

Attention Spiking Neural Networks

Man Yao, Guangshe Zhao, Hengyu Zhang, Yifan Hu, Lei Deng, Yonghong Tian, Bo Xu, and Guoqi Li

Abstract—Benefiting from the event-driven nature and sparse spiking communication of the brain, spiking neural networks (SNNs) are becoming a promising energy-efficient alternative to traditional artificial neural networks (ANNs). However, the performance gap between SNNs and ANNs has been a significant hindrance to deploying SNNs ubiquitously for a very long period of time. To leverage the full potential of SNNs, we study the effect of attention mechanisms in SNNs, which makes them focus on important information. We first present our idea of attention in SNNs with a plug-and-play combined module kit, termed the Multi-dimensional Attention (MA) module. Then, a new attention SNN architecture with end-to-end training called "MA-SNN" is proposed, which infers attention weights along the temporal dimension, channel dimension, as well as spatial dimension separately or simultaneously. Based on the existing neuroscience theories, we exploit the attention weights to optimize membrane potentials, which in turn regulate the spiking response in a data-dependent way. At the cost of negligible additional parameters, MA facilitates vanilla SNNs to achieve sparser spiking activity, better performance, and energy efficiency concurrently. Experiments are conducted in event-based DVS128 Gesture/Gait action recognition and ImageNet-1k image classification. On Gesture/Gait, the spike counts are reduced by 84.9%/81.6%, the task accuracy and energy efficiency are improved by 5.9%/4.7% and $3.4\times/3.2\times$. On ImageNet-1K, we achieve top-1 accuracy of 75.92% and 77.08% on single/4-step Res-SNN-104, which are state-of-the-art results in SNNs. Compared with counterpart Res-ANN-104, the performance gap becomes $-0.95/+0.21$ percent and has $31.8\times/7.4\times$ better energy efficiency. To our best knowledge, this is for the first time, that the SNN community achieves comparable or even better performance compared with its ANN counterpart in the large-scale dataset. To analyze and support the effectiveness of MA-SNN, we theoretically prove that the spiking degradation or the gradient vanishing, which usually holds in general SNNs, can be resolved by introducing the block dynamical isometry theory. We also analyze the efficiency of MA-SNN based on our proposed spiking response visualization method. Our work lights up SNN's potential as a general backbone to support various applications in the field of SNN research, with a great balance between effectiveness and efficiency.

Index Terms—Spiking neural network, Attention mechanism, Neuromorphic computing, Efficient neuromorphic inference

1 INTRODUCTION

As the most remarkable neural network, the human brain is incredibly efficient and capable of performing complex pattern recognition tasks, and has always been a source of innovation for artificial neural networks (ANNs) or conventional deep learning models [1], [2]. In the recent past, by reasonably emulating the deep hierarchy structure of the visual cortex, deep ANNs obtained powerful representation and brought amazing successes in a myriad of artificial intelligence applications, e.g., compute vision (CV) [3], natural language processing (NLP) [4], medical diagnosis [5], game playing [6], etc. Unfortunately, ANNs pay enormous computational costs to achieve such feats. For example, a standard computer performing only recognition among 1,000 different kinds of objects (ImageNet-1K dataset) expends about 250 watts [7]. By contrast, the human brain can operate with only nearly 20 watts consumption for various impressive achievements (such as simultaneous

recognition, reasoning, control, and movement) [8]. Many real-world platforms, e.g., smartphones, Internet-of-Things devices among others, have resources and battery constraints, which restrict the implementation of deep ANN [9]. To enable intelligence on such platforms, how to exploit the inherent efficient computation paradigm of the biological neural systems to achieve low-power of implementation of neural networks, is of great value.

Spiking neural networks (SNNs) offer an alternative for enabling energy-efficient intelligence, which emulate biological neuronal functionality by adopting binary spiking signals (0-nothing or 1-spiking event) to complete inter-neuron communication [10]. As a kind of neuromorphic computing algorithm, SNNs can be smoothly executed on the sparse neuromorphic chip, by only handling spike-based accumulate (AC) operations, and can avoid computing the zero values of input or activation (i.e., *event-driven*) [7]. Thus, SNNs consume much lower power than ANNs that are dominated by energy-hungry multiply-and-accumulate (MAC) operations on conventional dense computing hardware such as GPUs. With the release of neuromorphic chips like Tianjic [11], TrueNorth [12], and Loihi [13], we are not very far from neuromorphic processors becoming a part of everyday life.

It remains a challenge to directly train large-scale SNNs to achieve comparable performance with counterpart ANNs for real-world pattern recognition tasks. A recent study of directly training builds advancing residual learning to construct large-scale SNNs, and alleviates the performance gap between deep SNNs and ANNs [14]. However, the performance gap still exists. On the other hand, computation

- M. Yao is with the School of Automation Science and Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China, and also with Peng Cheng Laboratory, China.
 - G. Zhao is with the School of Automation Science and Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi, China.
 - H. Zhang is with Tsinghua Shenzhen International Graduate School, Tsinghua University, Shenzhen, China.
 - Y. Hu, and L. Deng are with Center for Brain-Inspired Computing Research, Department of Precision Instrument, Tsinghua University, Beijing, China.
 - Y. Tian is with Institute for Artificial Intelligence, Peking University, Beijing, China, and also with Peng Cheng Laboratory, China
 - B. Xu and G. Li are with Institute of Automation, Chinese Academy of Sciences, Beijing, China.
- The corresponding author: Guoqi Li (E-mail: guoqi.li@ia.ac.cn).

over multi-time steps¹ in deep SNNs [14], [15], [16] not only boosts the training time and simulation hardware costs, but also incurs high inference latency, more overall energy budget, and memory access overhead of fetching membrane potentials. These limitations prohibit the potential effective algorithm design and lessen the energy benefits of SNNs. To break the ice, we urgently need to take novel inspiration from how the brain works and classic deep learning to build more effective and efficient SNNs.

Humans can naturally and effectively find salient regions in complex scenes [17]. Motivated by this observation, attention mechanisms have been introduced into deep learning and achieved remarkable success in a wide spectrum of application domains. Current attention in deep learning generally exists in two ways. One is posing a fundamental paradigm shift in the way of executing meta-operator such as using self-attention conduct Transformer [18]. The other prefers integrating with the existing classic deep ANNs that work as auxiliary recalibration modules to increase the representation power of the basic model, such as attention convolutional neural network (attention CNNs) [19]. Recently, apart from the classic application in the NLP, the Transformer structure made its grand debut in the CV, and quickly set off an overwhelming wave of pure attention architecture design by its glaring performance in various tasks [20]. The success of self-attention facilitates researchers' understanding of different deep learning architectures, including Transformer, CNN, multi-layer perceptron (MLP), etc., and sparks more effective universal network architecture design [21]. Moreover, the above two attention practices can be combined together, such as using attention as an independent auxiliary module for the Transformer to focus on informative features or patches [22].

In contrast to the rapid development of attention mechanisms in ANNs, the application of attention in the SNN domain remains to be exploited. Existing few works are totally different from the aforementioned attention practices in traditional deep learning. They focus on using SNN to simulate the attention mechanism [23], [24] or executing SNN model compression by attention [25]. We do not intend to shift the meta-operator of existing SNNs, e.g., replacing convolution (or fully connected) with self-attention, but try to apply the attention as an auxiliary unit in a simple and lightweight way to easily integrate with existing SNN architectures for improving representation power, like attention CNNs. Challenges in adapting attention to SNNs arise from three aspects. Firstly, we must keep the neuromorphic computing characteristic of SNNs, which is the basis of SNN's energy efficiency. Thus, implementing the attention while retaining SNN's event-driven is the primary consideration. Secondly, SNNs are used to process various applications, such as sequential event streams and static images. We need to diverse attention SNN design to cope with different scenarios. Thirdly, binary spiking activity makes deep SNNs suffer from spike degradation [15] and gradient vanishing [14], collectively referred to as the *degradation problem*, i.e., an accuracy drop would occur on both the training and test sets when the network deepens. Attention should not make

the case worse.

In visual neuroscience, attention enhances neuronal communication efficacy by modulating synaptic weights [26] and neuronal spiking activity rate [27] in the noisy sensory environment. To emulate attention in the brain, we employ attention to facilitate optimizing the membrane potential of spiking neurons, which can be equivalent to synaptic alteration and would not disrupt the event-driven nature of SNNs. Our design philosophy is clear, exploiting attention to regulate membrane potentials, i.e., focusing on important features and suppressing unnecessary ones, which in turn affects the spiking activity. In contrast, attention is applied to refine activations in CNNs [19]. The underlying reason is that neurons in CNNs communicate with each other using activations coded in continuous values rather than brain-like spiking activations. To adapt attention SNNs to a variety of application scenarios, we merge multi-dimensional attention with SNN (MA-SNN), including *temporal*, *channel*, and *spatial* dimensions, to learn 'when', 'what' and 'where' to attend, respectively. These attention dimensions are exploited separately or simultaneously according to specific task metric requirements such as latency, accuracy, and energy cost. Classic convolutional block attention module (CBAM) [28] is adopted as the basic module to construct MA-SNN. Furthermore, attention residual SNNs are designed to process the large-scale ImageNet-1K. We exploit the MS-Res-SNN [14] as the backbone because of its higher accuracy and shortcut connection manner. We argue that membrane-shortcut in MS-Res-SNN is identical to our motivation for introducing the attention, which can also be seen as a way to optimize the membrane potentials.

The advantages of MA-SNN exist in three folds. Firstly, by emulating attention in brain, we propose the MA-SNN. Extensive experimental results on various tasks show that optimization of membrane potential in a data-dependent manner by attention can lead to sparser spiking responses and incurs better performance and energy efficiency concurrently, like the human brain. Secondly, we uncover the attention mechanism in MA-SNN. We answer one key problem: how can both effectiveness and energy efficiency be achieved simultaneously in MA-SNN. To address this issue, a new spiking response visualization method is proposed to observe the effect of attention-optimized membrane potential on spiking response. We show that the effectiveness of MA-SNN mainly stems from the proper focusing, just the same as previous CNN works [19], [29], [30], [31]. At the efficiency aspect, MA adaptively inhibits the membrane potentials of the background noise, then these spiking neurons would not be activated. With this point of view, we could explain why a much lower spiking activity rate can be achieved in attention SNN with great energy efficiency. Thirdly, we prove that the degradation problem, which holds in general deep SNNs, can be resolved when adding attention to MS-Res-SNN (i.e., Att-Res-SNN). Specifically, we prove the gradient norm equality [32] can be achieved in our attention residual learning by introducing the block dynamical isometry theory, which means that one could train very deep Att-Res-SNN in the same way as in MS-Res-SNN. To summarize, the main contributions of this work are as follows:

1. A time step is the unit of time taken by each input frame to be processed through all layers of the model.

- **Multi-dimensional Attention SNN:** Inspired by the attention mechanisms in neuroscience, we present our idea of attention SNN and propose the MA-SNN, which merges multi-dimensional attention with SNN and inherits the event-driven nature, including temporal, channel, and spatial dimensions to learn ‘when’, ‘what’, and ‘where’ to attend. The sparse spiking activity, performance, and energy efficiency of MA-SNN are verified on the multiple benchmarks under the multi-scale constraints of output latency. Based on the proposed model, now the SNN community is able to achieve comparable or even better performance compared with its ANN counterpart in the large-scale dataset.
- **Understanding and Visualizing of Attention:** Through the proposed spiking response visualization method, it is shown that the effectiveness of MA-SNN mainly stems from proper focusing, and efficiency comes from the improvement of sparsity by inhibiting the membrane potentials of the background noise. Thus, we explain why both the effectiveness and efficiency of attention can be achieved concurrently in SNNs.
- **Gradient Norm Equality of Att-Res-SNN:** We prove the gradient norm equality [32] can be achieved in Att-Res-SNN based on the block dynamical isometry theory. The degradation problem that holds in general deep SNNs can then be resolved when adding attention to MS-Res-SNN. Thus, we are able to train very deep Att-Res-SNN to enhance the potential of SNNs.

The rest of the paper is organized as follows. Section 2 reports preliminaries. Section 3 introduces our MA-SNN. Section 4 gives how to evaluate the energy cost of attention SNNs. Section 5 verifies the effectiveness and efficiency of our methods. Section 6 conducts ablation studies to comprehend the design of MA-SNN. Section 7 understands and visualizes the effectiveness and efficiency of attention SNNs. Section 8 concludes this work.

2 PRELIMINARIES

Training Methods of SNNs. ANN-to-SNN conversion and directly training an SNN are two main routines to train deep SNNs. The basic idea of ANN-to-SNN is that the activation values in a ReLU-based ANN can be approximated by the average firing rates of an SNN under the rate-coding scheme. There is a trade-off issue of accuracy and latency in ANN-to-SNN methods that need sufficient time steps for rate-coding to alleviate approximation errors [33]. Although the converted SNN can obtain the smallest accuracy gap with ANN in some large-scale structures, such as VGG and ResNet [34], [35], they need a longer time step or complicated training methods, which increases the SNN’s latency and restricts the practical application. Directly training an SNN is another training mode of SNN, which constitutes a continuous relaxation of the non-smooth spiking to enable backpropagation with a surrogate gradient [36]. Compared with ANN-to-SNN, it has a great advantage in the number of time steps and can also be applied to temporal tasks, e.g., event-based datasets. Direct training algorithms are diverse

in the selection of coding schemes such as time-coding [37] and rate-coding [38]. Time-coding has limited the network scale, and we use the rate-coding direct training method to obtain large-scale SNNs in this paper.

Event-based Vision. Dynamic vision sensor (DVS), which encodes the time, location, and polarity of the brightness changes for each pixel into event streams with a μs level temporal resolution, poses a new paradigm shift in visual information acquisition. Compared with conventional cameras, the advantages of DVS include [39]: requiring fewer resources since the events are only triggered when the intensity changes; a high temporal resolution which can avoid motion blur; a very high dynamic range which makes the DVS able to acquire information from challenging illumination conditions. These characteristics promote the application of DVS in various scenarios, such as high-speed object tracking [40], autonomous driving [41], low-latency interaction [42], etc. Processing events one by one is limited to performance because a single event has little information. The general method is to group event streams with a certain temporal window as alternative representations, e.g., frame-based [43], graph-based [44], etc. In this paper, we adopt the frame-based representation that transforms event streams into high-rate videos, where each frame has many blank (zero) areas. SNN is suitable to process event frames since it can skip the computation of the zero areas in each input frame [7].

Attention in CNNs. Depth, width, and cardinality are three important factors to get rich representation power in CNN architecture design. Apart from these factors, attention is another different aspect of architecture design that increases representation power by focusing on important information. Its significance has been studied extensively in the previous literature. Hu *et al.* [19] pioneered the attention module in CNNs, which first proposed the concept of channel attention and presented SENet for this purpose. Motivated by different channels that usually represent different objects, SENet models the relationship between channels to adaptively refine the weight of each channel, i.e., determining what to pay attention to. Since convolution operations extract informative features by blending cross-channel and spatial information together, Woo *et al.* [28] proposed CBAM that sequentially applies both channel and spatial attention modules to determine what and where to pay attention to concurrently. There are many optimizations from various aspects for the modeling of the channel and spatial attention, such as effectiveness [45], space complexity [31], computational complexity [46], etc. Please refer to [47] for a comprehensive review of attention in CV.

3 MULTI-DIMENSIONAL ATTENTION SPIKING NEURAL NETWORKS

In this section, we introduce the network input preprocessing and the Conv-based SNN in Section 3.1 and Section 3.2, respectively. Then we design the MA module that learns temporal (when), channel (what), and spatial-wise (where) attention separately in Section 3.3. Finally, we give our design of attention residual learning for SNN in Section 3.4.

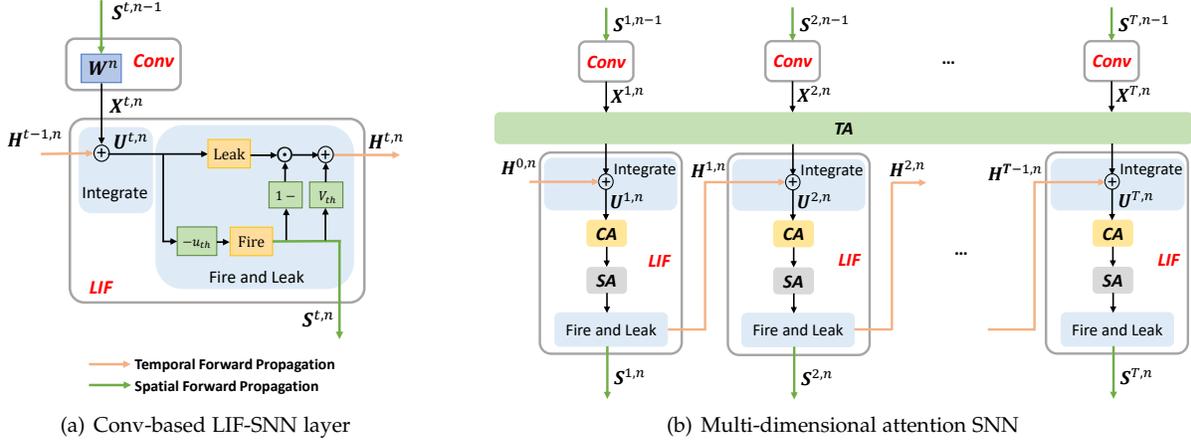


Fig. 1: The Conv-based SNN layer and the overview of MA-SNN.

3.1 Network Input

Event-based streams. We adopt the frame-based representation as the preprocessing method, which transforms event streams into high-rate frame sequences where each frame has many zero areas. Event stream comprises four dimensions: two spatial coordinates (x, y) , the timestamp, and the polarity of each single event. The polarity p indicates an increase (ON) or decrease (OFF) of brightness, where ON/OFF can be represented via +1/-1 values. Assume the initial temporal resolution of event stream is dt' (μ s level) and the spatial resolution is $h_0 \times w_0$, the spike pattern tensor $\mathbf{S}_{t'} \in \mathbf{R}^{2 \times h_0 \times w_0}$ is equal to events set $E_{t'} = \{e_i | e_i = [x_i, y_i, t', p_i]\}$ at timestamp t' . We can set a new millisecond-level temporal resolution $dt = dt' \times \alpha$, and consecutive α spike patterns can be grouped as a set

$$E_t = \{\mathbf{S}_{t'}\}, \quad (1)$$

where $t' \in [\alpha \times t, \alpha \times (t+1) - 1]$. Then, the frame for input layer at t time $\mathbf{S}^{t,0} \in \mathbf{R}^{2 \times h_0 \times w_0}$ based on dt can be got by

$$\mathbf{S}^{t,0} = q(E_t) = \sum_{t'=\alpha \times t}^{\alpha \times (t+1) - 1} \mathbf{S}_{t'}, \quad (2)$$

where $t \in \{1, 2, \dots, T\}$ is the *time step*, and $q(\cdot)$ is element-wise addition function. In this way, the event stream can be transformed into a sequence of real-valued frames with a new frame rate, e.g., $dt = 1ms$ corresponds to 10^3 frames per second. Fig. 7 shows an example of event frames.

Static images. For the analog-valued signal of pixel intensity in images, adding an encoding layer to generate spike signals globally is a generic method of SNN [48]. Since the SNN is a kind of spatio-temporal model, the image is copied and used as input frame at each time step when $T > 1$, i.e., $\mathbf{S}^{1,0} = \mathbf{S}^{2,0} = \dots = \mathbf{S}^{T,0}$.

3.2 Spiking Neural Networks

Spiking Neuron. Spiking neurons, the basic compute units of SNN, communicate through spikes coded in binary activations. The leaky integrate-and-fire (LIF) model is one of the most commonly used spiking neuron models, since it is a trade-off between the complex dynamic characteristics of

biological neurons and the simplified mathematical form. It is suitable for simulating large-scale SNN and can be described by a differential function [10]

$$\tau \frac{du(t)}{dt} = -u(t) + I(t), \quad (3)$$

where τ is a time constant, and $u(t)$ and $I(t)$ are the membrane potential of the postsynaptic neuron and the input collected from presynaptic neurons, respectively.

Conv-based LIF-SNN. Solving this differential equation, a simple iterative representation of LIF-SNN layer [36], [49] for easy inference and training is governed by

$$\begin{cases} \mathbf{U}^{t,n} = \mathbf{H}^{t-1,n} + \mathbf{X}^{t,n} \\ \mathbf{S}^{t,n} = \text{Hea}(\mathbf{U}^{t,n} - u_{th}) \\ \mathbf{H}^{t,n} = V_{reset} \mathbf{S}^{t,n} + (\beta \mathbf{U}^{t,n}) \odot (1 - \mathbf{S}^{t,n}), \end{cases} \quad (4)$$

where t and n denote the time step and layer, $\mathbf{U}^{t,n}$ means the membrane potential which is produced by coupling the spatial feature $\mathbf{X}^{t,n}$ and the temporal input $\mathbf{H}^{t-1,n}$, u_{th} is the threshold to determine whether the output spiking tensor $\mathbf{S}^{t,n}$ should be given or stay as zero, $\text{Hea}(\cdot)$ is a Heaviside step function that satisfies $\text{Hea}(x) = 1$ when $x \geq 0$, otherwise $\text{Hea}(x) = 0$, V_{reset} denotes the reset potential which is set after activating the output spiking, and $\beta = e^{-\frac{dt}{\tau}} < 1$ reflects the decay factor, and \odot denotes the element-wise multiplication.

In Eq. 4, spatial feature $\mathbf{X}^{t,n}$ can be extracted from the original input $\mathbf{S}^{t,n-1}$ through a convolution operation:

$$\mathbf{X}^{t,n} = \text{AvgPool} \left(\text{BN} \left(\text{Conv} \left(\mathbf{W}^n, \mathbf{S}^{t,n-1} \right) \right) \right), \quad (5)$$

where $\text{AvgPool}(\cdot)$, $\text{BN}(\cdot)$ and $\text{Conv}(\cdot)$ mean the average pooling, batch normalization [50] and convolution operation respectively, \mathbf{W}^n is the weight matrix, $\mathbf{S}^{t,n-1} (n \neq 1)$ is a spike tensor that only contains 0 and 1, and $\mathbf{X}^{t,n} \in \mathbb{R}^{c_n \times h_n \times w_n}$. To simplify the notation, bias terms are omitted. BN is a default operation following the Conv, we also omit it in the rest of this paper.

Fire and Leak Mechanism. As shown in Fig. 1(a), the CNN-based LIF-SNN layer consists of two parts, the Conv and LIF. The Conv module extracts spatial features from original input $\mathbf{S}^{t,n-1}$ firstly. Then the LIF module integrates

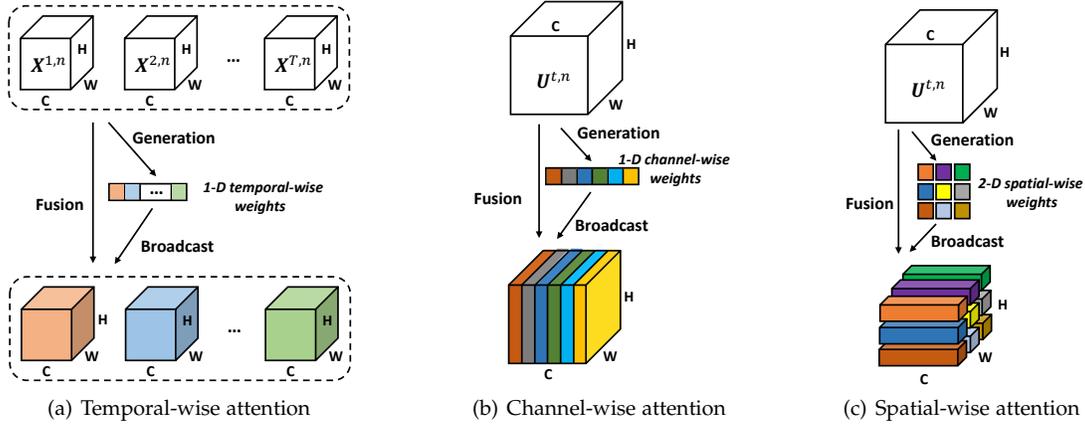


Fig. 2: Illustration of different attention dimensions.

the spatial feature $X^{t,n}$ and the temporal input $H^{t-1,n}$ into membrane potential $U^{t,n}$. Finally, the fire and leak mechanism is exploited to generate spatial spiking tensor for the next layer and the new cell states for the next time step. When the entries in $U^{t,n}$ are greater than the threshold u_{th} , the spatial output of spiking sequence $S^{t,n}$ will be activated, and the entries in $U^{t,n}$ will be reset to V_{reset} , then the temporal output $H^{t,n}$ should be decided by the $X^{t,n}$ since $1 - S^{t,n}$ must be 0. Otherwise, the decay of the $U^{t,n}$ will be used to transmit the $H^{t,n}$, since the $S^{t,n}$ is 0, which means there is no activated spiking output. Note that, after the convolution operation, all tensors in the LIF module have the same dimensions, i.e., $X^{t,n}, H^{t-1,n}, U^{t,n}, S^{t,n}, H^{t,n} \in \mathbb{R}^{c_n \times h_n \times w_n}$.

3.3 Multi-dimensional Attention for SNNs

Attention in Neuroscience. Selective attention is a powerful brain mechanism that enables enhanced processing of relevant information while preventing interference from distracting noise [51]. Many studies in humans and animals have investigated the effect of visual attention in single-neuron, and given two solid observations [26], [52]. The first conclusion is about how attention modulates neuronal communication with each other [26]. By altering dendritic spines (building blocks of synapse, i.e., synaptic weights [53]) in a dynamic and highly selective manner, neurons enhance the detection of salient information in the noisy sensory environment. Another consensus is that attention is associated with neuronal spiking activity rate, not only in local single-neuron [26] but also in more global visual cortex areas [54]. Inputs that carry salient sensory information enhance the spiking activity rates of neurons, while potentially redundant input information does the opposite effect.

Overview of Attention. Although the idea of attention in CNNs is the same as neuroscience, there are fundamental differences. Obviously, continuous activations of CNNs do not conform to the spiking activation properties in biological neurons, thus losing the potential energy efficiency earnings caused by attention. To mimic the attention that modulates the spiking activity of neurons in brain, we optimize the membrane potential of spiking neurons by attention in a data-dependent manner, and the spiking response of SNN

is consequently regulated. Generally, we can formulate attention processes as:

$$x_{Att} = f(g(x), x), \quad (6)$$

where x_{Att} is output with an attention mechanism, $g(x)$ is the function of generating attention weights which corresponds to the process of attending to the discriminative moments or regions. $f(g(x), x)$ means processing input x based on the attention weights $g(x)$.

As shown in Fig. 1(b), we design a multi-dimensional attention module that learns temporal (when), channel (what), and spatial (where) attention separately. Attention modules for each dimension can be exploited in an independently or jointly way. For MA of this paper, we adopt

$$x_{Att} = g(x) \cdot x, \quad (7)$$

and input x is usually an intermediate feature map.

Temporal-wise Attention (TA). Consider the event-based visual recognition in a real-time interaction scenario, where an event stream is divided into high-rate frames in sequence while a prediction can be retrieved after processing the data with $t_{lat} = dt \times T$ ms. Our previous research [43] observed that the accuracy of the SNN would not become worse even they masked half of the input event frames, and proposed a lightweight temporal-wise attention SNN to handle event streams by discarding irrelevant event frames. However, directly masking the input frames just keeps accuracy unchanged. By contrast, we use an advanced method that refines the interior features of the network in temporal dimension by exploiting the inter-time step relationship of feature blocks, which could obtain performance gain and energy reduction concurrently.

As depicted in Fig. 2(a), to compute the 1-D TA weights, we collect intermediate feature blocks of n -th layer at all time steps $\mathbf{X}^n = [\dots, \mathbf{X}^{t,n}, \dots] \in \mathbb{R}^{T \times c_n \times h_n \times w_n}$ as input. TA function $g_t(\cdot)$ infers a 1-D TA weight vector $g_t(\mathbf{X}^n) \in \mathbb{R}^{T \times 1 \times 1 \times 1}$ as

$$\mathbf{X}_{TA}^n = g_t(\mathbf{X}^n) \odot \mathbf{X}^n, \quad (8)$$

where $\mathbf{X}_{TA}^n = [\dots, \mathbf{X}_{TA}^{t,n}, \dots] \in \mathbb{R}^{T \times c_n \times h_n \times w_n}$ is the temporal-wise refined feature blocks. During multiplication, the TA weights are broadcasted (copied) accordingly, along both the channel and spatial dimension.

Inspired by the squeeze-and-excitation operation in previous classic channel-wise attention module designs [19], [28], we design TA function $g_t(\cdot)$ in a similar manner. We first aggregate spatial-channel information of a feature block at each time step by using both average-pooling and max-pooling operations, generating two different temporal context descriptors, which denote average-pooled features and max-pooled features respectively. Then, we transform both average-pooled and max-pooled features to a TA weight vector by a shared MLP network, i.e.,

$$g_t(\mathbf{X}^n) = \sigma(\mathbf{W}_{t1}^n(\text{ReLU}(\mathbf{W}_{t0}^n(\text{AvgPool}(\mathbf{X}^n)))) + \mathbf{W}_{t1}^n(\text{ReLU}(\mathbf{W}_{t0}^n(\text{MaxPool}(\mathbf{X}^n))))), \quad (9)$$

where $\text{AvgPool}(\mathbf{X}^n), \text{MaxPool}(\mathbf{X}^n) \in \mathbb{R}^{T \times 1 \times 1 \times 1}$ represent the results of average-pooling and max-pooling layer respectively, σ means the sigmoid function, $\mathbf{W}_{t0}^n \in \mathbb{R}^{\frac{T}{r_t} \times T}$ and $\mathbf{W}_{t1}^n \in \mathbb{R}^{T \times \frac{T}{r_t}}$ are the weights of linear layers in the shared MLP, r_t denotes the temporal dimension reduction factor used to control the computational burden of MLP.

Channel-wise Attention (CA). CA in CNNs only works on spatial features. Discriminatively, our CA design optimizes spatio-temporal fusion information of SNN, i.e., we adopt the CA function $g_c(\cdot)$ to directly refine the membrane potential of spiking neurons (more discusses of CA location design are given in Section 6.1). It is well known that each channel of feature maps corresponds to a certain visual pattern, and CA focuses on "what" are salient semantic attributes for the given input. Interestingly, we find that another key role of attention is the suppression of minor features, which is usually ignored in attention CNN but crucial for the efficiency of the SNN (details in Section 7).

We adopt the classic CBAM [28] to generate the CA vector for SNN at each time step as

$$g_c(\mathbf{U}^{t,n}) = \sigma(\mathbf{W}_{c1}^n(\text{ReLU}(\mathbf{W}_{c0}^n(\text{AvgPool}(\mathbf{U}^{t,n})))) + \mathbf{W}_{c1}^n(\text{ReLU}(\mathbf{W}_{c0}^n(\text{MaxPool}(\mathbf{U}^{t,n}))))), \quad (10)$$

where $g_c(\mathbf{U}^{t,n}) \in \mathbb{R}^{c_n \times 1 \times 1}$ is the 1-D CA weights, $\text{AvgPool}(\mathbf{U}^{t,n}), \text{MaxPool}(\mathbf{U}^{t,n}) \in \mathbb{R}^{c_n \times 1 \times 1}$, $\mathbf{W}_{c0}^n \in \mathbb{R}^{\frac{c_n}{r_c} \times c_n}$, $\mathbf{W}_{c1}^n \in \mathbb{R}^{c_n \times \frac{c_n}{r_c}}$, and r_c represents the channel dimension reduction factor. To reduce parameters, the MLP weights, \mathbf{W}_{c0}^n and \mathbf{W}_{c1}^n , are shared for different time steps. The refined feature $\mathbf{U}_{CA}^{t,n} \in \mathbb{R}^{c_n \times h_n \times w_n}$ is computed as

$$\mathbf{U}_{CA}^{t,n} = g_c(\mathbf{U}^{t,n}) \odot \mathbf{U}^{t,n}. \quad (11)$$

Note that, the only difference between TA and CA is the inputs of attention, the former is \mathbf{X}^n , the latter is $\mathbf{U}^{t,n}$ (see Fig. 2). Actually, $g_t(\cdot)$ and $g_c(\cdot)$ can be acted by various attention modules, including classic SE [19], energy-efficient ECA [46], and parameter-free SimAM [31], etc. We conduct ablation studies of these choices in Section 6.2.

Spatial-wise Attention (SA). Different from the above TA and CA, the SA focuses on "where" is an informative part. We adopt the SA part of CBAM [28] as our SA function $g_s(\cdot)$, which is described as:

$$g_s(\mathbf{U}^{t,n}) = \sigma(f^{7 \times 7}([\text{AvgPool}(\mathbf{U}^{t,n}); \text{MaxPool}(\mathbf{U}^{t,n})])), \quad (12)$$

where $\text{AvgPool}(\mathbf{U}^{t,n}), \text{MaxPool}(\mathbf{U}^{t,n}) \in \mathbb{R}^{1 \times h_n \times w_n}$, $g_s(\mathbf{U}^{t,n}) \in \mathbb{R}^{1 \times h_n \times w_n}$ is the 2-D SA attention weights, and $f^{7 \times 7}$ represents a convolution operation with the filter size

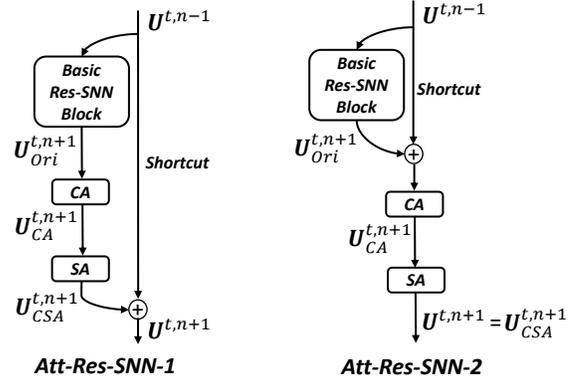


Fig. 3: Attention residual block contains three parts: basic Res-SNN block (MS-Res-SNN [14], details in Fig.5), shortcut, and CSA module. We exploit MA on the basic Res-SNN block outputs (i.e., membrane potential of spiking neurons) in each block. We recommend Att-Res-SNN-1 as the scheme for attention residual learning.

of 7×7 (default hyper-parameter setting in CBAM). Then, the refined feature $\mathbf{U}_{SA}^{t,n}$ (see Fig. 2(c)) is computed as

$$\mathbf{U}_{SA}^{t,n} = g_s(\mathbf{U}^{t,n}) \odot \mathbf{U}^{t,n}, \quad (13)$$

and the convolution operation $f^{7 \times 7}$ is shared for each time step to save parameters.

Finally, when we adopt the above three dimensions of attention concurrently (i.e., TCSA), compared with $\mathbf{U}^{t,n}$ of vanilla SNN in Eq. 4, the *new membrane potential* behaviors of TCSA-SNN layer follow

$$\begin{aligned} \mathbf{X}_{TA}^n &= g_t(\mathbf{X}^n) \odot \mathbf{X}^n, \\ \mathbf{U}_{CA}^{t,n} &= g_c(\mathbf{H}^{t-1,n} + \mathbf{X}_{TA}^{t,n}) \odot (\mathbf{H}^{t-1,n} + \mathbf{X}_{TA}^{t,n}), \\ \mathbf{U}^{t,n} &= g_s(\mathbf{U}_{CA}^{t,n}) \odot \mathbf{U}_{CA}^{t,n}. \end{aligned} \quad (14)$$

3.4 Attention Residual Learning of SNNs

SNNs have been theoretically proved its computational power can be matched with ANNs [10]. However, the limited scale in SNNs restricts the network representation power, which consequently impedes SNN's practice applications and intensifies the performance gap between SNNs and ANNs. Residual learning [2] becomes a milestone work in deep learning through attaching an identity skip connection throughout the entire network to achieve "very deep" neural networks. Unfortunately, directly copying the classic residual structure to SNNs, there will still be a degradation problem in that the deeper SNNs have higher train loss than the shallower SNNs. There are three mainstream residual structures of SNN, vanilla Res-SNN [15], SEW-Res-SNN [16], and MS-Res-SNN [14]. The main difference among these Res-SNN works is the construction of the basic residual blocks (see Fig.5), and currently there is no uniformly standard residual block scheme in the SNN community.

MA can be integrated into existing residual SNN architectures without constraints, and we consistently exploit attention to optimize membrane potential of spiking neurons. In this paper, we adopt the MS-Res-SNN [14] as the basic residual block. As shown in Fig. 3, the channel-spatial

attention module can be integrated with existing residual block in two ways. For Att-Res-SNN-1, we have

$$\begin{aligned} \mathbf{U}_{CA}^{t,n+1} &= g_c(\mathbf{U}_{Ori}^{t,n+1}) \odot \mathbf{U}_{Ori}^{t,n+1}, \\ \mathbf{U}_{CSA}^{t,n+1} &= g_s(\mathbf{U}_{CA}^{t,n+1}) \odot \mathbf{U}_{CA}^{t,n+1}, \\ \mathbf{U}^{t,n+1} &= \mathbf{U}_{CSA}^{t,n+1} + \mathbf{U}^{t,n-1}, \end{aligned} \quad (15)$$

and Att-Res-SNN-2 can be described as

$$\begin{aligned} \mathbf{U}_{CA}^{t,n+1} &= g_c(\mathbf{U}_{Ori}^{t,n+1} + \mathbf{U}^{t,n-1}) \odot (\mathbf{U}_{Ori}^{t,n+1} + \mathbf{U}^{t,n-1}), \\ \mathbf{U}_{CSA}^{t,n+1} &= g_s(\mathbf{U}_{CA}^{t,n+1}) \odot \mathbf{U}_{CA}^{t,n+1}, \\ \mathbf{U}^{t,n+1} &= \mathbf{U}_{CSA}^{t,n+1}, \end{aligned} \quad (16)$$

where $\mathbf{U}_{Ori}^{t,n+1}$ is the original output of the basic Res-SNN block, $\mathbf{U}_{CA}^{t,n+1}$ and $\mathbf{U}_{CSA}^{t,n+1}$ are the output of the CA and CSA module, respectively, $\mathbf{U}^{t,n+1}$ is the final output of Att-Res-SNN block. To keep event-driven nature and avoid degradation deficiency in deep SNNs concurrently, the design of attention location is also important in residual SNNs. We recommend Att-Res-SNN-1 as a scheme for attention residual learning. Analysis of basic Res-SNN block selection and ablation studies of attention residual SNNs are conducted in Section 6.4. Backpropagation gradient involvement of Att-Res-SNN-1 and Att-Res-SNN-2 are discussed in Section 7.1.

4 ANALYSIS OF ENERGY CONSUMPTION

Times of floating-point operations (FLOPs) are used to estimate computational burden in CNNs, where almost all FLOPs are MAC. For an SNN, the measure of energy cost is relatively complicated because FLOPs of the first encoder layer are MAC while all other Conv or FC layers are AC, please see Table S1 in Supplementary Materials (SM). For MA-SNN, we use MA to regulate membrane potentials, which in turn drops the spiking activity. Thus, the energy increase comes from MAC operations due to attention. The energy decrease comes from the drop of AC operations caused by sparser spiking activity. We here focus on evaluating the energy shift between vanilla and attention SNNs.

Energy Cost of Vanilla SNNs. In the encoder layer of vanilla SNNs ($n = 1$), FLOPs are MAC operations that are the same as CNNs, because the work of this layer is to transform analog inputs into spikes. In addition, all other Conv and FC layers transfer spikes and execute AC operations to accumulate weights of postsynaptic neurons. Thus, the inference energy cost of a vanilla SNN E_{Base} can be quantified as

$$\begin{aligned} E_{Base} &= E_{MAC} \cdot FL_{SNNConv}^1 \\ &+ E_{AC} \cdot \left(\sum_{n=2}^N FL_{SNNConv}^n + \sum_{m=1}^M FL_{SNNFC}^m \right), \end{aligned} \quad (17)$$

where N and M are the total number of layers of Conv and FC, E_{MAC} and E_{AC} represent the energy cost of MAC and AC operation, $FL_{SNNConv}^n$ and FL_{SNNFC}^m are the FLOPs of n -th Conv and m -th FC layer, respectively. Refer to previous SNN works [14], [55], [56], we assume the data for various operations are 32-bit floating-point implementation in 45nm technology [57], in which $E_{MAC} = 4.6pJ$ and $E_{AC} = 0.9pJ$.

Additional Model and Computational Complexity. Additional parameters and computational burden induced by

three dimensions of attention modules are shown in Table S2 (see SM). The additional parameters are solely from the two FC layers (TA, CA) or one Conv layer (SA), and therefore constitute a small fraction of the total network capacity. The additional computation burden includes two parts, Δ_{MAC1} and Δ_{MAC2} , where the former comes from generating attention weights and the latter derives from refinement membrane potentials.

Energy Shift of Attention SNNs. By optimizing the membrane potential, the attention mechanism drops the spiking activity of SNNs in both Conv and FC layers. The computation formulas of AC reduction are shown in Table S2 of SM. We can estimate the shift of the energy cost versus the additional computational burden $\Delta_{MAC} = \Delta_{MAC1} + \Delta_{MAC2}$ and the decreased AC operations Δ_{AC} to demonstrate the energy efficiency of attention SNNs. The absolute energy shift between vanilla and attention SNNs can be computed as

$$\Delta_E = E_{MAC} \cdot \Delta_{MAC} - E_{AC} \cdot \Delta_{AC}. \quad (18)$$

We term the attention SNN energy consumption as E_{Att} . With the vanilla SNN as the anchor, the energy efficiency of an attention SNN is defined as

$$r_{EE} = \frac{E_{Base}}{E_{Att}} = \frac{E_{Base}}{E_{Base} + \Delta_E}. \quad (19)$$

The higher the r_{EE} , the greater the energy efficiency. Generally, we represent the r_{EE} of baseline model as $1 \times$.

Network Average Spiking Activity Rate (NASAR). r_{EE} can estimate energy shift in a fine granularity manner, while it is difficult to use for the principle analysis of attention mechanisms. We want to design a simple and intuitive indicator to show the shift in energy cost and explore the underlying reasons.

Spike counts of SNNs are positively correlated with energy cost and can be used to roughly measure energy cost. A spike corresponds to some AC operations, and how much is associated with the design of network architecture (i.e., $FL_{SNNConv}^n$ and FL_{SNNFC}^m depend on architecture). Definitely, the lower the spike counts, the better the energy efficiency. To compute the spike counts, we define the network average spiking activity rate (NASAR) as follows: at time step t , a spiking network's spiking activity rate (NSAR) is the ratio of spikes produced over all the neurons to the total number of neurons in this step; then we define the NASAR which averages NSAR across all time steps T . Spike count is equal to $NASAR \cdot neuron\ number \cdot T$. Once the network architecture and time step are determined, $neuron\ number$ and T are fixed constants. So the spike counts are only related to NASAR. We mainly exploit NASAR and NSAR to approximately represent energy costs and explore how the attention module brings the energy shift for SNNs.

5 EXPERIMENTS

In this section, we investigate the *effectiveness* and *efficiency* of MA-SNN across a range of tasks, datasets, and architectures. We conduct extensive experiments on the datasets: DVS128 Gesture/Gait for event-based action recognition in Section 5.1; ImageNet-1K for static image classification in Section 5.2.

TABLE 1: Comparison with previous works on Gesture and Gait. The numbers or percentages in brackets denote the performance improvement or reduction proportions of spiking activity over the re-implementation baseline [43] and [38]. The format of “ $A_{\pm a}$ ” represents the mean and variance of the accuracy of repeat ten independent experiments.

Datasets	Work	$dt \times T$	Acc. (%)	r_{EE}	Neuron number	Roughly Energy Evaluation	
						NASAR	Spike Counts ($\times 10^6$)
DVS128 Gesture	Amir <i>et al.</i> 2017 [42]	1×120	92.59	-	-	-	-
	Shrestha <i>et al.</i> 2019 [58]	5×300	93.64	-	-	-	-
	Zheng <i>et al.</i> 2021 [15]	30×40	96.87	-	229,376	0.146	1.337
	Fang <i>et al.</i> 2021 [16]	375×16	97.92	-	1,922,304	0.053	1.636
	Yao <i>et al.</i> 2021 [43] (Vanilla SNN)	15×60	$90.63_{\pm 0.58}$	$1 \times$	106,763	0.173	1.108
	+ TCSA (This Work)		$96.53_{\pm 0.57}$ (+5.9)	$3.4 \times$	106,763	0.026 (-84.9%)	0.167 (-84.9%)
DVS128 Gait	Fang <i>et al.</i> [38] (Vanilla SNN)	300×20	$95.58_{\pm 1.00}$	$1 \times$	2,793,472	0.074	4.134
	+ TCSA (This Work)		$98.23_{\pm 0.46}$ (+2.7)	$1.8 \times$	2,793,472	0.011 (-85.1%)	0.615 (-85.1%)
	Wang <i>et al.</i> 2019 [59]	4400×1	89.9	-	-	-	-
	Wang <i>et al.</i> 2021 [44]	4400×1	94.9	-	-	-	-
	Yao <i>et al.</i> 2021 [43] (Vanilla SNN)	15×60	$87.59_{\pm 0.47}$	$1 \times$	106,772	0.245	1.570
	+ TCSA (This Work)		$92.29_{\pm 1.14}$ (+4.7)	$3.2 \times$	106,772	0.045 (-81.6%)	0.288 (-81.6%)
DVS128 Gait	Fang <i>et al.</i> [38] (Vanilla SNN)	220×20	$89.87_{\pm 1.32}$	$1 \times$	2,793,472	0.051	2.849
	+ TCSA (This Work)		$92.78_{\pm 0.79}$ (+2.9)	$1.7 \times$	2,793,472	0.013 (-74.5%)	0.726 (-74.5%)

To perform better apple-to-apple comparisons, we first re-implementation all the reported performance of networks [14], [38], [43] in the PyTorch framework and set them as our baselines. When training the baseline (vanilla) models or MA-integrated models, we follow their training schemes (i.e., hyper-parameter settings), if not otherwise specified. Throughout all experiments, we verify that MA outperforms all the vanilla models in both performance and efficiency without bells and whistles, then the general applicability of MA across different architectures and as well as different tasks are demonstrated. We here directly give our recommended combination of attention dimensions. Ablation studies of attention design are given in Section 6.

5.1 Event-based Action Recognition

5.1.1 Experimental Setup

Datasets. DVS128 Gesture [42] and DVS128 Gait [59] are both event stream datasets captured by the DVS camera, which has the same μs level temporal resolution and 128×128 spatial resolution but have different visual information. They capture the natural human motion, gestures and gaits. Gesture contains 11 kinds of hand gestures from 29 subjects under 3 kinds of illumination conditions, and there are 1,342 samples that each gesture has an average duration of 6 seconds. Gait contains various gaits from 21 volunteers (15 males and 6 females) under 2 kinds of viewing angles, and it records 4,200 samples where each gait has an average duration of 4.4 seconds.

Learning. We adopt the same experimental setup for vanilla SNNs and TCSA-SNN (here, we select TCSA as the MA), including network structure, hyper-parameters, learning methods, etc. MA is a plug-and-play module, all we have done is just add it to the vanilla models. We used *identical experimental setup* to analyze the effectiveness and efficiency of the same network on different datasets. All learning methods follow [43], and details are given in Section S2.1. Two kinds of vanilla structures are performed for Gesture and Gait to evaluate the energy shift induced by different backbones. One has three Conv layers, following [43]. The other has five Conv layers, following [38].

Moreover, we set various output latency $t_{lat} = dt \times T$ to investigate the effect of TCSA-SNN under the multi-scale constraints of latency.

5.1.2 Results and Analysis

Effectiveness. In Table 1, we report the accuracy of each vanilla model and its attention counterpart, and compare TCSA-SNN with previous works. We observe that in every comparison, TCSA-SNN outperforms the vanilla architectures, suggesting that the benefits of attention modules are not confined to a single event-based dataset, limited base architecture, or fixed output latency. TCSA-SNN performs better performance on the three-layer vanilla model (following [43]), which yields good gains of 5.9% and 4.7% on Gesture and Gait, respectively. The gains are consistent with the five-layer vanilla model (following [38]). Thus, attention can effectively facilitate the performance of lightweight plain SNN models for event-based tasks.

Efficiency. From Table 1, we observe a significant improvement in inference efficiency between our TCSA-SNN and vanilla models. On Gesture, three-layer TCSA-SNN (following [43]) and five-layer TCSA-SNN (following [38]) increase the energy efficiency by up to $3.4 \times$ and $1.8 \times$, respectively. We see a similar trend with regard to the effect of TCSA on Gait, finding that the TCSA-SNN outperforms our re-implemented counterpart baselines with both effectiveness and energy efficiency. Moreover, we observe that the NASAR is associated with the model size. Networks with more spiking neurons usually have sparser spiking activity, such as the NASAR of the three-layer and five-layer baseline on Gesture are 0.214 and 0.074, respectively. Incredibly, on Gesture, the NASAR of five-layer SNN drops to 0.011 with the help of TCA, which means that *only 1.1% spiking neurons are activated at each time step*.

It should be noted that only when the benefits outweigh the costs, i.e., $E_{AC} \cdot \Delta_{AC} > E_{MAC} \cdot \Delta_{MAC}$, the energy cost can be dropped, and all of our attention designs in this paper meet this condition (see Section 6.3). Moreover, the energy efficiency of attention SNNs is associated with the backbone structure (E_{base}). On Gait, the spike counts of

TABLE 2: Comparison with previous works on ImageNet-1K. * The input crops are enlarged to 288×288 in inference.

Methods	Work	Model	Time step (T)	Top-1 Acc.(%)	Energy Efficiency (r_{EE})
ANN-to-SNN	Sengupta <i>et al.</i> 2019 [60]	VGG-16	2500	69.96	-
		ResNet-34	2500	65.47	-
	Han <i>et al.</i> 2020 [61]	VGG-16	4096	73.09	-
		ResNet-34	4096	69.89	-
	Wu <i>et al.</i> 2021 [33]	AlexNet	16	55.19	-
		VGG-16	16	65.08	-
		VGG-16	2048	75.32	-
Li <i>et al.</i> 2021 [35]	VGG-16	2048	75.32	-	
Stöckl <i>et al.</i> 2021 [34]	ResNet-50	500	75.10	-	
Bu <i>et al.</i> 2022 [62]	VGG-16	512	74.69	-	
Direct Training	Zheng <i>et al.</i> 2021 [15]	ResNet-50	6	64.88	-
		Wide-ResNet-34	6	67.05	-
	Fang <i>et al.</i> 2021 [16]	ResNet-34	4	67.04	-
		ResNet-101	4	68.76	-
	Deng <i>et al.</i> 2022 [63]	ResNet-34	4	68.00	-
Direct Training	Hu <i>et al.</i> 2021 [14] (SOTA backbone)	ResNet-18	6	63.10	$3.3 \times$
Backpropagation	This Work (+CSA)	ResNet-18	1	63.97	$29.7 \times$
Backpropagation	ANN	ResNet-18	-	69.76	$1 \times$
Direct Training	Hu <i>et al.</i> 2021 [14] (SOTA backbone)	ResNet-34	6	69.42	$3.8 \times$
Backpropagation	This Work (+CSA)	ResNet-34	1	69.15	$29.6 \times$
Backpropagation	ANN	ResNet-34	-	73.30	$1 \times$
Direct Training	Hu <i>et al.</i> 2021 [14] (SOTA backbone)	ResNet-104*	5	76.02	$5.3 \times$
Backpropagation	This Work (+CSA)	ResNet-104*	4	77.08	$7.4 \times$
Backpropagation	This Work (+CSA)	ResNet-104*	1	75.92	$31.8 \times$
Backpropagation	ANN [14]	ResNet-104	-	76.87	$1 \times$

TABLE 3: Comparison with baselines on ImageNet-1K (inference spatial resolution is 224×224). Based on the MS-Res-SNN [14], we re-implement various baseline structures and their attention counterpart with $T = 1$ and report corresponding performance and energy shift. The case of $T = 1$ can be regarded as pre-training of multi-time step SNNs [64]. Note that, compared with attention CNNs, the accuracy improvement of the attention on MS-Res-SNN is very significant, e.g., using identical CSA for Res-CNN-34 [28] and Res-SNN-34 can improve the accuracy by +0.7 and +5.0 percent, respectively.

Model ($T = 1$)	Top-1 Acc. (%)	r_{EE}	NASAR
MS-Res-SNN-18 [14]	61.70	$1 \times$	0.224
+ CA (This work)	63.42 (+1.7)	$1.4 \times$	0.165 (-26.3%)
+ CSA (This work)	63.97 (+2.3)	$1.5 \times$	0.148 (-33.9%)
MS-Res-SNN-34 [14]	64.13	$1 \times$	0.203
+ CA (This work)	67.96 (+3.8)	$1.2 \times$	0.167 (-17.7%)
+ CSA (This work)	69.15 (+5.0)	$1.3 \times$	0.153 (-24.6%)
MS-Res-SNN-104 [14]	71.57	$1 \times$	0.218
+ CA (This work)	72.03 (+0.5)	$1.1 \times$	0.195 (-10.6%)
+ CSA (This work)	73.82 (+2.3)	$1.1 \times$	0.201 (-7.8%)

the three-layer baseline are reduced by 81.6% resulting in $3.2 \times$ better energy efficiency. By contrast, in the five-layer baseline, the reduction of spike counts and the value of r_{EE} are 74.5% and $1.7 \times$, respectively. The underlying reason is that there are more MAC operations in the first encoding layer of the five-layer baseline, which is induced by the architecture design and leads to higher E_{base} .

5.2 Static Image Classification

5.2.1 Experimental Setup

Datasets. ImageNet-1K [65] is the most typical static image dataset, which is widely used in the field of image classification. It provides a large-scale natural image dataset containing a total of 1,000 categories, which consists of 1.28 million training images and 50k test images.

Learning. Similar to event-based recognition, each baseline architecture and its corresponding MA counterpart for ImageNet is trained with identical optimization schemes. We opted to use MS-ResNet-SNN [14] architectures as strong baselines to assess the effectiveness and efficiency of attention modules and follow the training and evaluation protocols described in [14]. All models are trained from re-implement scratch. All learning details are given in Section S2.2 of SM. We exploit channel and spatial attention for ImageNet-1K, whose temporal information is weak.

5.2.2 Results and Analysis

Effectiveness. We would first set $T = 1$ for all experiments which can be regarded as pre-training of multi-time step SNNs [64] and the results of top-1 accuracy on ImageNet-1K are reported in Table 3. We begin by comparing Att-Res-SNN against Res-SNN baselines with different depths. We observe that attention modules consistently improve performance across different depths. Remarkably, CSA-Res-SNN-34 exceeds Res-SNN-34 by +5.0 percent. By contrast, an identical CSA module used for Res-CNN-34 induces only +0.7 percent performance gain [28]. In Table 2, we make a comprehensive comparison with prior works on ImageNet-1K, including ANN-to-SNN, direct training SNN, and ANN.

At the same network depth, we see that CSA-Res-SNN with single-time step can exceed or approach prior SOTA results [14], which uses multi-time steps. For example, CSA-Res-SNN-18, with an accuracy of 63.97%, is better than the 63.10% of its counterpart 6-time step backbone. The single-time step accuracy advantage of Res-SNN brought by the attention can also be extended to the deeper backbones (Res-SNN-34 and Res-SNN-104). Single-time step CSA-Res-SNN-34/104 has an accuracy of 69.15%/75.92%, which is slightly inferior to its 6/5-time step counterpart Res-SNN-34/104 backbone (69.42%/76.02%). We further extend CSA-Res-SNN-104 to 4-time steps and obtain an accuracy of 77.08% on ImageNet-1K, which is the SOTA result in the SNN domain (including direct training and ANN-to-SNN), and better than the ANN with the same architecture(76.87%). In contrast to SOTA ANN-to-SNN, our model have better accuracy (77.08% vs. 75.32%) and lower latency (4-time step vs. 2048-time step).

Inference Efficiency. From Table 3 ($T = 1$), we observe that employing the attention can drop the energy cost and improve performance simultaneously. For example, compared with Res-SNN-18, the NASAR of CSA-Res-SNN-18 can be reduced from 0.224 to 0.148, which means that each neuron emits only a 0.15 spike on average. Meanwhile, the task accuracy is improved from 61.70% to 63.97%, and the energy efficiency is up to $1.5\times$. Different from the plain SNNs in event-based tasks, additional energy cost ($E_{MAC} \cdot \Delta_{MAC}$) caused by the attention module in Res-SNNs can be approximate ignored since the high E_{base} , e.g., 33.9% decrease of spiking activity in Res-SNN-18 incurs 31.7% reduction of energy cost ($1.5\times$ better energy efficiency). In Table 2, we set the energy cost of ANN as the E_{base} ($1\times$) and assess the energy cost of ANN/SNN with the same structures. As opposed to the Res-ANN-104, 5-time step Res-SNN-104 and 4-time step CSA-Res-SNN-104 achieve up to $5.3\times$ and $7.4\times$ better energy efficiency, respectively. Dramatically, single-times step CSA-Res-SNN-104 has up to $31.8\times$ better compute energy efficiency compared to the Res-ANN-104. The results reported in Table 2 show that deep SNNs have significant energy efficiency advantages over deep ANNs. Especially the single-step deep SNNs can amplify this advantage.

Training Efficiency. In this paper, we first train large-scale Att-Res-SNN by a single-time step and found that it can yield better or comparable accuracy than multi-time steps Res-SNN. Meanwhile, in terms of training efficiency, we hold a $6.4\times$ training acceleration and less hardware requirement (details in Section S2.3 of SM). To the best of our knowledge, there are currently only two works [14], [16] that directly train SNN with more than 100 layers on ImageNet-1K. Because the long training time and enormous hardware resources caused by large-scale and multi-time steps are prohibitive. Existing work attempts to alleviate this dilemma by designing complex training tricks [66] or complicated new spiking neuron [67]. Our experimental results show another simple potential solution to break the predicament, which includes two steps. First, the SNN community can focus on the single-step simulation algorithm design of large-scale SNNs, because it requires lower training time and less hardware. Our experimental results show that this is feasible. Attention can help single-step SNNs outperform

TABLE 4: Effect of attention locations in three-layer SNN [43] at different dimensions on Gesture with $dt = 15$, $T = 60$

Model	Attention Location	Acc. (%)	NASAR
Vanilla SNN [43]	-	90.63 \pm 0.58	0.174
TA-SNN	Conv-PRE	91.01 \pm 0.69	0.074
	Conv-POST	92.60 \pm 0.47	0.073
CA-SNN	Conv-PRE	91.84 \pm 0.60	0.086
	Conv-POST	91.58 \pm 0.29	0.097
	Activate-PRE	93.88\pm0.34	0.072
SA-SNN	Conv-PRE	92.85 \pm 0.78	0.060
	Conv-POST	92.57 \pm 0.28	0.062
	Activate-PRE	92.50 \pm 0.47	0.058

or approach multi-step SNNs, which also demonstrates the great potential of SNNs. Then, the researchers separately focus on how to effectively extend single-step SNNs to multi-step SNNs, such as using the pre-train method [64]. We believe that reducing training time and hardware costs will help the development of SNN fields.

6 ABLATION STUDY

We conduct ablation experiments to gain a better understanding of the effectiveness and efficiency of adopting different configurations on the attention design for SNNs. All ablation experiments are performed by following the experimental setup in section 5, if not otherwise specified.

6.1 Attention Locations in Plain SNNs

Lightweight plain SNNs are suitable for real scenarios requiring latency, accuracy, and energy cost. We explore the effectiveness and spiking activity of SNNs with various attention locations based on the Gesture. The baseline architecture follows the three-layer SNN in [43]. The attention location design of SNNs is more sophisticated than CNNs. Firstly, the event-driven nature of SNNs will be destroyed if we copy the attention location in CNN (behind the activate operation, see Fig. 4). Secondly, the goal of attention is to optimize the membrane potentials, which can be achieved by exploiting various attention dimensions separately or simultaneously. Thirdly, we need to examine the effect of the attention dimension on its location. Comprehensively, we consider three possible attention locations: (1) Conv-PRE, in which the attention module is inserted before the Conv operation; (2) Conv-POST, in which the attention module is moved between the Conv and the integration operation and (3) Activate-PRE, in which the attention unit acts on the integrated membrane potential. These variants are illustrated in Fig. 4. Individual analyses and assessments are performed on these variants according to attention dimensions.

TA Location. Firstly, TA location cannot be Activate-PRE, i.e., TA cannot work on the membrane potential which has already aggregated spatio-temporal information, because it is impossible to recalibrate the state that has occurred. Performance of Conv-PRE and Conv-POST for TA are reported in rows 3 and 4 of Table 4. We observe that both Conv-PRE and Conv-POST can improve performance and drop NASAR. The usage of the Conv-POST leads to better

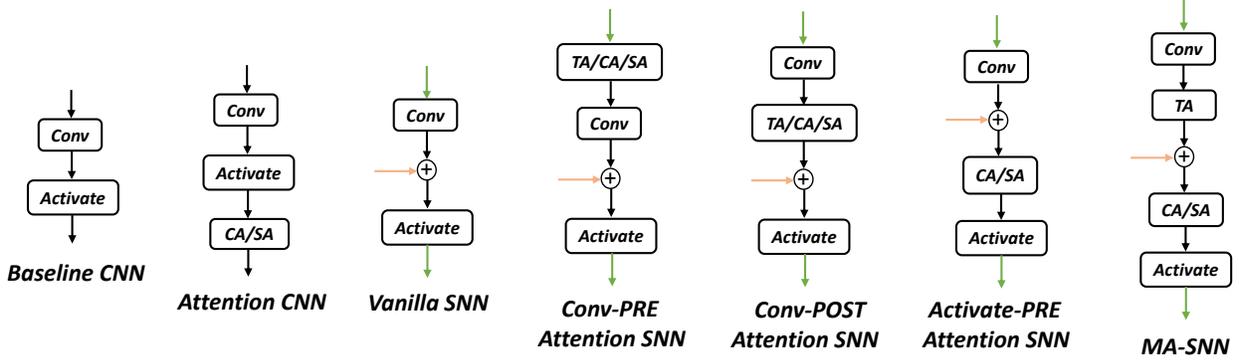


Fig. 4: Attention Locations in Plain SNN. From left to right: baseline CNN, processing spatial information. Attention CNN, performing attention module after activation [19], [28]. Vanilla SNN, processing spatio-temporal information. Conv-PRE Attention SNN, inserting attention modules before the Conv operation. Conv-POST, acting attention modules after the Conv operation while before the spatio-temporal integration. Activate-PRE, executing attention modules on the integrated membrane potential. MA-SNN, our recommended attention locations of three dimensions in plain SNN.

performance with a similar NASAR to Conv-PRE. So we select Conv-POST as the TA location.

CA Location. All location variants can be used to insert CA. We assess these locations and report results in Table 4. We see that incorporating the Activate-PRE CA for SNNs can obtain the best performance and lowest NASAR. This suggests that using CA to directly optimize spatial-temporal fused membrane potential is more effective and efficient than only optimizing spatial features.

SA Location. We finally explore the effect of using SA alone. The results are reported in the last three rows of Table 4. We observe that individual SA can also improve performance and drop NASAR without sensitivity to the location. In the practice of attention CNNs, SA is usually performed together with CA and follows CA [28]. We adopt this serial setting in this paper and set the Activate-PRE location for SA.

The comparison in Table 4 shows that effectiveness and efficiency are *robust* at various dimensions to a range of attention locations. Performance always goes up no matter which location is chosen. While specific gains of accuracy and NASAR vary from location to location. These experimental results also confirm our point that optimized membrane potential can induce sparser spiking activity and better task performance simultaneously. Comprehensively considering effectiveness and efficiency, our recommended attention locations are shown at the rightmost in Fig. 4.

6.2 Different Attention Variants

Attention modules generally concern three metrics for model assessment and comparison: accuracy, additional computational burden, and parameters. Many attention modules are designed to meet various metrics limitations in different scenes. Here we investigate three typical attention modules and apply them to SNNs. One of the representative works is CBAM [28], which stacks channel and spatial attention in series to achieve higher accuracy. To reduce the computational burden of the two-layer FC operation in classic CBAM, ECANet [46] uses a 1D convolution to model the interaction between channels. In addition, attention can be viewed as the process of feature optimization,

TABLE 5: Effect of Different attention modules in three-layer SNN [43] on Gesture with $dt = 15$, $T = 60$

Design	Acc. (%)	Params (\uparrow)	NASAR	Δ_{MAC} (\uparrow)
Baseline SNN [43]	90.63 \pm 0.58	-	0.173	-
CBAM [28]-SNN	93.88\pm0.34	+4,608	0.072	+276,480
ECA [46]-SNN	92.81 \pm 0.40	+48	0.071	+76,800
SimAM [31]-SNN	91.98 \pm 0.17	0	0.083	+25,559,043

e.g., SimAM [31] proposes an attention module based on mathematics and neuroscience theories with parameter-free. We explore the influence of the above three modules by directly integrating them into three-layer baseline SNN [43], where we adopt the same Activate-PRE location.

The metrics of attention SNNs are compared to vanilla SNNs in Table 5. We note that CBAM, ECA, and SimAM all perform well to improve performance and simultaneously lead to sparser spiking activity. CBAM brings the highest performance gains. SimAM has no additional parameters and obtains a maximum additional computational burden. ECA has the smallest Δ_{MAC} and negligible additional parameters. The NASARs of these attention SNNs are close. This experiment suggests that the benefits of the performance improvements and sparser spiking activity produced by optimizing membrane potentials are fairly *robust* to different attention mechanisms. In practice, we can choose the attention module according to the requirement of specific scenarios. To achieve better performance, we use CBAM as the basis attention module of SNNs in this paper.

6.3 Combinations of Attention Dimension

We perform an ablation study to assess the effectiveness and energy cost of the combination of various attention dimensions when executing them together. Results are reported in Table 6. Simple superposition of various attention dimensions can further boost the representation power and energy efficiency of SNNs, e.g., three-dimensional attention TC-SA-SNN achieves the highest accuracy gain of +5.9 and the high energy efficiency of 3.4 \times than vanilla SNN. The performance of TC-SA-SNN also exceeds all two-dimensional or

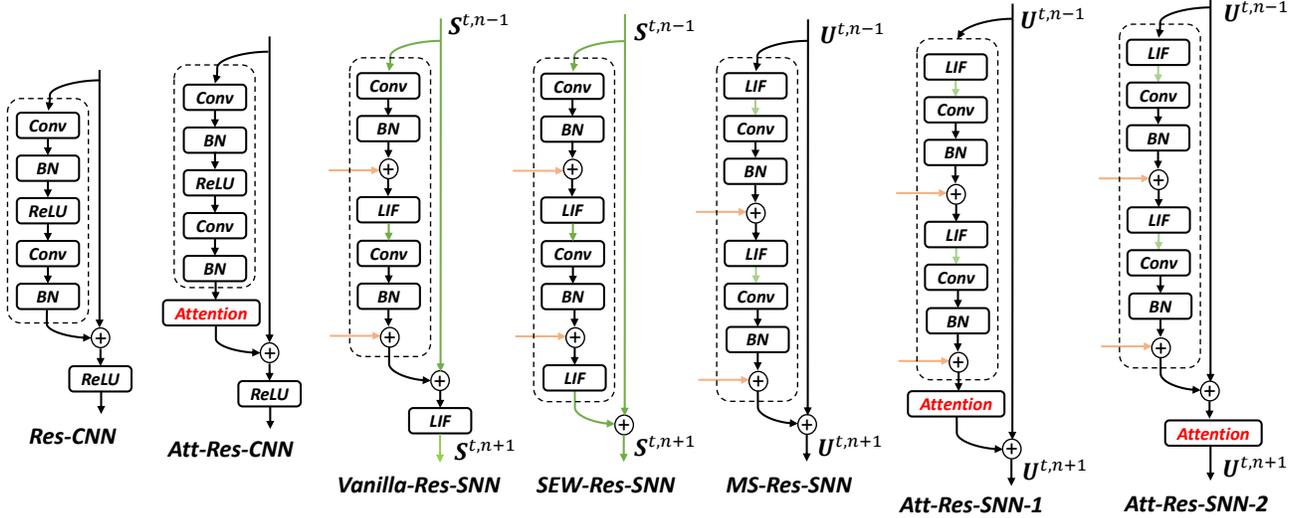


Fig. 5: Attention Residual Learning for SNNs. From left to right: Res-CNN [2]. Classic Att-Res-CNN [19], [28]. Vanilla Res-SNN [15], executing the same shortcut and residual block as Res-CNN. SEW-Res-SNN [16], mainly building shortcuts between spikes from different layers. MS-Res-SNN [14], constructing a shortcut among membrane potential of spiking neurons in different layers. Att-Res-SNN-1 (our recommended method) and Att-Res-SNN-2, both select MS-Res-SNN as the backbone. The former performs attention between the residual block and shortcut connection, which is the same as classic Att-Res-CNN. The latter executes attention after the shortcut connection.

TABLE 6: Effect of attention dimensions in three-layer Conv-based SNN [43] on Gesture with $dt = 15$, $T = 60$

Model	Acc. (%)	r_{EE}	NASAR
Vanilla SNN [43]	90.63 \pm 0.58	1 \times	0.173
+ TA	92.60 \pm 0.47 (+2.0)	2.2 \times	0.073
+ CA	93.88 \pm 0.34 (+3.3)	2.2 \times	0.072
+ SA	92.50 \pm 0.47 (+1.9)	2.5 \times	0.058
+ TCA	95.73 \pm 0.49 (+5.1)	3.5 \times	0.034
+ TSA	94.31 \pm 0.50 (+3.7)	3.4 \times	0.030
+ CSA	94.79 \pm 0.82 (+4.2)	3.4 \times	0.030
+ TCSA	96.53\pm0.57 (+5.9)	3.4 \times	0.026

TABLE 7: Effect of Different residual attention locations in Res-SNNs with $T = 1$ on ImageNet-1K

Model	Acc. (%)	NASAR
Res-SNN-18 [14]	61.70	0.224
CSA-Res-SNN-18-1	63.97(+2.3)	0.148
CSA-Res-SNN-18-2	63.49(+1.8)	0.137
Res-SNN-34 [14]	64.13	0.203
CSA-Res-SNN-34-1	69.15(+5.0)	0.153
CSA-Res-SNN-34-2	68.36(+4.2)	0.131

single-dimensional attention SNNs. Similarly, effectiveness and efficiency are *robust* at various attention dimension combinations.

6.4 Attention Residual Learning SNNs

As discussed in Section 3.4, Vanilla Res-SNN [15], SEW-Res-SNN [16] and MS-Res-SNN [14] are the only three kinds of residual learning of SNN by direct training, which dedicate to conquering the degradation problem of deep SNNs.

Referring to our previous works [14], [16], a deeper model should have a training error no greater than its shallower counterpart if the added layers implement the identity mapping. Vanilla Res-SNN copies the experience of Res-CNN but can not obtain identity mapping since it makes a mismatching shortcut connection on membrane potential and spikes. The underlying reason is that CNNs activate analog values while SNNs activate spikes. To address this, SEW-Res-SNN and MS-Res-SNN respectively establish shortcuts between spikes or membrane potentials of different layers, where both can obtain identity mapping (see Fig. 5).

In fact, we think the shortcut connection in MS-Res-SNN is identical to our motivation for introducing the attention, which can also be seen as a way to optimize the membrane potentials. Furthermore, given that MS-Res-SNN has higher accuracy, we choose it as the backbone model. In addition to the proposed attention residual design that performs attention between the residual block and shortcut connection, we consider another variant, Att-Res-SNN-2, in which the attention is moved after the shortcut. These variants are illustrated in Fig.5 and the performance of each variant is reported in Table 7. We observe that both Att-Res-SNN-1 (our recommended method) and Att-Res-SNN-2 perform well (i.e., have *robustness*) on effectiveness and efficiency concretely. Furthermore, both of them can overcome the degradation problem in general deep SNNs, i.e., they can achieve dynamical isometry (details in Section 7.1). Moreover, Att-Res-SNN-1 is better in the effectiveness aspect, and Att-Res-SNN-2 has sparser spiking activity. Although it is beyond the scope of this work, we anticipate that further effectiveness and efficiency gains will be achievable simultaneously by tailoring backbone SNNs and attention module usage for specific tasks.

7 UNDERSTANDING AND VISUALIZING ATTENTION

We have shown that attention can concurrently boost the effectiveness and efficiency of a plain or deep SNN for various vision tasks. Here we provide an in-depth analysis of how an MA-integrated model (i.e., three-layer SNN + TCSA, Res-SNN-34 + CSA) may differ from its vanilla counterpart (i.e., three-layer SNN, Res-SNN-34). We first explore the gradient norm equality of attention deep SNNs by using the dynamical isometry framework proposed by Chen *et al.* [32]. Next, we provide visualization results of that the MA-SNN success to correctly classifying but the vanilla model fails. Then, we explore the effectiveness and efficiency of MA-SNN by the proposed spiking response visualization method. Finally, we investigate the change in spiking activity rate induced by attention.

7.1 Gradient Evolvement in Att-Res-SNNs

In recent years, dynamical isometry has been developed as a theoretical explanation of well-behaved neural networks. When a deep neural network is dynamical isometry, it can avoid gradient vanishing or explosion and every singular value of its input-output Jacobian matrix remains close to one. In this subsection, we analyze that both of Att-Res-SNN-1 and Att-Res-SNN-2 in Fig. 3 can achieve gradient norm equality with the help of block dynamical isometry framework [32]. In a nutshell, Att-Res-Net-1 and Att-Res-Net-2 can avoid the drawback of degradation problem and attain great stability constituting a much shallower network in effect than it appears to be for gradient norm.

Without loss of generality, a neural network can be viewed as a serial of blocks:

$$f(x_0) = f_{\theta^L}^L \circ f_{\theta^{L-1}}^{L-1} \circ \dots \circ f_{\theta^1}^1(x_0), \quad (20)$$

where θ^i is the parameter matrix of the i -th layer. For simplicity, we denote $\frac{\partial f^j}{\partial f^{j-1}}$ as \mathbf{J}_j . Let $\phi(\mathbf{J})$ be the expectation of $\text{tr}(\mathbf{J})$, and $\varphi(\mathbf{J})$ be $\phi(\mathbf{J}^2) - \phi(\mathbf{J})^2$.

Definition 1 (Block Dynamical Isometry). (*Definition 3.1 in [32]*) Consider a neural network that can be represented as Eq. 20 and the j -th block's Jacobian matrix is denoted as \mathbf{J}_j . If $\forall j$, $\phi(\mathbf{J}_j \mathbf{J}_j^T) \approx 1$ and $\varphi(\mathbf{J}_j \mathbf{J}_j^T) \approx 0$, the network achieves block dynamical isometry.

Lemma 1 (Shallow Network Trick). (*Proposition 5.8 in [32]*) Assuming that for each of L sequential blocks in a neural network, we have $\phi(\mathbf{J}_j \mathbf{J}_j^T) = \omega + \tau \phi(\widetilde{\mathbf{J}}_j \widetilde{\mathbf{J}}_j^T)$ where \mathbf{J}_j is its Jacobian matrix, ω and τ are two constants determined by the network structure. For a shallow λ -layer network, given $\lambda \in \mathbb{N}^+ < L$, if $C_L^\lambda (1 - \omega)^\lambda$ and $C_L^\lambda \tau^\lambda$ are small enough, the L -block network would be as stable as a λ -layer network when both networks have $\forall j$, $\phi(\mathbf{J}_j \mathbf{J}_j^T) \approx 1$.

Based on the definition of block dynamical isometry (Definition 1) and the shallow network trick (Lemma 1), we can judge whether Att-Res-SNN-1&2 can achieve gradient norm equality or not. The serial and parallel connections in neural networks are denoted as Lemma S1 (multiplication theorem) and Lemma S2 (addition theorem) in Section S3 of SM. Some commonly used network components in Att-Res-SNNs are summarized in Table 8 (Details are in Lemma S3 and Lemma S4 of Section S3).

TABLE 8: $\phi(\mathbf{J}\mathbf{J}^T)$ and $\varphi(\mathbf{J}\mathbf{J}^T)$ of ReLU, Conv, Orthogonal, and sigmoid operations in neural networks. Results are collected from Lemma S3 and Lemma S4 in Section S3.

Part	$\phi(\mathbf{J}\mathbf{J}^T)$	$\varphi(\mathbf{J}\mathbf{J}^T)$
ReLU	p	$p - p^2$
Conv	$c_{in} k_h k_w \epsilon^2$	-
Orthogonal	γ^2	0
Sigmoid	$\frac{1}{16}$	0

Theorem 1 (Gradient Norm Equality of Att-Res-SNNs). Assuming two kinds of Att-Res-SNN designs in Fig.3 consisting of L sequential blocks, they can both achieve block dynamical isometry that Att-Res-SNNs could be as stable as a λ -layer network which satisfies $\phi(\mathbf{J}_j \mathbf{J}_j^T) \approx 1$ and $\lambda \in \mathbb{N}^+ < L$.

Proof. According to Lemma S1, the whole network's Jacobian matrix can be decomposed into the multiplication of its blocks' Jacobian matrices. After integrating the attention module into the basic Res-SNN block, we expect that each Att-Res-SNN block should satisfy $\phi(\mathbf{J}_j \mathbf{J}_j^T) \approx 1$, which provides stable gradient evolvement. We first analyze the CA and SA blocks. Then we evaluate the CSA block consisting of a serial connection of CA and SA. Finally, we discuss Att-Res-SNN-1&2 that integrate CSA blocks.

The Jacobian matrix of channel (function $g_c(\cdot)$) and spatial (function $g_s(\cdot)$) attention block are denoted as \mathbf{J}_{CA} and \mathbf{J}_{SA} , respectively. Assuming \mathbf{W}_{c1}^n and \mathbf{W}_{c0}^n in Eq. 10 are satisfy the Haar orthogonal initialization method. According to Eq. 10, Eq. 12, and Lemma S1, $\phi(\mathbf{J}_{CA} \mathbf{J}_{CA}^T)$ and $\phi(\mathbf{J}_{SA} \mathbf{J}_{SA}^T)$ can be decomposed as follow:

$$\begin{aligned} \phi(\mathbf{J}_{CA} \mathbf{J}_{CA}^T) &= 2\phi(\mathbf{J}_{Sig} \mathbf{J}_{Sig}^T) \phi(\mathbf{J}_{W_{c1}} \mathbf{J}_{W_{c1}}^T) \\ &\quad \phi(\mathbf{J}_{ReLU} \mathbf{J}_{ReLU}^T) \phi(\mathbf{J}_{W_{c0}} \mathbf{J}_{W_{c0}}^T). \end{aligned} \quad (21)$$

and

$$\phi(\mathbf{J}_{SA} \mathbf{J}_{SA}^T) = \phi(\mathbf{J}_{Sig} \mathbf{J}_{Sig}^T) \phi(\mathbf{J}_{Conv} \mathbf{J}_{Conv}^T). \quad (22)$$

According to Table 8, we have

$$\phi(\mathbf{J}_{CA} \mathbf{J}_{CA}^T) = \frac{\gamma_{w1}^2 \gamma_{w0}^2 p}{8}, \quad (23)$$

and

$$\phi(\mathbf{J}_{SA} \mathbf{J}_{SA}^T) = \frac{c_{in} k_h k_w \epsilon^2}{16}. \quad (24)$$

Then we evaluate the CSA block. For simplicity, we suppose the input of $g_c(\cdot)$ is U , the output of channel refinement is U_c . The channel refinement of membrane potentials can be described as

$$U_c = g_c(U) \odot U = g_c(U)(IU), \quad (25)$$

where I is the identity tensor, and we have

$$\frac{\partial U_c}{\partial U} = I g_c(U) + \frac{\partial g_c}{\partial U}(IU). \quad (26)$$

Similarly, for spatial refinement of membrane potentials, we have

$$U_s = g_s(U) \odot U = g_s(U)(IU), \quad (27)$$

and

$$\frac{\partial U_s}{\partial U} = I g_s(U) + \frac{\partial g_s}{\partial U}(IU), \quad (28)$$

where U_s is the output of spatial refinement. Linking channel and spatial attention in tandem

$$U_{cs} = g_s(g_c(U) \odot U) \odot (g_c(U) \odot U), \quad (29)$$

where U_{cs} is the output of channel-spatial refinement. The Jacobian matrix of channel-spatial attention block \mathbf{J}_{CSA} is

$$\begin{aligned} \mathbf{J}_{CSA} &= \frac{\partial U_{cs}}{\partial U} = \frac{\partial U_{cs}}{\partial U_c} \frac{\partial U_c}{\partial U} \\ &= (I_{g_s}(U_c) + \frac{\partial g_s}{\partial U_c}(IU_c))(I_{g_c}(U) + \frac{\partial g_c}{\partial U}(IU)). \end{aligned} \quad (30)$$

Based on Lemma S1 and Table 8, we have

$$\begin{aligned} \phi(\mathbf{J}_{CSA}\mathbf{J}_{CSA}^T) &= \\ \phi((I_{g_s}(U_c) + \frac{\partial g_s}{\partial U_c}(IU_c))(I_{g_s}(U_c) + \frac{\partial g_s}{\partial U_c}(IU_c))^T) & \\ \phi((I_{g_c}(U) + \frac{\partial g_c}{\partial U}(IU))(I_{g_c}(U) + \frac{\partial g_c}{\partial U}(IU))^T) & \\ = \phi(\mathbf{J}_{SA}\mathbf{J}_{SA}^T)\phi((IU_c)(IU_c)^T) + \phi(I_{g_s}(U_c)(I_{g_s}(U_c))^T) & \\ (\phi(\mathbf{J}_{CA}\mathbf{J}_{CA}^T)\phi((IU)(IU)^T) + \phi(I_{g_c}(U)(I_{g_c}(U))^T)). & \end{aligned} \quad (31)$$

According to Table 8, the means of outputs from sigmoid function is $\frac{1}{2}$. Thus, $\phi(I_{g_c}(U)(I_{g_c}(U))^T)$, $\phi(I_{g_s}(U_c)(I_{g_s}(U_c))^T)$, and $\phi((IU_c)(IU_c)^T)$ are equal to $\frac{1}{4}$. The means of U are 0 and C (a constant) in Att-Res-SNN-1 (BN output) and Att-Res-SNN-2 (BN output plus attention output), respectively. Bring Eq. 24 and Eq. 23 into Eq. 31, then for Att-Res-SNN-1:

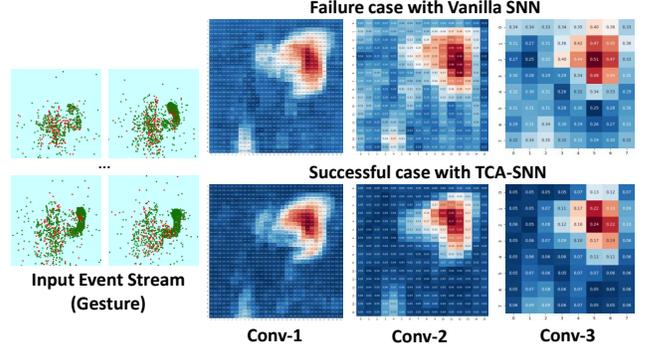
$$\begin{aligned} \phi(\mathbf{J}_{CSA}\mathbf{J}_{CSA}^T) &= \\ (\frac{1}{4} + \frac{c_{in}k_hk_w\epsilon^2}{16}\frac{1}{4})(\frac{1}{4} + \frac{\gamma_{w1}^2\gamma_{w0}^2p}{8} \times 0), & \end{aligned} \quad (32)$$

and for Att-Res-SNN-2:

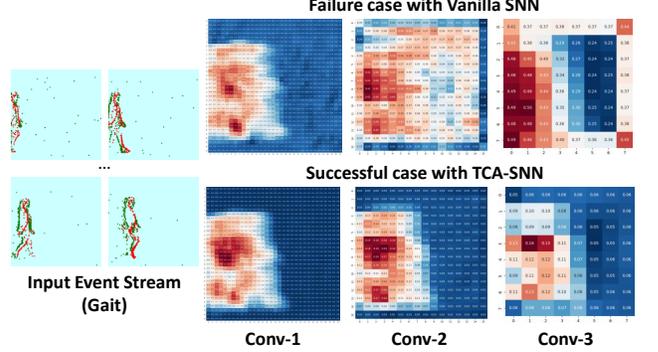
$$\begin{aligned} \phi(\mathbf{J}_{CSA}\mathbf{J}_{CSA}^T) &= \\ (\frac{1}{4} + \frac{c_{in}k_hk_w\epsilon^2}{16}\frac{1}{4})(\frac{1}{4} + \frac{\gamma_{w1}^2\gamma_{w0}^2p}{8} \times C). & \end{aligned} \quad (33)$$

Thus, for attention blocks (CA, SA, CSA), $\phi(\mathbf{J}\mathbf{J}^T) = 1$ can be achieved when ϵ , γ_{w1}^2 , γ_{w0}^2 , and p are set appropriately.

Finally, we consider the Att-Res-SNN-1 and Att-Res-SNN-2 that Res-SNN block and CSA block form the new attention residual block in Fig 3. For simplicity, we denote the Jacobian matrix of attention residual block and shortcut as \mathbf{J}_j and $\tilde{\mathbf{J}}_j$ respectively. According to Lemma 1, for Att-Res-SNN-1, we have $\phi(\mathbf{J}_j\mathbf{J}_j^T) = 1 + \zeta^2\phi(\mathbf{J}_{CSA}\mathbf{J}_{CSA}^T)\phi(\tilde{\mathbf{J}}_j\tilde{\mathbf{J}}_j^T)$ where ζ is from the linear scale transformation $\zeta x + \beta$ within the normalization at the bottom of the basic Res-SNN block. Att-Res-SNN-1 can be viewed as an extreme example of Lemma 1 with $(1 - \omega) \rightarrow 0$. Therefore $\forall \lambda$, $C_L^\lambda(1 - \omega)^\lambda$ is close to zero, and $C_L^\lambda\zeta^\lambda$ can be small enough for a given λ if ζ is initialized as a relative small value. In this way, the non-optimal block's error will be influential only within λ layers, and the Att-Res-SNN-1 will be as stable as a much shallower λ -layer network. For Att-Res-SNN-2, we can obtain $\phi(\mathbf{J}_j\mathbf{J}_j^T) = \phi(\mathbf{J}_{CSA}\mathbf{J}_{CSA}^T)(1 + \zeta^2\phi(\tilde{\mathbf{J}}_j\tilde{\mathbf{J}}_j^T))$. So if the basic residual block can achieve dynamical isometry, Att-Res-SNN-2 can also do it. \square



(a) Case study on DVS128 Gesture



(b) Case study on DVS128 Gait

Fig. 6: Case study on event-based tasks. We can observe that attention drives SNNs to focus on the target while the vanilla model shows more decentralized spiking activations.

7.2 Case Study

Understanding attention mechanisms by visualizing intermediate features [30] or attention heat maps [29] of a single sample are the two most common methods in CNNs. Activation values can directly generate the former, and the latter is produced by class activation mapping (CAM) [29]. Compared with these traditional methods, SNN's visualization is more natural since it has only two active statuses. In this paper, for a single event-based sample, we averaged all the 4D $([T, C, H, W])$ spiking maps of SNN into a 2D map $([H, W])$ over the temporal and channel dimension at each layer. Then we plot the 2D feature, which represents the average spiking response of every layer for this sample.

To visualize the effectiveness and efficiency of attention SNNs, we select two examples with regard to the case of the vanilla SNN failing in recognition but the attention SNN succeeds, where one from Gesture and the other from Gait. As shown in Fig. 6, each feature indicates the average spiking response of a layer of SNN. We make the following three observations about the effect of attention on SNN. First, the spiking activity is more concentrated in TCA-SNN, i.e., the red area of TCA-SNN is smaller and more focused. This suggests that attention is good for focusing on the important spatial region of intermediate channels. The second observation is that in the background region, the spiking response is suppressed. We see that attention darkens the color of the light blue area (background). The bluer the pixel, the closer the spiking activity rate is to

0. Finally, in all network layers, besides the obvious focus and suppression phenomena between vanilla and attention SNNs, we see a decrease in the overall spiking response induced by attention. For example, in Conv-3 of vanilla SNN and TCA-SNN on Gesture, the highest values are 0.51 and 0.24, respectively. These observations are consistent with the NASAR values of vanilla and attention SNNs in Table 6. Thus, by optimizing the membrane potential of spiking neurons, attention induces sparser spiking activity of SNNs.

7.3 Analysis of Overall Spiking Response

Classical CNN visualization methods before-mentioned generally can only be exploited to analyze a single sample such as one image, providing an intuitive feel of the attention mechanism. Previous work [68] extended CNN visualization methods to the SNN community, but did not get rid of the single-sample analysis limitation. On the other hand, existing attention works usually neglect the suppression part of the attention, which could dominate the network efficiency in SNNs. Although some works [28], [30], [69] have observed that attention would diminish background responses, they didn't realize the significance of this phenomenon. The underlying reason is that even if the background is suppressed to zero, it still needs to be computed and consumes energy for ANN on the GPU. By contrast, suppressing background noise is critical to the efficiency of SNNs because we observe that the background part has a higher spiking activity rate.

In this paper, we propose the average spiking response visualization (ASRV) method to demonstrate the spiking response distribution of SNNs on various datasets. Specifically, we first compute the spiking tensor $\mathbf{S}^{t,n}$ (spiking feature maps, only 0 or 1) for each sample of the validation set. Then average all spiking tensors to get average spiking response feature $\bar{\mathbf{S}}^{t,n} \in \mathbf{R}^{c_n \times h_n \times w_n}$, where each element of $\bar{\mathbf{S}}^{t,n}$ represents the spiking activity rate of a neuron on the validation set. Finally, we plot $\bar{\mathbf{S}}^{t,n}$ to visualize the spiking response.

Analysis of overall spiking response on DVS128 Gait.

In Fig. 7, we visualize average spiking response $\bar{\mathbf{S}}^{t,n}$ (take $n = 1, t = 1, 25, 49$ as examples) for various models based on Gait with $dt = 15, T = 60$, including vanilla SNN, TA-SNN, CA-SNN and TCA-SNN. Each pixel on the feature represents the spiking activity rate of one neuron over the whole validation set. We can clearly observe that attention modules drive the network to focus on the target and suppress the redundant background channels. For example, we see vanilla SNN has nine channels with a large area of red (spiking activity rates of these neurons are very close to 1) when $T = 49$, which means that these channels focus on background information. After integrating with attention modules, the background channels are significantly suppressed. Thereby the NASAR of SNNs is greatly reduced. Specifically, the NASAR of TA-SNN and CA-SNN has been respectively dropped from 0.245 to 0.091 and 0.103. Combining TA and CA (i.e., TCA) can further drop NASAR to 0.045 without invalid background information. In short, we discover that attention can significantly suppress unimportant background channels that contain a very high spiking activity rate, which in turn drops the energy

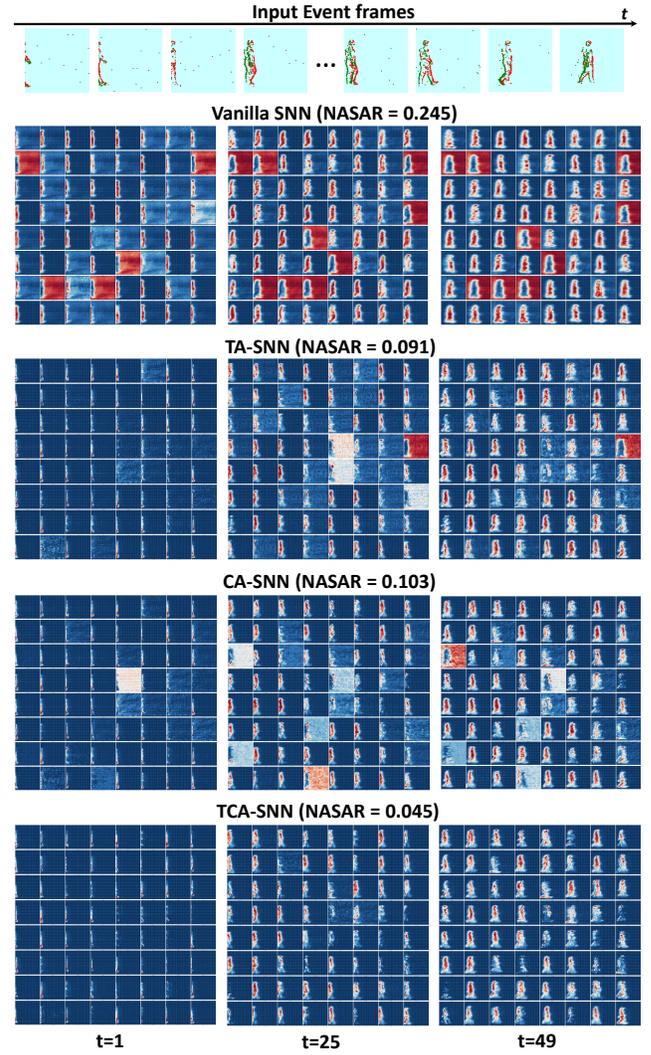


Fig. 7: Visualization of overall spiking response on Gait. From top to bottom: input event frames on Gait. Visualization of spiking response features in vanilla SNN, TA-SNN, CA-SNN, and TCA-SNN respectively, where $t = 1, 25, 49$ and $n = 1$. Each pixel on the feature represents the spiking activity rate of one neuron over the whole validation set. For a single channel, the redder the pixel, the higher spiking activity rate; the bluer the pixel, the closer the spiking activity rate is to 0. The dimension of $\bar{\mathbf{S}}^{t,1}$ (first layer) is $(32, 32, 64)$ at each time step, and we depict all 64 channels. We can clearly observe that attention drives the network to focus on the target and suppress the redundant background channels. Masking the background channel will significantly reduce the SNN energy cost, since the background contains a high spiking activity rate.

cost of SNNs. Crucially, the lower the NASAR, the higher the energy efficiency because the neuromorphic chip can skip the computation of zeros.

Analysis of overall spiking response on ImageNet-1K.

Each sample in the Gait is a person walking in front of a DVS camera, the only difference being that each person's gait was different. So in this single task dataset, we can easily observe the effect of attention on shallow plain SNNs.

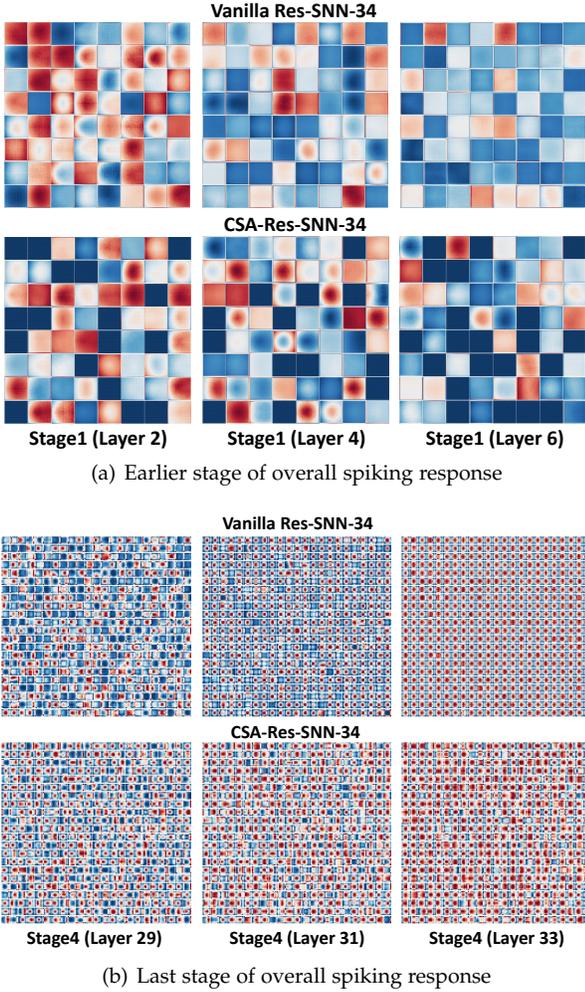


Fig. 8: Visualization of overall spiking response on ImageNet-1K. (a) Earlier stage. In contrast to the vanilla model, some dark blue (all neurons in these channels are not spiking at all) channels appear in the stage 1 of CSA-Res-SNN, which induces great energy efficiency. For example, there are 17, 13, 21 dark blue channels in layer 2, 4, and 6 respectively. (b) Last stage. With the assistance of attention, features (red regions) of the last Conv layer (layer 33) become more diverse, which helps to enhance the class selectivity of the deep model.

By contrast, when we compute the $\bar{S}^{t,n}$ over large-scale datasets such as ImageNet-1K with 1,000 categories, it is difficult to observe specific objects in the channels. Nevertheless, we execute our ASRV method on ImageNet-1K and give the visualization results in Fig. 8, and still can observe the suppression of background channels caused by attention in deep SNNs. It is well known that at the early stage, the filters of deep networks tend to extract class-agnostic low-level features while extracting class-specific high-level features at the last stage. To investigate the effect of attention on low-level and high-level feature extraction of deep SNNs, we plot average spiking features at stage 1 (layers 2, 4, and 6 with 64 channels) and stage 4 (layers 29, 31, and 33 with 512 channels).

We first observe stage 1 (earlier stage). The spiking

response of the SNN is averaged from 50,000 samples. It is hard to see specific objects in each channel. Here we focus on the shift of spiking response caused by attention. In contrast to the baseline model, some dark blue channels appear in CSA-Res-SNN. For example, there are 17, 13, and 21 dark blue channels in layers 2, 4, and 6, respectively. We check all these dark blue channels carefully, finding that all neurons in these channels have a firing rate of 0. That is to say, these dark blue channels are suppressed to zero by attention. This phenomenon is interesting, which indicates attention can suppress some low-level features (probably useless background noise information) and induces sparser spiking activity at the early stage of deep SNNs.

Then we observe stage 4 (last stage). At great depth, we see that the types of objects (red region in a channel) within each channel become diverse in Att-Res-SNN-34. Especially in the last Conv layer (layer 33), the red regions of objects in the channel of the attention model are significantly richer than the vanilla model counterpart. Previous works [19], [29] demonstrate that features at the last Conv layer are critical to correct classification. Namely, earlier layer (e.g., stage 1) features are typically more general and re-used within the network. While the later layer (e.g., stage 4) features exhibit great levels of specificity. We conjecture that the attention module helps class selectivity by diversifying features of the last Conv layer, which brings a significant performance gain to Res-SNN-34 (+5.0 percent).

7.4 Spiking Response of Attention SNNs

We plot the spiking response of vanilla SNN and TCA-SNN on Gait (Fig. S1 in Section S4 of SM), and we observe that the NASR of vanilla SNN is almost unchanged at each time step, which means SNN responds similarly to various inputs. This phenomenon is unreasonable because event streams are sparse and non-uniform [43]. With the help of data-dependent attention, the NASR of TCA-SNN is uneven and small at the temporal axis, which induces a much lower NASAR than vanilla SNN. Similar experimental results could also be found in Gesture. Furthermore, we count the NASAR values of vanilla Res-SNN and Att-Res-SNN on ImageNet-1K with $T = 1$. We find that the variance of NASAR is very small in vanilla Res-SNN, indicating the network’s response to different images is almost invariant. In contrast, the variance of NASAR in Att-Res-SNN is more prominent. These observations demonstrate that attention can produce instance-specific dynamic responses.

Actually, making the spiking activity of SNNs sparser is always a fascinating topic because the human brain is a paragon model of sparse and efficiency [7]. Current SNN models mainly drop NASAR by activity regularization [70] or network compression [71]. But forcing sparsity too much may hurt predictive performance for an equal number of neurons since the effective capacity of the model might be reduced [72]. Thus, in these *parameter regularization* methods, the reduction of NASAR is limited, or the performance only holds a slight improvement. By contrast, our strategy is to optimize the membrane potential of spiking neurons in a *data-dependent* way. We argue that reasonable membrane potential optimization based on specific features can naturally induce sparser spiking activity and better performance of SNNs concurrently.

8 CONCLUSION

In this work, we propose a lightweight attention module for SNNs, named multi-dimensional attention (MA), to boost both the performance and energy efficiency of SNNs. MA module is a plug-and-play that can be easily implemented and integrated with existing Conv-based SNNs. Inspired by attention theories in the neuroscience, our module optimizes the membrane potential of spiking neurons by learning when, what and where to focus and suppress through three separate pathways, which in turn drops the spiking activity and improves the performance. A wide range of experiments show the effectiveness and efficiency of MA, which achieve state-of-the-art performance and significant energy efficiency across multiple datasets and tasks, including event-based DVS128 Gesture/Gait and ImageNet-1K. We analyze how and why sparser spiking activity caused by attention is better, by visualizing the spiking response of vanilla and attention SNNs. We argue that effectiveness and efficiency exist as two sides of an elegant coin that should go hand in hand, which can be naturally symbiotic in SNNs, like coexistence in the human brain. The attention mechanisms can help us achieve this point.

ACKNOWLEDGMENTS

This work was partially supported by Beijing Natural Science Foundation for Distinguished Young Scholars (JQ21015) and National Key R&D Program of China (2018AAA0102600), and Beijing Academy of Artificial Intelligence (BAAI) and Pengcheng Lab.

REFERENCES

- [1] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [2] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [3] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information processing systems*, vol. 25, 2012.
- [4] J. Hirschberg and C. D. Manning, "Advances in natural language processing," *Science*, vol. 349, no. 6245, pp. 261–266, 2015.
- [5] A. Esteva, B. Kuprel, R. A. Novoa, J. Ko, S. M. Swetter, H. M. Blau, and S. Thrun, "Dermatologist-level classification of skin cancer with deep neural networks," *Nature*, vol. 542, no. 7639, pp. 115–118, 2017.
- [6] D. Silver, A. Huang, C. J. Maddison, A. Guez, L. Sifre, G. Van Den Driessche, J. Schrittwieser, I. Antonoglou, V. Panneershelvam, M. Lanctot *et al.*, "Mastering the game of go with deep neural networks and tree search," *Nature*, vol. 529, no. 7587, pp. 484–489, 2016.
- [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] S. B. Laughlin and T. J. Sejnowski, "Communication in neuronal networks," *Science*, vol. 301, no. 5641, pp. 1870–1874, 2003.
- [9] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [10] W. Maass, "Networks of spiking neurons: The third generation of neural network models," *Neural Networks*, vol. 10, no. 9, pp. 1659–1671, 1997.
- [11] J. Pei, L. Deng *et al.*, "Towards artificial general intelligence with hybrid tianjin chip architecture," *Nature*, vol. 572, no. 7767, pp. 106–111, 2019.
- [12] P. A. Merolla, J. V. Arthur, R. Alvarez-Icaza, A. S. Cassidy, J. Sawada, F. Akopyan, B. L. Jackson, N. Imam, C. Guo, Y. Nakamura *et al.*, "A million spiking-neuron integrated circuit with a scalable communication network and interface," *Science*, vol. 345, no. 6197, pp. 668–673, 2014.
- [13] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain *et al.*, "Loihi: A neuromorphic manycore processor with on-chip learning," *Ieee Micro*, vol. 38, no. 1, pp. 82–99, 2018.
- [14] Y. Hu, Y. Wu, L. Deng, and G. Li, "Advancing residual learning towards powerful deep spiking neural networks," *arXiv preprint arXiv:2112.08954*, 2021.
- [15] H. Zheng, Y. Wu, L. Deng, Y. Hu, and G. Li, "Going deeper with directly-trained larger spiking neural networks," *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, pp. 11 062–11 070, 2021.
- [16] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, and Y. Tian, "Deep residual learning in spiking neural networks," *Thirty-fifth Conference on Neural Information Processing Systems (NIPS2021)*, 2021.
- [17] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 20, no. 11, pp. 1254–1259, 1998.
- [18] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [19] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [20] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 10 012–10 022.
- [21] K. Han, Y. Wang, H. Chen, X. Chen, J. Guo, Z. Liu, Y. Tang, A. Xiao, C. Xu, Y. Xu, Z. Yang, Y. Zhang, and D. Tao, "A survey on vision transformer," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2022.
- [22] L. Yuan, Y. Chen, T. Wang, W. Yu, Y. Shi, Z.-H. Jiang, F. E. Tay, J. Feng, and S. Yan, "Tokens-to-token vit: Training vision transformers from scratch on imagenet," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 558–567.
- [23] S. Chevallier and P. Tarroux, "Covert attention with a spiking neural network," in *International conference on computer vision systems*. Springer, 2008, pp. 56–65.
- [24] K. C. Neokleous, M. N. Avraamides, C. K. Neocleous, and C. N. Schizas, "Selective attention and consciousness: investigating their relation through computational modelling," *Cognitive Computation*, vol. 3, no. 1, pp. 321–331, 2011.
- [25] S. Kundu, G. Datta, M. Pedram, and P. A. Beerel, "Spike-thrift: Towards energy-efficient deep spiking neural networks by limiting spiking activity via attention-guided compression," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2021, pp. 3953–3962.
- [26] F. Briggs, G. R. Mangun, and W. M. Usrey, "Attention enhances synaptic efficacy and the signal-to-noise ratio in neural circuits," *Nature*, vol. 499, no. 7459, pp. 476–480, 2013.
- [27] H. Spitzer, R. Desimone, and J. Moran, "Increased attention enhances both behavioral and neuronal performance," *Science*, vol. 240, no. 4850, pp. 338–340, 1988.
- [28] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "Cbam: Convolutional block attention module," in *Proceedings of the European conference on computer vision (ECCV)*, 2018, pp. 3–19.
- [29] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, "Learning deep features for discriminative localization," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2921–2929.
- [30] J. Park, S. Woo, J.-Y. Lee, and I. S. Kweon, "A simple and lightweight attention module for convolutional neural networks," *International Journal of Computer Vision*, vol. 128, no. 4, pp. 783–798, 2020.
- [31] L. Yang, R.-Y. Zhang, L. Li, and X. Xie, "Simam: A simple, parameter-free attention module for convolutional neural networks," in *International Conference on Machine Learning*. PMLR, 2021, pp. 11 863–11 874.
- [32] Z. Chen, L. Deng, B. Wang, G. Li, and Y. Xie, "A comprehensive and modularized statistical framework for gradient norm equality

- in deep neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.
- [33] J. Wu, C. Xu, X. Han, D. Zhou, M. Zhang, H. Li, and K. C. Tan, "Progressive tandem learning for pattern recognition with deep spiking neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [34] C. Stöckl and W. Maass, "Optimized spiking neurons can classify images with high accuracy through temporal coding with two spikes," *Nature Machine Intelligence*, vol. 3, no. 3, pp. 230–238, 2021.
- [35] Y. Li, S. Deng, X. Dong, R. Gong, and S. Gu, "A free lunch from ann: Towards efficient, accurate coding in spiking neural networks calibration," in *International Conference on Machine Learning*. PMLR, 2021, pp. 6316–6325.
- [36] E. O. Neftci, H. Mostafa, and F. Zenke, "Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51–63, 2019.
- [37] I. M. Comsa, K. Potempa, L. Versari, T. Fischbacher, A. Gesmundo, and J. Alakuijala, "Temporal coding in spiking neural networks with alpha synaptic function," in *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2020, pp. 8529–8533.
- [38] 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," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 2661–2671.
- [39] G. Gallego, T. Delbruck, G. Orchard, C. Bartolozzi, B. Tabá, A. Censi, and et al., "Event-based vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2020.
- [40] H. Rebecq, R. Ranftl, V. Koltun, and D. Scaramuzza, "High speed and high dynamic range video with an event camera," *IEEE transactions on pattern analysis and machine intelligence*, pp. 1–1, 2019.
- [41] W. Cheng, H. Luo, W. Yang, L. Yu, S. Chen, and W. Li, "Det: A high-resolution dvs dataset for lane extraction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2019, pp. 1666–1675.
- [42] A. Amir, B. Tabá, 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.
- [43] M. Yao, H. Gao, G. Zhao, D. Wang, Y. Lin, Z. Yang, and G. Li, "Temporal-wise attention spiking neural networks for event streams classification," in *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021, pp. 10221–10230.
- [44] Y. Wang, X. Zhang, Y. Shen, B. Du, G. Zhao, L. C. Cui Lizhen, and H. Wen, "Event-stream representation for human gaits identification using deep neural networks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–1, 2021.
- [45] Z. Gao, J. Xie, Q. Wang, and P. Li, "Global second-order pooling convolutional networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 3024–3033.
- [46] Q. Wang, B. Wu, P. Zhu, P. Li, W. Zuo, and Q. Hu, "Eca-net: Efficient channel attention for deep convolutional neural networks," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [47] M.-H. Guo, T.-X. Xu, J.-J. Liu, Z.-N. Liu, P.-T. Jiang, T.-J. Mu, S.-H. Zhang, R. R. Martin, M.-M. Cheng, and S.-M. Hu, "Attention mechanisms in computer vision: A survey," *Computational Visual Media*, pp. 1–38, 2022.
- [48] L. Deng, Y. Wu, X. Hu, L. Liang, Y. Ding, G. Li, and et al, "Rethinking the performance comparison between snns and anns," *Neural Networks*, vol. 121, pp. 294–307, 2020.
- [49] 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.
- [50] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *International conference on machine learning*. PMLR, 2015, pp. 448–456.
- [51] J. Moran and R. Desimone, "Selective attention gates visual processing in the extrastriate cortex," *Science*, vol. 229, no. 4715, pp. 782–784, 1985.
- [52] J. H. Maunsell, "Neuronal mechanisms of visual attention," *Annual review of vision science*, vol. 1, pp. 373–391, 2015.
- [53] F. Engert and T. Bonhoeffer, "Dendritic spine changes associated with hippocampal long-term synaptic plasticity," *Nature*, vol. 399, no. 6731, pp. 66–70, 1999.
- [54] J. F. Mitchell, K. A. Sundberg, and J. H. Reynolds, "Spatial attention decorrelates intrinsic activity fluctuations in macaque area v4," *Neuron*, vol. 63, no. 6, pp. 879–888, 2009.
- [55] S. Kundu, M. Pedram, and P. A. Beerel, "Hire-snn: Harnessing the inherent robustness of energy-efficient deep spiking neural networks by training with crafted input noise," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 5209–5218.
- [56] B. Yin, F. Corradi, and S. M. Bohtë, "Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks," *Nature Machine Intelligence*, vol. 3, no. 10, pp. 905–913, 2021.
- [57] M. Horowitz, "1.1 computing's energy problem (and what we can do about it)," in *2014 IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*. IEEE, 2014, pp. 10–14.
- [58] S. B. Shrestha and G. Orchard, "Slayer: spike layer error reassignment in time," in *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, 2018, pp. 1419–1428.
- [59] Y. Wang, B. Du, Y. Shen, K. Wu, G. Zhao, J. Sun, and H. Wen, "Ev-gait: Event-based robust gait recognition using dynamic vision sensors," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6358–6367.
- [60] 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.
- [61] B. Han, G. Srinivasan, and K. Roy, "Rmp-snn: Residual membrane potential neuron for enabling deeper high-accuracy and low-latency spiking neural network," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [62] T. Bu, J. Ding, Z. Yu, and T. Huang, "Optimized potential initialization for low-latency spiking neural networks," *Proceedings of the AAAI Conference on Artificial Intelligence*, 2022.
- [63] S. Deng, Y. Li, S. Zhang, and S. Gu, "Temporal efficient training of spiking neural network via gradient re-weighting," in *International Conference on Learning Representations*, 2022.
- [64] Y. Lin, Y. Hu, S. Ma, G. Li, and D. Yu, "Rethinking pretraining as a bridge from anns to snns," *arXiv preprint arXiv:2203.01158*, 2022.
- [65] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "Imagenet: A large-scale hierarchical image database," in *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 2009, pp. 248–255.
- [66] S. S. Chowdhury, N. Rathi, and K. Roy, "One timestep is all you need: Training spiking neural networks with ultra low latency," *arXiv preprint arXiv:2110.05929*, 2021.
- [67] C. Xu, Y. Liu, and Y. Yang, "Direct training via backpropagation for ultra-low latency spiking neural networks with multi-threshold," *arXiv preprint arXiv:2112.07426*, 2021.
- [68] Y. Kim and P. Panda, "Visual explanations from spiking neural networks using inter-spike intervals," *Scientific reports*, vol. 11, no. 1, pp. 1–14, 2021.
- [69] F. Wang, M. Jiang, C. Qian, S. Yang, C. Li, H. Zhang, X. Wang, and X. Tang, "Residual attention network for image classification," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3156–3164.
- [70] L. Deng, Y. Wu, Y. Hu, L. Liang, G. Li, X. Hu, Y. Ding, P. Li, and Y. Xie, "Comprehensive snn compression using admm optimization and activity regularization," *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2021.
- [71] H.-H. Lien and T.-S. Chang, "Sparse compressed spiking neural network accelerator for object detection," *IEEE Transactions on Circuits and Systems I: Regular Papers*, 2022.
- [72] X. Glorot, A. Bordes, and Y. Bengio, "Deep sparse rectifier neural networks," in *Proceedings of the fourteenth international conference on artificial intelligence and statistics*. JMLR Workshop and Conference Proceedings, 2011, pp. 315–323.