SPIKE-DRIVEN TRANSFORMER V2: META SPIKING NEURAL NETWORK ARCHITECTURE INSPIRING THE DESIGN OF NEXT-GENERATION NEUROMORPHIC CHIPS

Man Yao¹, Jiakui Hu², Tianxiang Hu¹, Yifan Xu¹, Zhaokun Zhou^{2,3}, Yonghong Tian^{2,3}, Bo Xu¹, Guoqi Li^{1*}

¹Institute of Automation, Chinese Academy of Sciences, Beijing, China ²Peking University, Beijing, China ³Peng Cheng Laboratory, Shenzhen, Guangzhou, China

Abstract

Neuromorphic computing, which exploits Spiking Neural Networks (SNNs) on neuromorphic chips, is a promising energy-efficient alternative to traditional AI. CNN-based SNNs are the current mainstream of neuromorphic computing. By contrast, no neuromorphic chips are designed especially for Transformer-based SNNs, which have just emerged, and their performance is only on par with CNN-based SNNs, offering no distinct advantage. In this work, we propose a general Transformer-based SNN architecture, termed as "Meta-SpikeFormer", whose goals are: i) Lower-power, supports the spike-driven paradigm that there is only sparse addition in the network; ii) Versatility, handles various vision tasks; iii) *High-performance*, shows overwhelming performance advantages over CNN-based SNNs; iv) Meta-architecture, provides inspiration for future nextgeneration Transformer-based neuromorphic chip designs. Specifically, we extend the Spike-driven Transformer in Yao et al. (2023b) into a meta architecture, and explore the impact of structure, spike-driven self-attention, and skip connection on its performance. On ImageNet-1K, Meta-SpikeFormer achieves 80.0% top-1 accuracy (55M), surpassing the current state-of-the-art (SOTA) SNN baselines (66M) by 3.7%. This is the first direct training SNN backbone that can simultaneously supports classification, detection, and segmentation, obtaining SOTA results in SNNs. Finally, we discuss the inspiration of the meta SNN architecture for neuromorphic chip design. Source code and models are available at https://github.com/BICLab/Spike-Driven-Transformer-V2.

1 INTRODUCTION

The ambition of SNNs is to become a low-power alternative to traditional machine intelligence (Roy et al., 2019; Li et al., 2023). The unique *spike-driven* is key to realizing this magnificent concept, i.e., *only a portion of spiking neurons are ever activated to execute sparse synaptic ACcumulate (AC)* when SNNs are run on neuromorphic chips (Roy et al., 2019). Neuromorphic computing is essentially an algorithm-hardware co-design paradigm (Frenkel et al., 2023). Biological neurons are modeled as spiking neurons and somehow form SNNs at the algorithmic level(Maass, 1997a). Neuromorphic chips are then outfitted with spike-driven SNNs at the hardware level (Schuman et al., 2022).

CNN-based SNNs are currently the common spike-driven design. Thus, typical neuromorphic chips, such as TrueNorth (Merolla et al., 2014), Loihi (Davies et al., 2018), Tianjic (Pei et al., 2019), etc., all support spike-driven Conv and MLP operators. Nearly all CNN-era architectures, e.g., VGG (Simonyan & Zisserman, 2015), ResNet (He et al., 2016b), etc., can be developed into corresponding SNN versions (Wu et al., 2021). As Transformer (Vaswani et al., 2017) in ANNs has shown great potential in various tasks (Dosovitskiy et al., 2021), some Transformer-based designs have emerged in SNNs during the past two years (Zhang et al., 2022b;c; Han et al., 2023; Zhou et al., 2023).

^{*}Corresponding author, guoqi.li@ia.ac.cn

Most Transformer-based SNNs fail to take advantage of the low-power of SNNs because they are not spike-driven. Typically, they retain the energy-hungry Multiply-and-ACcumulate (MAC) operations dominated by vanilla Transformer, such as scaled dot-product (Han et al., 2023), softmax (Leroux et al., 2023), scale (Zhou et al., 2023), etc. Recently, Yao et al. (2023b) developed a spike-driven self-attention operator, integrating the spike-driven into Transformer for the first time. However, while the Spike-driven Transformer Yao et al. (2023b) with only sparse AC achieved SOTA results in SNNs on ImageNet-1K, it has yet to show clear advantages over Conv-based SNNs.

In this work, we advance the SNN domain by proposing Meta-SpikeFormer in terms of *performance* and *versatility*. Since Vision Transformer (ViT) (Dosovitskiy et al., 2021) showed that Transformer can perform superbly in vision, numerous studies have been produced (Han et al., 2022). Recently, Yu et al. (2022a;b) summarized various ViT variants and argued that there is general architecture abstracted from ViTs by not specifying the token mixer (self-attention). Inspired by this work, we investigate the meta architecture design in Transformer-based SNNs, involving three aspects: *network structure, skip connection (shortcut), Spike-Driven Self-Attention (SDSA)* with fully AC operations.

We first align the structures of Spike-driven Transformer in Yao et al. (2023b) with the CAFormer in Yu et al. (2022b) at the macro-level. Specifically, as shown in Fig. 2, the original four spiking encoding layers are expanded into four Conv-based SNN blocks. We experimentally verify that early-stage Conv blocks are important for the performance and versatility of SNNs. Then we design Conv-based and Transformer-based SNN blocks at the micro-level (see Table 5). For instance, the generation of spike-form Query (Q_S), Key (K_S), Value (V_S), three new SDSA operators, etc. Furthermore, we test the effects of three shortcuts based on the proposed Meta-SpikeFormer.

We conduct a comprehensive evaluation of Meta-SpikeFormer on four types of vision tasks, including image classification (ImageNet-1K (Deng et al., 2009)), event-based action recognition (HAR-DVS (Wang et al., 2022), currently the *largest* event-based human activity recognition dataset), object detection (COCO (Lin et al., 2014)), and semantic segmentation (ADE20K (Zhou et al., 2017), VOC2012 (Everingham et al., 2010)). The main contributions of this paper are as follows:

- SNN architecture design. We design a meta Transformer-based SNN architecture with only sparse addition, including macro-level Conv-based and Transformer-based SNN blocks, as well as some micro-level designs, such as several new spiking convolution methods, the generation of Q_S , K_S , V_S , and three new SDSA operator with different computational complexities, etc.
- **Performance.** The proposed Meta-SpikeFormer enables the performance of the SNN domain on ImageNet-1K to achieve 80% for the first time, which is 3.7% higher than the current SOTA baseline but with 17% fewer parameters (55M vs. 66M).
- Versatility. To the best of our knowledge, Meta-SpikeFormer is the first direct training SNN backbone that can handle image classification, object detection, semantic segmentation concurrently. We achieve SOTA results in the SNN domain on all tested datasets.
- **Neuromorphic chip design.** We thoroughly investigate the general components of Transformerbased SNN, including structure, shortcut, SDSA operator. And, Meta-SpikeFormer shows significant performance and versatility advantages over Conv-based SNNs. This will undoubtedly inspire and guide the neuromorphic computing field to develop Transformer-based neuromorphic chips.

2 RELATED WORK

Spiking Neural Networks can be simply considered as Artificial Neural Networks (ANNs) with bio-inspired spatio-temporal dynamics and spike (0/1) activations (Li et al., 2023). Spike-based communication enables SNNs to be spike-driven, but the conventional backpropagation algorithm (Rumelhart et al., 1986) cannot be applied directly because the spike function is non-differentiable. There are typically two ways to tackle this challenge. One is to discrete the trained ANNs into corresponding SNNs through neuron equivalence (Deng & Gu, 2021; Hu et al., 2023), i.e., ANN2SNN. Another is to train SNNs directly, using surrogate gradients (Wu et al., 2018; Neftci et al., 2019). In this work, we employ the direct training method due to its small timestep and adaptable architecture.

Backbone in Conv-based SNNs. The architecture of the Conv-based SNN is guided by residual learning in ResNet (He et al., 2016b;a). It can be roughly divided into three categories. Zheng et al. (2021) directly copied the shortcut in ResNet and proposed a tdBN method, which expanded SNN

from several layers to 50 layers. To solve the degradation problem of deeper Res-SNN, Fang et al. (2021) and Hu et al. (2024) proposed SEW-Res-SNN and MS-Res-SNN to raise the SNN depth to more than 100 layers. Then, the classic ANN architecture can have a corresponding SNN direct training version, e.g., attention SNNs (Yao et al., 2021; 2023d;a;c), spiking YOLO (Su et al., 2023), etc. Unfortunately, current CNN-based SNNs fail to demonstrate generality in vision tasks.

Vision Transformers. After ViT (Dosovitskiy et al., 2021) showed the promising performance, improvements and discussions on ViT have gradually replaced traditional CNNs as the mainstay. Some typical work includes architecture design (PVT (Wang et al., 2021a), MLP-Mixer (Tolstikhin et al., 2021)), enhancement of self-attention (Swin (Liu et al., 2021), Twins (Chu et al., 2021)), training optimization (DeiT (Touvron et al., 2021), T2T-ViT (Yuan et al., 2021)), efficient ViT (Katharopoulos et al., 2020; Xu et al., 2022), etc. In this work, we aim to reference and explore a meta spiking Transformer architecture from the cumbersome ViT variants available to bridge the gap between SNNs and ANNs, and pave the way for future Transformer-based neuromorphic chip design.

Neuromorphic Chips. Neuromorphic hardware is non-von Neumann architecture hardware whose structure and function are inspired by brains Roy et al. (2019). Typical neuromorphic chip features include collocated processing and memory, spike-driven computing, etc (Schuman et al., 2022). Functionally speaking, neuromorphic chips that are now on the market either support solely SNN or support hybrid ANN/SNN (Li et al., 2023). The former group consists of TrueNorth (Merolla et al., 2014), Loihi (Davies et al., 2018), Darwin (Shen et al., 2016), etc. The latter includes the Tianjic series (Pei et al., 2019; Ma et al., 2022), SpiNNaker 2 (Höppner et al., 2021). All of these chips enable CNN-based SNNs, but none of them are designed to support Transformer-based SNNs.

3 SPIKE-DRIVEN TRANSFORMER V2: META-SPIKEFORMER

3.1 THE CONCEPT OF META TRANSFORMER ARCHITECTURE IN ANNS

The self-attention (serves as a *token mixer*) mechanism for aggregating information between different spatial locations (tokens) has long been attributed to the success of Transformer. With the deepening of research, researchers have found that token mixer can be replaced by spatial Multi-Layer Perception (MLP) (Tolstikhin et al., 2021), Fourier Transform (Guibas et al., 2022), etc. Consequently, Yu et al. (2022a;b) argue that compared with a specific token mixer, a genera meta Transformer block (Fig. 1), is more essential than a specific token mixer for the model to achieve competitive performance.

Specifically, the input is first embedded as a feature sequence (tokens) (Vaswani et al., 2017; Dosovitskiy et al., 2021):

X =InputEmbedding(I),



Figure 1: Meta Transformer Block.

where $I \in \mathbb{R}^{3 \times H \times W}$ and $X \in \mathbb{R}^{N \times D}$. 3, *H*, and *W* denote channel, height and width of the 2D image. *N* and *D* represent token number and channel dimension respectively. Then the token sequence *X* is fed into repeated meta Transformer block, one of which can be expressed as (Fig. 1)

(1)

$$X' = X + \text{TokenMixer}(X), \tag{2}$$

$$X'' = X' + \text{ChannelMLP}(X'), \tag{3}$$

where $TokenMixer(\cdot)$ means token mixer mainly for propagating spatial information among tokens, ChannelMLP(\cdot) denotes a channel MLP network with two layers. $TokenMixer(\cdot)$ can be selfattention (Vaswani et al., 2017), spatial MLP (Touvron et al., 2021), convolution (Yu et al., 2022b), pooling (Yu et al., 2022a), linear attention (Katharopoulos et al., 2020), identity map (Wang et al., 2023b), etc, with different computational complexity, parameters and task accuracy.

3.2 SPIKING NEURON LAYER

Spiking neuron layer incorporates spatio-temporal information into membrane potentials, then converts them into binary spikes for spike-driven computing in the next layer. We adopt the standard



Figure 2: The overview of Meta-SpikeFormer. At the macro level, we refer to the general vision Transformer architecture in Yu et al. (2022a;b) and align Spike-driven Transformer (Yao et al., 2023b) with it. The main macro-level alteration is that we enlarge the spike coding layer from four Conv SNN layers to four Conv-based SNN blocks. At the micro level, we use the meta Transformer block in Fig. 1 as the basis to upgrade to Conv-based and Transformer-based SNN blocks (see Table 5), such as Channel Conv, SDSA operations, etc., to bring them more in line with SNN features.

Leaky Integrate-and-Fire (LIF) (Maass, 1997b) spiking neuron layer, whose dynamics are:

$$U[t] = H[t-1] + X[t],$$
(4)

$$S[t] = \text{Hea}\left(U[t] - u_{th}\right),\tag{5}$$

$$H[t] = V_{reset}S[t] + \left(\beta U[t]\right)\left(\mathbf{1} - S[t]\right),\tag{6}$$

where X[t] (X[t] can be obtained through spike-driven operators such as Conv, MLP, and selfattention) is the spatial input current at timestep t, U[t] means the membrane potential that integrates X[t] and temporal input H[t-1]. Hea(\cdot) is a Heaviside step function which equals 1 for $x \ge 0$ and 0 otherwise. When U[t] exceeds the firing threshold u_{th} , the spiking neuron will fire a spike S[t], and temporal output H[t] is reset to V_{reset} . Otherwise, U[t] will decay directly to H[t], where $\beta < 1$ is the decay factor. For simplicity, we mainly focus on Eq. 5 and re-write the spiking neuron layer as $SN(\cdot)$, with its input as membrane potential tensor U and its output as spike tensor S.

3.3 META-SPIKEFORMER

In SNNs, the input sequence $I \in \mathbb{R}^{T \times 3 \times H \times W}$, where T denote timestep. For example, images are repeated T times when dealing with a static dataset. To ease of understanding, we subsequently assume T = 1 when describing the architectural details of Meta-SpikeFormer.

Overall Architecture. Fig. 2 shows the overview of Meta-SpikeFormer, where Conv-based and Transformer-based SNN blocks are both variants of the meta Transformer block in Sec 3.1. In Spike-driven Transformer (Yao et al., 2023b), the authors exploited four Conv layers before Transformer-based blocks for encoding. By contrast, in the architectural form of Conv+ViT in ANNs, there are generally multiple stages of Conv blocks (Han et al., 2022; Xiao et al., 2021). We follow this typical design in ANNs, setting the first two stages to Conv-based SNN blocks, and using a pyramid structure (Wang et al., 2021a) in the last two Transformer-based SNN stages. Note, to control parameter number, we set the channels to 10C in stage 4 instead of the typical double (16C). Fig. 2 is our recommended architecture. Other alternatives and their impacts are summarized in Table 5.

Conv-based SNN Block uses the inverted separable convolution module $SepConv(\cdot)$ with 7×7 kernel size in MobileNet V2 (Sandler et al., 2018) as $TokenMixer(\cdot)$, which is consistent with (Yu et al., 2022b). But, we change $ChannelMLP(\cdot)$ with 1×1 kernel size in Eq. 3 to $ChannelConv(\cdot)$ with



Figure 3: Spike-Driven Self-Attention (SDSA) modules with different computational complexity. SDSA-1 is the operator used in Yao et al. (2023b). SDSA-2/3/4 is the newly designed operator in this paper. We exploit SDSA-3 by default. All SDSAs only have addition, no softmax and scale.

 3×3 kernel size. The stronger inductive is empirically proved to significantly improve performance (see Table 5). Specifically, the Conv-based SNN block is written as:

$$U' = U + \operatorname{SepConv}(U), \tag{7}$$

$$U'' = U' + \text{ChannelConv}(U'), \tag{8}$$

 $SepConv(U) = Conv_{pw2}(Conv_{dw}(\mathcal{SN}(Conv_{pw1}(\mathcal{SN}(U))))),$ (9)

$$ChannelConv(U') = Conv(\mathcal{SN}(Conv(\mathcal{SN}(U')))).$$
(10)

where $\operatorname{Conv}_{pw1}(\cdot)$ and $\operatorname{Conv}_{pw2}(\cdot)$ are pointwise convolutions, $\operatorname{Conv}_{dw}(\cdot)$ is depthwise convolution (Chollet, 2017), $\operatorname{Conv}(\cdot)$ is the normal convolution. $\mathcal{SN}(\cdot)$ is the spike neuron layer in Sec 3.2.

Transformer-based SNN Block contains an SDSA module and a two-layered ChannelMLP(·):

$$Q_S = \mathcal{SN}(\operatorname{RepConv}_1(U)), K_S = \mathcal{SN}(\operatorname{RepConv}_2(U)), V_S = \mathcal{SN}(\operatorname{RepConv}_3(U)), \quad (11)$$

$$U' = U + \operatorname{RepConv}_4(\operatorname{SDSA}(Q_S, K_S, V_S)), \tag{12}$$

$$U'' = U' + \text{ChannelMLP}(U'), \tag{13}$$

$$ChannelMLP(U') = SN(SN(U')W_1)W_2,$$
(14)

where $\operatorname{RepConv}(\cdot)$ is the re-parameterization convolutions (Ding et al., 2021) with kernel size 3×3 , $W_1 \in \mathbb{R}^{C \times rC}$ and $W_2 \in \mathbb{R}^{rC \times C}$ are learnable parameters with MLP expansion ratio r = 4. Note, both the input (Q_S, K_S, V_S) and output of $SDSA(\cdot)$ will be reshaped. We omit this for simplicity.

Spike-Driven Self-Attention (SDSA). The main difference of SDSA over vanilla self-attention with $O(N^2D)$ in Dosovitskiy et al. (2021) lies in *three* folds: i) Query, Key, Value are spiking tensors; ii) The operations among Q_S , K_S , V_S do not have softmax and scale; iii) The computational complexity of SDSA(·) is linear with the token number N. Four SDSA operators are given in Fig. 3. SDSA-1 is proposed in Yao et al. (2023b). SDSA-2/3/4 are new operators designed in this work. The main difference between them is the operation between Q_S , K_S , V_S . SDSA-1/2 primarily work with Hadamard product while SDSA-3/4 use matrix multiplication. Spike-driven matrix multiplication can be converted to additions via addressing algorithms (Frenkel et al., 2023). SDSA-1/2/3/4 all only have sparse addition. Details of SDSAs and energy evaluation are given in Appendix A and B.

In this work, we use SDSA-3 by default, which is written as:

$$SDSA_3(Q_S, K_S, V_S) = S\mathcal{N}_s\left(Q_S\left(K_S^{\mathrm{T}}V_S\right)\right) = S\mathcal{N}_s((Q_SK_S^{\mathrm{T}})V_S).$$
(15)

where $SN_s(\cdot)$ is $SN(\cdot)$ with the threshold $s \cdot u_{th}$. Note, $SDSA_3(\cdot)$ is inspired by the spiking selfattention $SN(Q_SK_S^TV_S * s)$ in Zhou et al. (2023). Because $Q_SK_S^TV_S$ yield large integers, a scale factor *s* for normalization is needed to avoid gradient vanishing in Zhou et al. (2023). In our SDSA-3, we directly merge the *s* into the threshold of the spiking neuron to circumvent the multiplication by *s*. Further, in SDSA-4, we set the threshold as a learnable parameter.

Methods	Architecture	Spike	Param (M)	Power (mJ)	Step	Acc.(%)
	ResNet-34 (Rathi et al., 2020)	1	21.8	-	250	61.5
ANN2SNN	VGG-16 (Wu et al., 2021)	1	-	-	16	65.1
	VGG-16 (Hu et al., 2023)	✓	138.4	44.9	7	73.0
	SEW-Res-SNN	X	25.6	4.9	4	67.8
	(Fang et al., 2021)	X	60.2	12.9	4	69.2
CNN-based	MS-Res-SNN	1	21.8	5.1	4	69.4
SNN	(Hu et al., 2024)	1	77.3	10.2	4	75.3
	Att-MS-Res-SNN	X	22.1	0.6	1	69.2
	(Yao et al., 2023d)	X	78.4	7.3	4	76.3
A NINI	RSB-CNN-152 (Wightman et al., 2021)	X	60	53.4	1	81.8
AININ	ViT (Dosovitskiy et al., 2021)	X	86	81.0	1	79.7
	SpikFormer	X	29.7	11.6	4	73.4
	(Zhou et al., 2023)	X	66.3	21.5	4	74.8
	Spike-driven Transformer	\checkmark	29.7	4.5	4	74.6
	(Yao et al., 2023b)	1	66.3	6.1	4	76.3
Transformer		1	15.1	4.0	1	71.8
-based SNN		1	15.1	16.7	4	74.1
	Meta-SpikeFormer	~	31.3	7.8	1	75.4
	(This Work)	1	31.3	32.8	4	77.2
		-	55.4	13.0	1	79.1*
		./	55 4	524	4	80.0*

Table 1: Performance on ImageNet-1K (Deng et al., 2009). The input crop is 224×224 . *We obtain these results by employing distillation training on method in DeiT (Touvron et al., 2021). When trained directly, the accuracy are 78.0% (T = 1) and 79.7% (T = 4). Note, "Spike", "Para", and "Step" in all Table headers of this paper denote "Spike-driven", "Parameters", and "Timestep".

Shortcuts. Residual learning in SNNs mainly considers two points: first, whether identity mapping (He et al., 2016a) can be realized, which determines whether there is a degradation problem; second, whether spike-driven computing can be guaranteed, which is the basis of SNNs' lowpower. There are three shortcuts in SNNs, see Fig. 4. Vanilla Shortcut (VS) (Zheng et al., 2021) execute a shortcut between membrane potential and spike that are consistent with those in Res-CNN (He et al., 2016b). It can be spike-driven, but cannot achieve identity mapping (Fang et al., 2021). Spike-Element-Wise (SEW) (Fang et al., 2021) exploits a shortcut to connect spikes in different layers. Identity mapping is possible with SEW, but spike addition results in integers.



Figure 4: Existing shortcut in SNNs.

Thus, SEW-SNN is an "integer-driven" rather than a spike-driven SNN. Membrane Shortcut (MS) makes a shortcut between the membrane potential of neurons, and can achieve identity mapping with spike-driven (Hu et al., 2024). We use MS in this work and report the accuracy of other shortcuts.

4 **EXPERIMENTS**

In the classification task, we set the timestep T = 1 for 200 epoch training to reduce training cost, then finetune it to T = 4 with an additional 20 epochs. We here mainly emphasize the network scale. Other details, such as training schemes and configurations, are in Appendix C.1. Moreover, we use the model trained on ImageNet classification to finetune the detection or segmentation heads. *This is the first time that the SNN domain has been able to process dense prediction tasks in a unified way.*

Methods	Architecture	Spike	Param(M)	Acc.(%)
ANN	SlowFast (Feichtenhofer et al., 2019) ACTION-Net (Wang et al., 2021b)	X X	33.6 27.9	46.5 46.9
	TimeSformer (Bertasius et al., 2021)	×	121.2	50.8
CNN-based SNN	Res-SNN-34 (Fang et al., 2021)	×	21.8	46.1
Transformer-based SNN	Meta-SpikeFormer (This Work)	1	18.3	47.5

Table 2. Performance of event-based action recognition on \mathbf{HAK} -DVS (wang et al., 2022	n recognition on HAR-DVS (Wang et al., 2022).	able 2: Performance of event-based action
--	---	---

4.1 IMAGE CLASSIFICATION

Setup. ImageNet-1K (Deng et al., 2009) is one of computer vision's most widely used datasets. It contains about 1.3M training and 50K validation images, covering common 1K classes. As shown in Fig. 2, changing the channel number can obtain various model scales. We set C = 32, 48, 64. Their parameters are 15.1M, 31.3M, and 55.4M, respectively. The energy cost of ANNs is FLOPs times E_{MAC} . The energy cost of SNNs is FLOPs times E_{AC} times spiking firing rate. $E_{MAC} = 4.6pJ$ and $E_{AC} = 0.9pJ$ are the energy of a MAC and an AC, respectively (more details in Appendix B).

Results. We comprehensively compare Meta-SpikeFormer with other methods in accuracy, parameter, and power (Table 1). We can see that Meta-SpikeFormer achieves SOTA in the SNN domain with significant accuracy advantages. For example, **Meta-SpikeFormer** vs. MS-Res-SNN vs. Spike-driven Transformer: Param, **55M** vs. 77M vs. 66M; Acc, **79.7%** vs. 75.3% vs. 76.3%. If we employ the distillation strategy in DeiT (Touvron et al., 2021), the accuracy of 55M Meta-SpikeFormer at T = 1 and T = 4 can be improved to 79.1% and 80.0%. It should be noted that after adding more Conv layers at stage 1/2, the power of Meta-SpikeFormer increases. This is still very cost-effective. For instance, **Meta-SpikeFormer** vs. MS-Res-SNN vs. Spike-driven Transformer: Power, **11.9mJ** (T = 1) vs. 6.1mJ (T = 4) vs. 10.2mJ (T = 4); Acc, **79.1%** vs. 75.3% vs. 76.3%. In summary, Meta-SpikeFormer obtained the first achievement of 80% accuracy on ImageNet-1K in SNNs.

4.2 EVENT-BASED ACTIVITY RECOGNITION

Event-based vision (also known as "neuromorphic vision") is one of the most typical application scenarios of neuromorphic computing (Indiveri & Douglas, 2000; Gallego et al., 2022; Wu et al., 2022). The famous neuromorphic camera, Dynamic Vision Sensors (DVS) (Lichtsteiner et al., 2008), encodes vision information into a sparse event (spike with address) stream only when brightness changes. Since the spike-driven nature, SNNs have the inherent advantages of low power and low latency when processing events. We use HAR-DVS to evaluate our method. HAR-DVS (Wang et al., 2022) is the largest event-based Human Activity Recognition (HAR) dataset currently, containing 300 classes and 107,646 samples, acquired by a DAVIS346 camera with a spatial resolution of 346×260 . The raw HAR-DVS is more than 4TB, and the authors convert each sample into frames to form a new HAR-DVS. We handle HAR-DVS in a direct training manner with T = 8. Meta-SpikeFormer achieves comparable accuracy to ANNs and is better than the Conv-based SNN baseline (Table 2).

4.3 **OBJECT DETECTION**

So far, no backbone with direct training in SNNs can handle classification, detection, and segmentation tasks concurrently. Only recently did the SNN domain have the first directly trained model to process detection (Su et al., 2023). We evaluate Meta-SpikeFormer on the COCO benchmark (Lin et al., 2014) that includes 118K training images (train2017) and 5K validation images (val2017). We first transform *mmdetection* (Chen et al., 2019) codebase into a spiking version and use it to execute our model. We exploit Meta-SpikeFormer with Mask R-CNN (He et al., 2017). ImageNet pre-trained weights are utilized to initialize the backbones, and Xavier (Glorot & Bengio, 2010) to initialize the added layers. Results are shown in Table 3. We can see that Meta-SpikeFormer achieves SOTA results in the SNN domain. Note, EMS-Res-SNN got performance close to ours using only 14.6M parameters, thanks to its direct training strategy and special network design. In contrast, we only use a fine-tuning strategy, which results in lower design and training costs. To be fair, we also tested directly trained Meta-SpikeFormer + Yolo and achieved good performance (Appendix C.2).

Methods	Architecture	Spike	Param(M)	Power(mJ)	Step	mAP@0.5(%)
ANN	ResNet-18 (Yu et al., 2022a)	×	31.2	890.6	1	54.0
	PVT-Tiny (Wang et al., 2021a)	×	32.9	1040.5	1	56.7
ANN2SNN	Spiking-Yolo (Kim et al., 2020) Spike Calibration (Li et al., 2022)	√ √	10.2 17.1	-	3500 512	25.7 45.3
CNN-based	Spiking Retina (Zhang et al., 2023)	×	11.3	-	4	28.5
SNN	EMS-Res-SNN (Su et al., 2023)	✓	14.6		4	50.1
Transformer	Meta-SpikeFormer (This Work)	\	34.9	49.5	1	44.0
-based SNN		\	75.0	140.8	1	51.2

Table 3: Performance of object detection on COCO val2017 (Lin et al., 2014).

Table 4: Performance of semantic segmentation on ADE20K (Zhou et al., 2017).

Methods	Architecture	Spike	Param(M)	Power(mJ)	Step	MIoU(%)
	ResNet-18 (Yu et al., 2022a)	X	15.5	147.1	1	32.9
ANTNI	PVT-Tiny (Wang et al., 2021a)	X	17.0	152.7	1	35.7
AININ	PVT-Small (Wang et al., 2021a)	X	28.2	204.7	1	39.8
	DeepLab-V3 (Zhang et al., 2022a)	X	68.1	1240.6	1	42.7
		1	16.5	22.1	1	32.3
Transformer	Meta-SpikeFormer	1	16.5	88.1	4	33.6
-based SNN	(This Work)	-	58.9	46.6	1	34.8
		1	58.9	183.6	4	35.3

4.4 SEMANTIC SEGMENTATION

ADE20K (Zhou et al., 2017) is a challenging semantic segmentation benchmark commonly used in ANNs, including 20K and 2K images in the training and validation set, respectively, and covering 150 categories. No SNN has yet reported processing results on ADE20K. In this work, Meta-SpikeFormers are evaluated as backbones equipped with Semantic FPN (Kirillov et al., 2019). ImageNet trained checkpoints are used to initialize the backbones while Xavier (Glorot & Bengio, 2010) is utilized to initialize other newly added layers. We transform *mmsegmentation* (Contributors, 2020) codebase into a spiking version and use it to execute our model. Training details are given in Appendix C.3. We see that in lightweight models (16.5M in Table 4), Meta-SpikeFormer with lower power achieves comparable results to ANNs. For example, **Meta-SpikeFormer** (T = 1) vs. ResNet-18 vs. PVT-Tiny: Param, **16.5M** vs. 15.5M vs. 17.0M; MIOU, **32.3%** vs. 32.9% vs. 35.7%; Power, **22.1mJ** vs. 147.1mJ vs. 152.7mJ. To demonstrate the superiority of our method over other SNN segmentation methods, we also evaluate our method on VOC2012 and achieve SOTA results (Appendix C.4).

4.5 Ablation Studies

Conv-based SNN Block. In this block, We follow the ConvFormer in (Yu et al., 2022b), which uses SpeConv as token mixer in stage-1/2. However, we note that SpeConv in Meta-SpikeFormer seems less important. After removing SpeConv, the power is reduced by 29.5%, the accuracy is only lost by 0.3%. If we replace the channel Conv with the channel MLP in (Yu et al., 2022b), the accuracy will drop by up to 2%. Thus, the design of Conv-based SNN Blocks is important to SNNs' performance. Moreover, we experimentally verified (specific results are omitted) that keeping only one stage of the Conv-based block or using only four Conv layers leads to lower performance on downstream tasks.

Transformer-based SNN Block. Spike-driven Transformer in (Yao et al., 2023b) uses linear layers (i.e., 1×1 convolution) to generate Q_S, K_S, V_S . We find that replacing linear with RepConv can improve accuracy and reduce the parameter number, but energy costs will increase. The design of the SDSA operator and pyramid structure will also affect task accuracy. Overall, SDSA-3 has the highest computational complexity (more details in Appendix A), and its accuracy is also the best.

Shortcut. In our architecture, MS has the highest accuracy. Shortcut has almost no impact on power.

Ablation	Methods	Param(M)	Power(mJ)	Acc.(%)
	Meta-SpikeFormer (Baseline)	31.3	7.8	75.4
	Remove SepConv	30.9	5.5	75.1 (-0.3)
Conv-based SNN Block	Channel Conv -> Channel MLP	25.9	6.3	73.4 (-2.0)
	Stage 1 -> $2C \times \frac{H}{4} \times \frac{W}{4}$	31.8	7.4	75.2 (-0.2)
	RepConv-1/2/3 -> Linear	27.2	6.0	75.0 (-0.4)
Transformer based	RepConv-4 -> Linear	30.0	7.1	75.3 (-0.1)
SNIN block	SDSA-3 -> SDSA-1	31.3	7.2	74.6 (-0.8)
SININ DIOCK	SDSA-3 -> SDSA-2	28.6	6.3	74.2 (-1.2)
	SDSA-3 -> SDSA-4	31.3	7.7	75.4 (+0.0)
Shortout	Membrane Shortcut -> Vanilla shortcut	31.3	-	*
Shortcut	Membrane Shortcut -> SEW shortcut	31.3	7.8	73.5 (-1.9)
	Remove Pyramid (Stage4 = Stage 3)	26.9	7.4	74.7 (-0.7)
Architecture	Fully CNN-based SNN blocks	36.0	2.9	72.5 (-2.9)
	Fully Transformer-based SNN blocks	26.2	5.4	71.7 (-3.7)

Table 5: Ablation studies of Meta-SpikeFormer on ImageNet-1K. In each ablation experiment, we start with Meta-SpikeFormer (C = 48) with T = 1 as the baseline and modify just one point to track how the parameters, power, accuracy vary. * Does not converge.

Architecture. We change the network to fully Conv-based or Transformer-based blocks. Performance is significantly reduced in both cases. We note that compared to Meta-SpikeFormer, the power of fully spiking Transformer and spiking CNN are reduced. These observations can inspire future architectural designs to achieve multiple trade-offs in terms of parameter, power, and accuracy.

5 DISCUSSION AND CONCLUSION

Discussion: How does Meta-SpikeFormer inspire future neuromorphic chip design? The *technical* inspiration of Meta-SpikeFormer for chip design lies in *three* folds: i) *Conv+ViT design*. This hybrid progressive local-global modeling can leverage the strengths of both CNNs and Transformers (Guo et al., 2022), where the former models features and the latter captures long-range dependencies. We experimentally verify that this design is beneficial to the performance and versatility of SNNs. ii) *SDSA operator* is the core design of long-distance dependency modeling in Transformer-based SNN block, but this is a design that current neuromorphic chips lack. iii) *Meta architecture*. Given meta Conv-based and Transformer-based blocks, researchers can perform targeted optimization of the design details inside the meta SNN blocks according to their requirements in terms of accuracy, parameters, and power. As shown by our ablation experiments in Table 5.

The *significance* of Meta-SpikeFormer to chip design lies in *three* folds: i) *Algorithm-hardware co-design*. Most neuromorphic chip design begins from the bottom of the compute stack, i.e., the materials and devices (Schuman et al., 2022). The excellent features shown by our algorithm may attract and inspire algorithm-driven chip design. ii) *Confidence in large-scale neuromorphic computing*. Small-scale neuromorphic computing has shown significant power and performance advantages (Yin et al., 2021; Rao et al., 2022). We demonstrate the potential of larger-scale SNNs in performance and versatility. iii) *Reduce chip design costs*. Meta design facilitates follow-up optimization by subsequent researchers, helps the SNN field to quickly narrow the gap with ANNs, and reduces the cost of algorithm exploration required before algorithm-driven hardware design.

Conclusion. This paper investigates the meta design of Transformer-based SNNs, involving architecture, spike-driven self-attention, shortcut, etc. The proposed Meta-SpikeFormer is the first direct training SNN backbone that can perform classification, detection, and segmentation tasks concurrently, and we achieve state-of-the-art results on all tested datasets. Remarkably, for the first time, we advanced the accuracy of the SNN domain on ImageNet-1K to 80%, which is 3.7% higher than the prior SOTA result with 17% fewer parameters. This work paves the way for SNN to serve as a universal vision backbone and can inspire future Transformer-based neuromorphic chip designs.

ACKNOWLEDGEMENT

This work was partially supported by National Science Foundation for Distinguished Young Scholars (62325603), National Natural Science Foundation of China (62236009, U22A20103), Beijing Natural Science Foundation for Distinguished Young Scholars (JQ21015), and CAAI-MindSpore Open Fund, developed on OpenI Community.

REFERENCES

- Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, number 3, pp. 4, 2021.
- Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, and Judy Hoffman. Hydra attention: Efficient attention with many heads. In *Computer Vision–ECCV 2022 Workshops: Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VII*, pp. 35–49. Springer, 2023.
- Kai Chen, Jiaqi Wang, Jiangmiao Pang, Yuhang Cao, Yu Xiong, Xiaoxiao Li, Shuyang Sun, Wansen Feng, Ziwei Liu, Jiarui Xu, et al. Mmdetection: Open mmlab detection toolbox and benchmark. *arXiv preprint arXiv:1906.07155*, 2019.
- Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan L Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 40(4):834–848, 2017.
- François Chollet. Xception: Deep learning with depthwise separable convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1251–1258, 2017.
- Xiangxiang Chu, Zhi Tian, Yuqing Wang, Bo Zhang, Haibing Ren, Xiaolin Wei, Huaxia Xia, and Chunhua Shen. Twins: Revisiting the design of spatial attention in vision transformers. *Advances in Neural Information Processing Systems*, 34:9355–9366, 2021.
- MMSegmentation Contributors. Mmsegmentation: Openmmlab semantic segmentation toolbox and benchmark, 2020. URL https://github.com/open-mmlab/mmsegmentation.
- Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A neuromorphic manycore processor with on-chip learning. *IEEE Micro*, 38(1):82–99, 2018.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pp. 248–255. IEEE, 2009.
- Shikuang Deng and Shi Gu. Optimal conversion of conventional artificial neural networks to spiking neural networks. In *International Conference on Learning Representations*, 2021. URL https://openreview.net/forum?id=FZ10TwcXchK.
- Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. Repvgg: Making vgg-style convnets great again. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13733–13742, 2021.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2021.
- Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International Journal of Computer Vision*, 88: 303–338, 2010.
- Wei Fang, Zhaofei Yu, Yanqi Chen, Tiejun Huang, Timothée Masquelier, and Yonghong Tian. Deep residual learning in spiking neural networks. Advances in Neural Information Processing Systems, 34:21056–21069, 2021.

- Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6202–6211, 2019.
- Charlotte Frenkel, David Bol, and Giacomo Indiveri. Bottom-up and top-down approaches for the design of neuromorphic processing systems: Tradeoffs and synergies between natural and artificial intelligence. *Proceedings of the IEEE*, 2023.
- Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi, Stefan Leutenegger, Andrew J. Davison, Jörg Conradt, Kostas Daniilidis, and Davide Scaramuzza. Event-based vision: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(1):154–180, 2022. doi: 10.1109/TPAMI.2020.3008413.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pp. 249–256. JMLR Workshop and Conference Proceedings, 2010.
- John Guibas, Morteza Mardani, Zongyi Li, Andrew Tao, Anima Anandkumar, and Bryan Catanzaro. Efficient token mixing for transformers via adaptive fourier neural operators. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum? id=EXHG-A3jlM.
- Jianyuan Guo, Kai Han, Han Wu, Yehui Tang, Xinghao Chen, Yunhe Wang, and Chang Xu. Cmt: Convolutional neural networks meet vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12175–12185, 2022.
- Kai Han, Yunhe Wang, Hanting Chen, Xinghao Chen, Jianyuan Guo, Zhenhua Liu, Yehui Tang, An Xiao, Chunjing Xu, Yixing Xu, et al. A survey on vision transformer. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(1):87–110, 2022.
- Minglun Han, Qingyu Wang, Tielin Zhang, Yi Wang, Duzhen Zhang, and Bo Xu. Complex dynamic neurons improved spiking transformer network for efficient automatic speech recognition. *Thirty-Seventh AAAI Conference on Artificial Intelligence (AAAI 2023)*, 2023.
- Ali Hatamizadeh, Greg Heinrich, Hongxu Yin, Andrew Tao, Jose M Alvarez, Jan Kautz, and Pavlo Molchanov. Fastervit: Fast vision transformers with hierarchical attention. *arXiv preprint arXiv:2306.06189*, 2023.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (eds.), *Computer Vision ECCV 2016*, pp. 630–645, Cham, 2016a. Springer International Publishing.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016b.
- Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2961–2969, 2017.
- Sebastian Höppner, Yexin Yan, Andreas Dixius, Stefan Scholze, Johannes Partzsch, Marco Stolba, Florian Kelber, Bernhard Vogginger, Felix Neumärker, Georg Ellguth, et al. The spinnaker 2 processing element architecture for hybrid digital neuromorphic computing. *arXiv preprint arXiv:2103.08392*, 2021.
- Mark 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), pp. 10–14. IEEE, 2014.
- Yangfan Hu, Qian Zheng, Xudong Jiang, and Gang Pan. Fast-snn: Fast spiking neural network by converting quantized ann. *arXiv preprint arXiv:2305.19868*, 2023.
- Yifan Hu, Lei Deng, Yujie Wu, Man Yao, and Guoqi Li. Advancing spiking neural networks toward deep residual learning. *IEEE Transactions on Neural Networks and Learning Systems*, pp. 1–15, 2024.

Giacomo Indiveri and Rodney Douglas. Neuromorphic vision sensors. *Science*, 288(5469):1189–1190, 2000.

- A. Katharopoulos, A. Vyas, N. Pappas, and F. Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2020. URL https://fleuret.org/papers/katharopoulos-et-al-icml2020.pdf.
- Seijoon Kim, Seongsik Park, Byunggook Na, and Sungroh Yoon. Spiking-yolo: Spiking neural network for energy-efficient object detection. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34:11270–11277, Apr. 2020.
- Youngeun Kim, Joshua Chough, and Priyadarshini Panda. Beyond classification: Directly training spiking neural networks for semantic segmentation. *Neuromorphic Computing and Engineering*, 2 (4):044015, 2022.
- Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6399–6408, 2019.
- Nathan Leroux, Jan Finkbeiner, and Emre Neftci. Online transformers with spiking neurons for fast prosthetic hand control. *arXiv preprint arXiv:2303.11860*, 2023.
- Guoqi Li, Lei Deng, Huajing Tang, Gang Pan, Yonghong Tian, Kaushik Roy, and Wolfgang Maass. Brain inspired computing: A systematic survey and future trends. 2023.
- Yang Li, Xiang He, Yiting Dong, Qingqun Kong, and Yi Zeng. Spike calibration: Fast and accurate conversion of spiking neural network for object detection and segmentation. *arXiv preprint arXiv:2207.02702*, 2022.
- Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A 128×128 120 db 15 μ s latency asynchronous temporal contrast vision sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, 2008.
- Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer Vision– ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2117–2125, 2017.
- Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022, 2021.
- Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440, 2015.
- Songchen Ma, Jing Pei, Weihao Zhang, Guanrui Wang, Dahu Feng, Fangwen Yu, Chenhang Song, Huanyu Qu, Cheng Ma, Mingsheng Lu, et al. Neuromorphic computing chip with spatiotemporal elasticity for multi-intelligent-tasking robots. *Science Robotics*, 7(67):eabk2948, 2022.
- Wolfgang Maass. Networks of spiking neurons: The third generation of neural network models. *Neural Networks*, 10(9):1659–1671, 1997a.
- Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models. *Neural Networks*, 10(9):1659–1671, 1997b.
- Paul A Merolla, John V Arthur, Rodrigo Alvarez-Icaza, Andrew S Cassidy, Jun Sawada, Filipp Akopyan, Bryan L Jackson, Nabil Imam, Chen Guo, Yutaka Nakamura, et al. A million spikingneuron integrated circuit with a scalable communication network and interface. *Science*, 345 (6197):668–673, 2014.

- Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning convolutional neural networks for resource efficient inference. In *International Conference on Learning Representations*, 2017. URL https://openreview.net/forum?id=SJGCiw5ql.
- Emre O Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks. *IEEE Signal Processing Magazine*, 36(6):51–63, 2019.
- Priyadarshini Panda, Sai Aparna Aketi, and Kaushik Roy. Toward scalable, efficient, and accurate deep spiking neural networks with backward residual connections, stochastic softmax, and hybridization. *Frontiers in Neuroscience*, 14:653, 2020.
- Jing Pei, Lei Deng, et al. Towards artificial general intelligence with hybrid tianjic chip architecture. *Nature*, 572(7767):106–111, 2019.
- Arjun Rao, Philipp Plank, Andreas Wild, and Wolfgang Maass. A long short-term memory for ai applications in spike-based neuromorphic hardware. *Nature Machine Intelligence*, 4(5):467–479, 2022.
- Nitin Rathi, Gopalakrishnan Srinivasan, Priyadarshini Panda, and Kaushik Roy. Enabling deep spiking neural networks with hybrid conversion and spike timing dependent backpropagation. In *International Conference on Learning Representations*, 2020. URL https://openreview.net/forum?id=BlxSperKvH.
- Kaushik Roy, Akhilesh Jaiswal, and Priyadarshini Panda. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575(7784):607–617, 2019.
- D Rumelhart, G Hinton, and R Williams. Learning representations by back-propagating errors. *Nature*, 323:533–536, 1986.
- Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4510–4520, 2018.
- Catherine D Schuman, Shruti R Kulkarni, Maryam Parsa, J Parker Mitchell, Bill Kay, et al. Opportunities for neuromorphic computing algorithms and applications. *Nature Computational Science*, 2 (1):10–19, 2022.
- Juncheng Shen, De Ma, Zonghua Gu, Ming Zhang, Xiaolei Zhu, Xiaoqiang Xu, Qi Xu, Yangjing Shen, and Gang Pan. Darwin: A neuromorphic hardware co-processor based on spiking neural networks. *Science China Information Sciences*, 59(2):1–5, 2016.
- K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, pp. 1–14, San Diego, CA, United states, 2015.
- Qiaoyi Su, Yuhong Chou, Yifan Hu, Jianing Li, Shijie Mei, Ziyang Zhang, and Guoqi Li. Deep directly-trained spiking neural networks for object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 6555–6565, 2023.
- Ilya O Tolstikhin, Neil Houlsby, Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Thomas Unterthiner, Jessica Yung, Andreas Steiner, Daniel Keysers, Jakob Uszkoreit, et al. Mlp-mixer: An all-mlp architecture for vision. *Advances in Neural Information Processing Systems*, 34: 24261–24272, 2021.
- Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International conference on machine learning*, pp. 10347–10357. PMLR, 2021.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, pp. 5998–6008, 2017.

- Dingheng Wang, Bijiao Wu, Guangshe Zhao, Man Yao, Hengnu Chen, Lei Deng, Tianyi Yan, and Guoqi Li. Kronecker cp decomposition with fast multiplication for compressing rnns. *IEEE Transactions on Neural Networks and Learning Systems*, 34(5):2205–2219, 2023a.
- Jiahao Wang, Songyang Zhang, Yong Liu, Taiqiang Wu, Yujiu Yang, Xihui Liu, Kai Chen, Ping Luo, and Dahua Lin. Riformer: Keep your vision backbone effective but removing token mixer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 14443–14452, 2023b.
- Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 568–578, 2021a.
- Xiao Wang, Zongzhen Wu, Bo Jiang, Zhimin Bao, Lin Zhu, Guoqi Li, Yaowei Wang, and Yonghong Tian. Hardvs: Revisiting human activity recognition with dynamic vision sensors. *arXiv preprint arXiv:2211.09648*, 2022.
- Zhengwei Wang, Qi She, and Aljosa Smolic. Action-net: Multipath excitation for action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13214–13223, 2021b.
- Ross Wightman, Hugo Touvron, and Hervé Jégou. Resnet strikes back: An improved training procedure in timm. *arXiv preprint arXiv:2110.00476*, 2021.
- Jibin Wu, Chenglin Xu, Xiao Han, Daquan Zhou, Malu Zhang, Haizhou Li, and Kay Chen Tan. Progressive tandem learning for pattern recognition with deep spiking neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(11):7824–7840, 2021.
- Yang Wu, Ding-Heng Wang, Xiao-Tong Lu, Fan Yang, Man Yao, Wei-Sheng Dong, Jian-Bo Shi, and Guo-Qi Li. Efficient visual recognition: A survey on recent advances and brain-inspired methodologies. *Machine Intelligence Research*, 19(5):366–411, 2022.
- Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in Neuroscience*, 12:331, 2018.
- Tete Xiao, Mannat Singh, Eric Mintun, Trevor Darrell, Piotr Dollár, and Ross Girshick. Early convolutions help transformers see better. *Advances in Neural Information Processing Systems*, 34: 30392–30400, 2021.
- Yifan Xu, Zhijie Zhang, Mengdan Zhang, Kekai Sheng, Ke Li, Weiming Dong, Liqing Zhang, Changsheng Xu, and Xing Sun. Evo-vit: Slow-fast token evolution for dynamic vision transformer. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 36, pp. 2964–2972, 2022.
- Helin Yang, Kwok-Yan Lam, Liang Xiao, Zehui Xiong, Hao Hu, Dusit Niyato, and H Vincent Poor. Lead federated neuromorphic learning for wireless edge artificial intelligence. *Nature Communications*, 13(1):4269, 2022.
- Man Yao, Huanhuan Gao, Guangshe Zhao, Dingheng Wang, Yihan Lin, Zhaoxu Yang, and Guoqi Li. Temporal-wise attention spiking neural networks for event streams classification. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pp. 10221–10230, 2021.
- Man Yao, Jiakui Hu, Guangshe Zhao, Yaoyuan Wang, Ziyang Zhang, Bo Xu, and Guoqi Li. Inherent redundancy in spiking neural networks. In *Proceedings of the IEEE/CVF International Conference* on Computer Vision, pp. 16924–16934, 2023a.
- Man Yao, JiaKui Hu, Zhaokun Zhou, Li Yuan, Yonghong Tian, Bo XU, and Guoqi Li. Spike-driven transformer. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023b. URL https://openreview.net/forum?id=9FmolyOHi5.
- Man Yao, Hengyu Zhang, Guangshe Zhao, Xiyu Zhang, Dingheng Wang, Gang Cao, and Guoqi Li. Sparser spiking activity can be better: Feature refine-and-mask spiking neural network for event-based visual recognition. *Neural Networks*, 166:410–423, 2023c.

- Man Yao, Guangshe Zhao, Hengyu Zhang, Yifan Hu, Lei Deng, Yonghong Tian, Bo Xu, and Guoqi Li. Attention spiking neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 45(8):9393–9410, 2023d.
- Bojian Yin, Federico Corradi, and Sander M Bohté. Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks. *Nature Machine Intelligence*, 3(10):905–913, 2021.
- Weihao Yu, Mi Luo, Pan Zhou, Chenyang Si, Yichen Zhou, Xinchao Wang, Jiashi Feng, and Shuicheng Yan. Metaformer is actually what you need for vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 10819–10829, 2022a.
- Weihao Yu, Chenyang Si, Pan Zhou, Mi Luo, Yichen Zhou, Jiashi Feng, Shuicheng Yan, and Xinchao Wang. Metaformer baselines for vision. *arXiv preprint arXiv:2210.13452*, 2022b.
- Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 558–567, 2021.
- Hang Zhang, Chongruo Wu, Zhongyue Zhang, Yi Zhu, Haibin Lin, Zhi Zhang, Yue Sun, Tong He, Jonas Mueller, R Manmatha, et al. Resnest: Split-attention networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2736–2746, 2022a.
- Hong Zhang, Yang Li, Bin He, Xiongfei Fan, Yue Wang, and Yu Zhang. Direct training highperformance spiking neural networks for object recognition and detection. *Frontiers in Neuroscience*, 17, 2023.
- Jiqing Zhang, Bo Dong, Haiwei Zhang, Jianchuan Ding, Felix Heide, Baocai Yin, and Xin Yang. Spiking transformers for event-based single object tracking. In *Proceedings of the IEEE/CVF* conference on Computer Vision and Pattern Recognition, pp. 8801–8810, 2022b.
- Jiyuan Zhang, Lulu Tang, Zhaofei Yu, Jiwen Lu, and Tiejun Huang. Spike transformer: Monocular depth estimation for spiking camera. In *European Conference on Computer Vision*, pp. 34–52. Springer, 2022c.
- Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li. Going deeper with directly-trained larger spiking neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 11062–11070, 2021.
- Bolei Zhou, Hang Zhao, Xavier Puig, Sanja Fidler, Adela Barriuso, and Antonio Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 633–641, 2017.
- Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang, Shuicheng Yan, Yonghong Tian, and Li Yuan. Spikformer: When spiking neural network meets transformer. In *The Eleventh International Conference on Learning Representations*, 2023.