. 11, p. 682, 2017.
- [14] A. Sengupta, Y. Ye, R. Wang, C. Liu, and K. Roy, "Going deeper in spiking neural networks: Vgg and residual architectures," *Frontiers in neuroscience*, vol. 13, p. 95, 2019.
- [15] S. Kim, S. Park, B. Na, and S. Yoon, "Spiking-yolo: spiking neural network for energy-efficient object detection," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 34, no. 07, 2020, pp. 11 270–11 277.
- [16] W. Maass, T. Natschl ger, and H. Markram, "Real-time computing without stable states: A new framework for neural computation based on perturbations," *Neural computation*, vol. 14, no. 11, pp. 2531–2560, 2002.
- [17] S. Guo, L. Qu, L. Wang, S. Tian, S. Li, and W. Xu, "Exploration of input patterns for enhancing the performance of liquid state machines," 2020.
- [18] A. Kugele, T. Pfeil, M. Pfeiffer, and E. Chicca, "Efficient processing of spatio-temporal data streams with spiking neural networks," *Frontiers in Neuroscience*, vol. 14, p. 439, 2020.
- [19] S. B. Shrestha and G. Orchard, "Slayer: Spike layer error reassignment in time," *arXiv preprint arXiv:1810.08646*, 2018.
- [20] H. Zheng, Y. Wu, L. Deng, Y. Hu, and G. Li, "Going deeper with directly-trained larger spiking neural networks," *arXiv preprint arXiv:2011.05280*, 2020.
- [21] Y. Wu, L. Deng, G. Li, J. Zhu, and L. Shi, "Spatio-temporal backpropagation for training high-performance spiking neural networks," *Frontiers in neuroscience*, vol. 12, p. 331, 2018.
- [22] Y. Wu, L. Deng, G. Li, J. Zhu, Y. Xie, and L. Shi, "Direct training for spiking neural networks: Faster, larger, better," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 33, no. 01, 2019, pp. 1311–1318.
- [23] W. Fang, Y. Chen, J. Ding, and other contributors, "Spikingjelly," <https://github.com/fangwei123456/spikingjelly>, 2020.
- [24] H. Li, H. Liu, X. Ji, G. Li, and L. Shi, "Cifar10-dvs: an event-stream dataset for object classification," *Frontiers in neuroscience*, vol. 11, p. 309, 2017.
- [25] W. Fang, Z. Yu, Y. Chen, T. Masquelier, T. Huang, and Y. Tian, "Incorporating learnable membrane time constant to enhance learning of spiking neural networks," *arXiv preprint arXiv:2007.05785*, 2020.
- [26] 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.