

Temporal-adaptive Weight Quantization for Spiking Neural Networks

Han Zhang, Qingyan Meng, Jiaqi Wang, Baiyu Chen, Zhengyu Ma, Xiaopeng Fan *Senior Member, IEEE*

Abstract—Weight quantization in spiking neural networks (SNNs) could further reduce energy consumption. However, quantizing weights without sacrificing accuracy remains challenging. In this study, inspired by astrocyte-mediated synaptic modulation in the biological nervous systems, we propose Temporal-adaptive Weight Quantization (TaWQ), which incorporates weight quantization with temporal dynamics to adaptively allocate ultra-low-bit weights along the temporal dimension. Extensive experiments on static (e.g., ImageNet) and neuromorphic (e.g., CIFAR10-DVS) datasets demonstrate that our TaWQ maintains high energy efficiency (4.12M, 0.63mJ) while incurring a negligible quantization loss of only 0.22% on ImageNet.

Index Terms—Spiking neural network, weight quantization, astrocyte, energy-efficient.

I. INTRODUCTION

Inspired by biological nervous systems, spiking neural networks (SNNs) are regarded as the third-generation neural networks [1]. Their event-driven nature, which transmits information through binary spikes [2], [3], enables accumulation (AC) operations to substitute multiply-accumulate (MAC) in the neural network. This paradigm shift in computation results in low energy consumption, making SNNs conducive to being deployed on resource-constrained devices.

Although SNNs already improve energy efficiency relative to artificial neural networks, they still fall short of the remarkable efficiency achieved by biological nervous systems. For instance, the human brain sustains about 8.6×10^6 neurons and over 1×10^{14} synapses while operating at approximately 20 watts [4]–[6]. In contrast, the energy efficiency of SNNs deployed on neuromorphic chips remains significantly inferior to this biological benchmark [7]–[9]. In the nervous system, a presynaptic spike triggers the release of discrete,

neurotransmitter-filled vesicles into the synaptic cleft. Since each vesicle carries a fixed, quantal amount of neurotransmitter, the synaptic strength depends on the discrete number of vesicles released [10]–[12]. This intrinsic transmitter quantization provides a biological analog for weight quantization strategies in SNNs: substituting full-precision weights with ultra-low-bit ones that count discrete quanta [13], yielding exponential energy savings.

Quantization of weights has been generally exploited in artificial neural networks (ANNs), where full-precision weights are quantized to no more than 8-bit, and in some cases, as low as 1-bit [14], [15]. In contrast, weight quantization of SNNs is only beginning to emerge. There are only a few related studies demonstrating that coupling low-bit weights with event-driven spikes can drastically cut energy consumption: the QSD-Transformer [16] quantizes weights to 4-bit precision, while Q-SNNs [13], BESTformer [17], and AGMM [18] employ 1-bit quantization. Nevertheless, weight quantization in SNNs without sacrificing accuracy remains a significant challenge, underscoring the need for biologically grounded strategies that narrow the efficiency gap with the nervous systems.

Astrocytes are widely distributed in the biological nervous system, forming tripartite synapses with excitatory or inhibitory presynaptic and postsynaptic neurons [19], and exhibit the ability to modulate synaptic strength across time [20], [21], eliminate synapses [22], [23], and facilitate synapse formation [24], which critically depend on intracellular calcium concentration oscillations. As illustrated in Fig. 1(d), we abstract astrocyte-mediated synaptic modulation into three principles: i) Under the modulation of astrocytes, synaptic strengths vary across time with invariant signs; ii) Astrocytes can eliminate some synapses, that is, convert excitatory or inhibitory synapses to asynaptic state; iii) Astrocytes can secrete thrombospondin and Hevin to facilitate the synapse formation, that is, convert the asynaptic state to excitatory or inhibitory synapses. The asynaptic state occurs either after the elimination of synapses or prior to synapse formation.

Inspired by the tripartite synapse, we propose a novel quantization method termed Temporal-adaptive Weight Quantization (TaWQ), which integrates mechanisms with temporal dynamics into the quantization process, enabling weights to adopt distinct values across timesteps, so as to mitigate the challenge of performance degradation in weight quantization. This emulates the role of astrocytes in the modulation of synaptic strength by following the three abstract principles. This method constrains synaptic connectivity to

This work was supported by National Science and Technology Innovation 2030 Major Project (No. 2025ZD0215501). This work was also supported in part by the National Natural Science Foundation of China (NSFC) under grant U22B2035. (Corresponding authors: Zhengyu Ma, Xiaopeng Fan.)

Han Zhang and Xiaopeng Fan are with the Faculty of Computing, Harbin Institute of Technology, Harbin 150001, China, and with Pengcheng Laboratory, Nanshan, Shenzhen 518000, China. Xiaopeng Fan is also with Suzhou Research Institute, Harbin Institute of Technology, Suzhou 215104, China (e-mail: 23b303002@stu.hit.edu.cn, fxp@hit.edu.cn).

Jiaqi Wang is with the Institute of Computing and Intelligence (ICI), Harbin Institute of Technology, Shenzhen 518055, China, and with Pengcheng Laboratory, Nanshan, Shenzhen 518000, China.

Baiyu Chen is with Institute of Automation, Chinese Academy of Sciences, Beijing 100190, China, and with Pengcheng Laboratory, Nanshan, Shenzhen 518000, China.

Qingyan Meng and Zhengyu Ma are with Pengcheng Laboratory, Nanshan, Shenzhen 518000, China (e-mail: mengqy@pcl.ac.cn, mazhy@pcl.ac.cn).

The code is available at <https://github.com/ZhangHanN1/TaWQ>

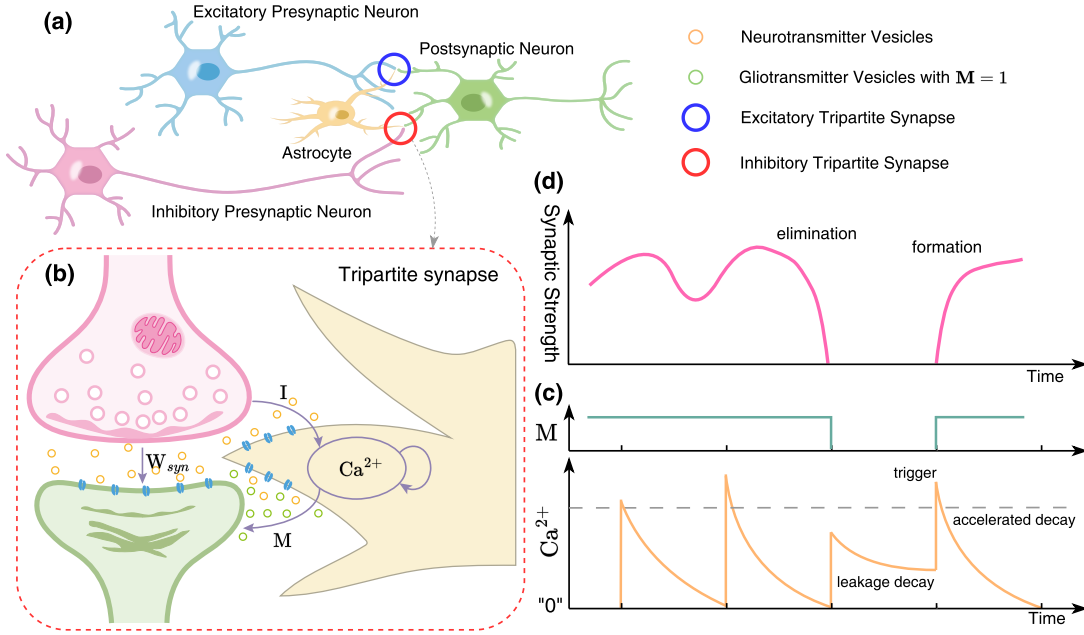


Fig. 1. Schematic illustration of tripartite synapses. (a) An excitatory or inhibitory presynaptic neuron and a postsynaptic neuron, together with an astrocyte, form a tripartite synapse. (b) Schematic of the tripartite synapse structure, W_{syn} is the synaptic strength, I represents the stimulus received by astrocytes, Ca^{2+} denotes the calcium concentration, M is a symbol (not spike) designating whether astrocytes modulate synapses. (c) Calcium dynamics curve, triggering the M upon exceeding the threshold. (d) The synaptic strength varies over time under the modulation of astrocytes.

three discrete states: excitatory synapses, asynaptic state, and inhibitory synapses, thereby converting full-precision floating-point weights into 1.58-bit ternary values $\{+1, 0, -1\}$.

Our main contributions are as follows. (1) First, we focus on the intracellular calcium concentration oscillations within astrocytes, extracting and modeling key characteristics of calcium dynamics, as astrocytic function is critically dependent on fluctuations in calcium concentration. (2) Based on the calcium dynamics model, we develop the Temporal-adaptive Weight Quantization (TaWQ) method, which integrates temporal dynamics into the weight quantization. Without additional trainable parameters, TaWQ quantizes full-precision floating-point weights into time-varying 1.58-bit ternary values $\{+1, 0, -1\}$, corresponding to three synaptic states: excitatory, asynaptic, and inhibitory. (3) Extensive experiments on both static and neuromorphic datasets demonstrate that our TaWQ narrows the accuracy gap between full-precision and ultra-low-bit SNNs while enhancing their intrinsic energy advantages.

II. RELATED WORKS

Ultra-Low-Bit Weight Quantization in Artificial Neural Networks. XNOR-Net [14] is a classical ultra-low-bit quantization method that compresses full-precision weights and activations into $\{+1, -1\}$. The dot product between binary vectors is implemented via XNOR-bitcounting operations, while the introduction of scaling factors mitigates accuracy degradation caused by quantization. DoReFa-Net [25] pioneers the comprehensive ultra-low-bit quantization of weights (1-bit), activations (2-bit), and gradients (6-bit), significantly reducing overhead and enhancing training efficiency. Bi-Real Net [26] enhances the network's representational capacity by introducing an optimized residual structure and a refined parameter updating algorithm. ReActNet [27] proposes generalized activation functions and a distributional loss to narrow

the performance gap to full-precision baselines. Although these methods achieve remarkable hardware efficiency, they incur a noticeable performance degradation. BitNet [15] and its ternary extension BitNet-1.58 [28] demonstrate that ultra-low-bit weight quantization can preserve high performance while yielding substantial energy savings. However, their activations remain at 8-bit, leaving considerable potential for further efficiency refinement when both weights and activations are pushed into the binary-ternary regime.

Ultra-Low-Bit Weight Quantization in Spiking Neural Networks. Event-driven SNNs naturally achieve low-power inference since information is transmitted only through binary spikes, utilizing sparse and additive synaptic operations instead of multiply-accumulate operations. Weight quantization could further amplify these efficiency advantages. QP-SNNs [29] employs a weight rescaling strategy that efficiently utilizes bit-width to enhance uniform quantization methods, thereby strengthening the model's representational capability, and the accuracy on ImageNet is only 61.36% with 8-bit weights. QSD-Transformer [16] employs an information-enhanced LIF neuron and fine-grained distillation to mitigate performance degradation. The accuracy on ImageNet is 80.3%. However, 4-bit quantization retains potential for further bit-width reduction, as some works have implemented 1-bit weight quantization. The Q-SNNs [13] binarize synaptic weights into $\{+1, -1\}$, and introduce a loss function to constrain neuronal firing rates approaching 0.5, and further reduce memory footprint by implementing 8-bit quantization on membrane potentials, but it exhibits noticeable performance degradation of 0.82% and 0.75% on the CIFAR-10 and CIFAR-100 datasets, respectively. AGMM [18] proposes an adaptive gradient modulation to ensure the performance after quantization, with 64.67% accuracy on ImageNet. In addition to

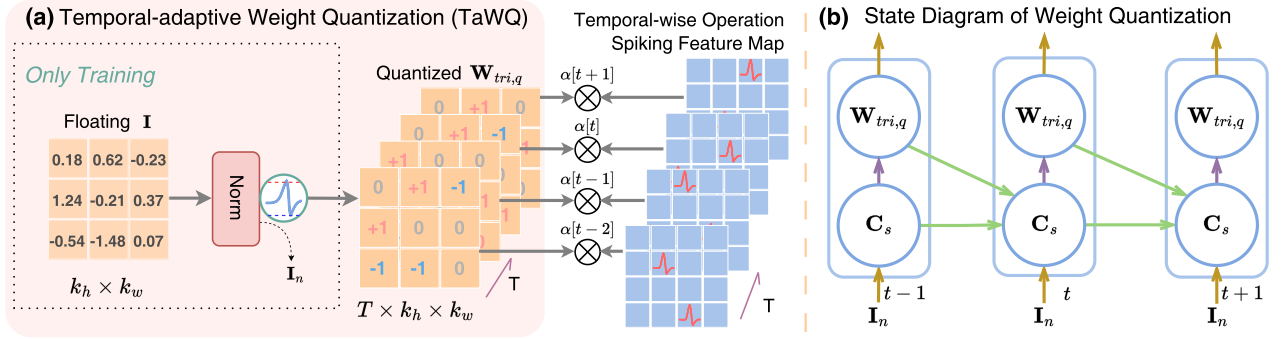


Fig. 2. Schematic illustration of TaWQ. (a) Weights are quantized into time-varying 1.58-bit values $\{+1, 0, -1\}$, followed by temporal-wise operation. (b) The state diagram in the weight quantization process, \mathbf{I}_n , \mathbf{C}_s , and $\mathbf{W}_{tri,q}$ are the normalized stimulus, intermediate variable, and quantized weight, respectively.

the network weights, BESTformer [17] further quantizes the self-attention map to 1-bit, with the resultant performance degradation being compensated by the proposed Coupled Information Enhancement strategy, the accuracy is 63.46% on ImageNet. While existing methods attain efficient computation and model compression via integrating SNNs with ultra-low-bit quantization, they still require further optimization to narrow the performance gap with full-precision networks.

III. METHOD

A. Calcium Dynamics in Astrocytes

Synaptic states in biological nervous systems persistently evolve throughout development, and even in mature systems, they remain dynamic to enable complex functionalities. This phenomenon underscores the necessity of integrating SNNs with temporally evolving synaptic states, such integration significantly enhances the representational capacity of SNNs, thus mirroring the biological characteristic wherein synaptic dynamics facilitate sophisticated neural functions. Astrocytes are widely distributed throughout biological nervous systems and actively participate in the modulation of synapses. The astrocytic modulation of synapses represents a critical mechanism for enabling changes in synaptic states. As illustrated in Fig. 1(b), this modulation process relies on intracellular calcium concentration oscillations within astrocytes.

These calcium dynamics exhibit three key characteristics: 1) calcium concentration accumulates over time [30], 2) calcium concentration undergoes leakage decay [31], and 3) accelerated decay occurs upon exceeding a calcium concentration threshold [32]. These principles form the foundation of our motivation. We deliberately model astrocytic calcium dynamics by collapsing the temporal aspect, expressed mathematically as follows:

$$\begin{cases} \text{Ca}^{2+}[t+1] = \kappa_1 \cdot \text{Ca}^{2+}[t](1 - \mathbf{M}[t]) + \kappa_2 \cdot |\mathbf{I}|, \\ \mathbf{M}[t+1] = \mathcal{H}(\text{Ca}^{2+}[t+1], C_{th}), \\ \mathbf{W}_{tri}[t+1] = \mathcal{F}(\mathbf{M}[t+1], \mathbf{W}_{syn}), \end{cases} \quad (1)$$

where κ_1 and κ_2 are scaling factors, $\text{Ca}^{2+}[t]$ denotes the calcium concentration at timestep t , $\mathbf{M}[t]$ is a symbol designating whether astrocytes modulate synapses at timestep t , not a spike. \mathbf{I} represents the stimulus received by astrocytes. A non-negative value for $\text{Ca}^{2+}[t]$ is guaranteed by taking the absolute value of \mathbf{I} . $\mathcal{H}(\text{Ca}^{2+}, C_{th})$ is the Heaviside step

function, which equals 1 if $\text{Ca}^{2+} \geq C_{th}$, otherwise equals 0. And \mathbf{W}_{syn} denotes the synaptic strength, \mathbf{W}_{tri} represents the astrocyte-modulated synaptic strength, and $\mathcal{F}(\cdot)$ signifies the modulation function, as illustrated in Fig. 1(a).

Eq. (1) models key characteristics of calcium dynamics. The calcium concentration accumulates over time, with the scaling factor κ_1 controlling the rate of leakage decay, and when $\mathbf{M} = 1$, it induces an accelerated decrease, as illustrated in Fig. 1(c), in which the label "0" denotes the calcium concentration in the absence of stimulus. Although the calcium concentration is not absolutely constant without stimulus, we simplify the model by setting it to 0. This means that when the stimulus $\mathbf{I} = 0$, the resulting calcium concentration $\text{Ca}^{2+}[t] = 0$.

B. Temporal-adaptive Weight Quantization

To narrow the performance gap induced by weight quantization, we draw inspiration from tripartite synapses and introduce dynamics into the quantization based on the calcium dynamics model described in Eq. (1), proposing a novel SNNs-specific weight quantization method termed Temporal-adaptive Weight Quantization (TaWQ). The quantized weights exhibit temporal variability, adaptively switching between synaptic (excitatory or inhibitory) and asynaptic states. This behavior aligns with astrocyte-mediated synaptic modulation.

1.58-bit TaWQ. We first normalize the stimulus \mathbf{I} to ensure that they are subject to a distribution $\mathcal{N}(0, 1)$, which follows [43], and it can be formulated as:

$$\mathbf{I}_n = \frac{\mathbf{I} - \mu_I}{\sqrt{\sigma_I^2 + \epsilon}}. \quad (2)$$

The stimulus \mathbf{I} and the weight \mathbf{W}_{syn} share the same shape, indicating that each synapse receives its respective stimulus input. $\mathbf{W}_{syn}, \mathbf{I} \in \mathbb{R}^{C_o \times C_i \times k_h \times k_w}$ and $\mathbf{I}_n \in \mathbb{R}^{C_o \times C_i \times k_h \times k_w}$ denote the stimulus before and after normalization, where C_o and C_i represent the number of output and input channels in convolution layers (Conv), k_h and k_w denote the height and width of the convolution kernel, and μ_I and σ_I denote the mean and standard deviation of \mathbf{I} , respectively. ϵ is an infinitesimal constant introduced to avoid division-by-zero scenarios. The stimulus generated by excitatory synapses should correspondingly induce excitation, and the same logic applies to inhibitory synapses. Therefore, the sign of \mathbf{I}_n should be consistent with that of \mathbf{W}_{syn} , that is, $\text{sign}(\mathbf{I}_n) = \text{sign}(\mathbf{W}_{syn})$.

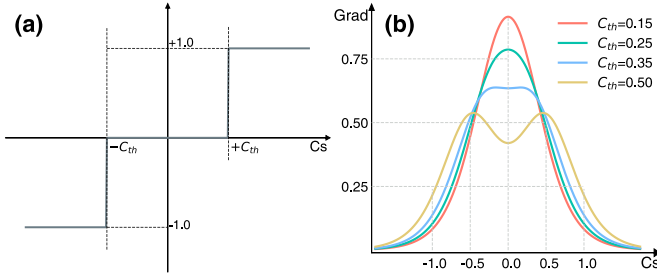


Fig. 3. Curve of the quantization function and surrogate gradient, "Cs" on the horizontal axis is C_s in Eq. (6). (a) The quantization function converts floating-point values into 1.58-bit ternary values $\{+1, 0, -1\}$. (b) Surrogate gradient under varying thresholds C_{th} .

Subsequently, we define the $\mathcal{F}(\mathbf{M}, \mathbf{W}_{syn}) = \mathbf{M} \cdot \mathbf{W}_{syn}$. To focus solely on the presence or absence of synaptic connections rather than their specific strengths, the matrix \mathbf{W}_{syn} is replaced with $\text{sign}(\mathbf{W}_{syn})$, which effectively reduces the function to:

$$\mathcal{F}(\mathbf{M}, \mathbf{W}_{syn}) = \mathbf{M} \cdot \text{sign}(\mathbf{W}_{syn}) = \mathbf{M} \cdot \text{sign}(\mathbf{I}_n). \quad (3)$$

Given the constraint $\kappa_1 + \kappa_2 = 1$, if $\kappa_1 = \lambda$, then $\kappa_2 = 1 - \lambda$. The expression for quantized \mathbf{W}_{tri} can be rewritten as:

$$\begin{cases} \mathbf{Ca}^{2+}[t+1] = \lambda \cdot \mathbf{Ca}^{2+}[t](1 - \mathbf{W}_{tri,q}[t]) + (1 - \lambda) \cdot |\mathbf{I}_n|, \\ \mathbf{W}_{tri,q}[t+1] = \mathcal{H}(\mathbf{Ca}^{2+}[t+1], C_{th}) \cdot \text{sign}(\mathbf{I}_n). \end{cases} \quad (4)$$

The $\mathcal{H}(\mathbf{Ca}^{2+}, C_{th}) \cdot \text{sign}(\mathbf{I}_n)$ is equivalent to the quantization function $\mathcal{S}(\mathbf{Ca}^{2+} \cdot \text{sign}(\mathbf{I}_n), C_{th})$, which is defined as Eq. (5) and illustrated in Fig. 3(a):

$$\mathcal{S}(\mathbf{Ca}^{2+} \cdot \text{sign}(\mathbf{I}_n), C_{th}) = \begin{cases} +1, & \mathbf{Ca}^{2+} \cdot \text{sign}(\mathbf{I}_n) > +C_{th}, \\ -1, & \mathbf{Ca}^{2+} \cdot \text{sign}(\mathbf{I}_n) < -C_{th}, \\ 0, & \text{else.} \end{cases} \quad (5)$$

To distinguish the effects of excitatory/inhibitory stimulus on calcium concentration while simplifying the expression, we use intermediate variable C_s in the real domain to denote the $\mathbf{Ca}^{2+} \cdot \text{sign}(\mathbf{I}_n)$, thereby obtaining the expression for the 1.58-bit TaWQ as follows:

$$\begin{cases} C_s[t+1] = \lambda \cdot C_s[t](1 - |\mathbf{W}_{tri,q}[t]|) + (1 - \lambda) \cdot \mathbf{I}_n, \\ \mathbf{W}_{tri,q}[t+1] = \mathcal{S}(C_s[t+1], C_{th}), \end{cases} \quad (6)$$

The forward propagation state diagram of C_s is illustrated in Fig. 2(b). $\mathbf{W}_{tri,q} \in \mathbb{R}^{T \times C_o \times C_i \times k_h \times k_w}$ represents the quantized 1.58-bit weight, T is timesteps, $\lambda = 0.5$ serves as the scalar scaling coefficient, and C_{th} indicates the quantization threshold. For Linear layers, both k_h and k_w are set to 1.

The TaWQ in SNNs quantizes weights \mathbf{W}_{syn} into 1.58-bit time-varying $\mathbf{W}_{tri,q}$. Following the TaWQ, we subsequently introduce a temporal-wise scaling factor to compensate for performance loss caused by weight quantization. Let $\mathbf{X}_i \in \mathbb{R}^{T \times C_i \times H \times W}$ denote the input to the Quantized Conv Layer (Q-Conv), where H and W represent the height and width of \mathbf{X}_i , respectively. The quantized weights of Q-Conv are $\mathbf{W}_{tri,q}$. As illustrated in Fig. 2(a), the computation of the Q-Conv can be formulated as:

$$\mathbf{X}_o[t] = (\alpha[t] \odot \mathbf{W}_{tri,q}[t]) \otimes \mathbf{X}_i[t] = \alpha[t] \odot (\mathbf{W}_{tri,q}[t] \otimes \mathbf{X}_i[t]). \quad (7)$$

In the above expression, $\mathbf{X}_o \in \mathbb{R}^{T \times C_o \times H' \times W'}$ is output feature map of Q-Conv, \otimes denotes matrix product, \odot is an element-wise product, the parameter $\alpha \in \mathbb{R}^{T \times C_o}$ is the temporal-wise scaling factor. After training, α becomes fixed and can be folded into subsequent computational steps during inference. α defined as:

$$\alpha[t, c] = \phi\left(\frac{1}{C_i k_h k_w} \sum_{i=0}^{C_i-1} \sum_{j=0}^{k_h-1} \sum_{k=0}^{k_w-1} |\mathbf{W}_{tri,q}[t, c, i, j, k]|\right). \quad (8)$$

The $\phi(\cdot)$ denotes the reciprocal function. For multiple timesteps, the computation of the Q-Conv is shown as follows:

$$\mathbf{X}_o = \text{Stack}(\mathbf{X}_o[0], \mathbf{X}_o[1], \dots, \mathbf{X}_o[T-1]), \quad (9)$$

During the inference phase, the temporal-wise scaling factor α of Q-Conv and the trained parameters of batch normalization (BN) can be jointly folded into the spiking neuron. If the LIF neurons are employed in SNNs, when the reset membrane potential is 0, their charging process can be described as:

$$\begin{aligned} \mathbf{U}[t] &= (1 - \frac{1}{\tau})\mathbf{U}[t-1] + \frac{1}{\tau}(\gamma \frac{\alpha[t-1] \odot \mathbf{X}_q[t-1] - \mu}{\sqrt{\sigma^2 + \epsilon}} + \beta) \\ &= (1 - \frac{1}{\tau})\mathbf{U}[t-1] + \rho[t-1] \odot \mathbf{X}_q[t-1] + \delta, \end{aligned} \quad (10)$$

$$\rho[t-1] = \frac{\gamma \alpha[t-1]}{\tau \sqrt{\sigma^2 + \epsilon}}, \quad \delta = \frac{1}{\tau}(\beta - \frac{\gamma \mu}{\sqrt{\sigma^2 + \epsilon}}). \quad (11)$$

$\mathbf{U}[t]$ denotes the accumulated membrane potential, $\mathbf{X}_q[t] = \mathbf{W}_{tri,q}[t] \otimes \mathbf{X}_i[t]$, τ represents the membrane time constant, μ and σ are the mean and standard deviation of $\alpha \odot \mathbf{X}_q$, while γ and β correspond to the scaling and shift factors of BN.

Backpropagation. Gradients of \mathbf{I} are calculated as $\frac{\partial L}{\partial \mathbf{I}} = \frac{\partial L}{\partial \mathbf{I}_n} \frac{\partial \mathbf{I}_n}{\partial \mathbf{I}}$, and $\frac{\partial L}{\partial \mathbf{I}_n}$ is as follows:

$$\begin{aligned} \frac{\partial L}{\partial \mathbf{I}_n} &= \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{W}_{tri,q}[t]} \frac{\partial \mathbf{W}_{tri,q}[t]}{\partial C_s[t]} \left(\frac{\partial C_s[t]}{\partial \mathbf{I}_n} \right. \\ &\quad + \sum_{j < t} \prod_{i=1}^{t-j} \left(\frac{\partial C_s[t-i+1]}{\partial C_s[t-i]} \right. \\ &\quad \left. \left. + \frac{\partial C_s[t-i+1]}{\partial \mathbf{W}_{tri,q}[t-i]} \frac{\partial \mathbf{W}_{tri,q}[t-i]}{\partial C_s[t-i]} \right) \frac{\partial C_s[j]}{\partial \mathbf{I}_n} \right), \end{aligned} \quad (12)$$

where the term $\frac{\partial \mathbf{W}_{tri,q}[t]}{\partial C_s[t]}$ denotes the derivative of $\mathcal{S}(C_s, C_{th})$, which is non-differentiable and will be replaced by the surrogate gradient in Fig. 3(b), formulated as follows:

$$\frac{\partial \mathbf{W}_{tri,q}}{\partial C_s} = \frac{1}{2}(\theta'(4(C_s + C_{th})) + \theta'(4(C_s - C_{th}))). \quad (13)$$

The value of $\frac{\partial \mathbf{W}_{tri,q}}{\partial C_s}$ is composed of two surrogate gradient, where $\theta'(4(C_s - C_{th}))$ and $\theta'(4(C_s + C_{th}))$ are derivative of Sigmoid. The convergence analysis for Eq. (12) is provided in the Supplementary Material. The update rule for \mathbf{I} is:

$$\mathbf{I}^+ = \mathbf{I}^- - \eta \cdot \mathbf{G}, \quad (14)$$

where η denotes the learning rate, and \mathbf{G} represents a gradient-related term whose formulation varies depending on the optimizer employed.

TABLE I

RESULTS ON THE IMAGENET DATASET. "POWER (MJ)" INDICATES ENERGY CONSUMPTION, AND "ACC (%)" IS THE TOP-1 ACCURACY. THE POWER "0.99(2.54%)" DENOTES AN ENERGY CONSUMPTION OF 0.99MJ, WHICH REPRESENTS ONLY 2.54% OF THE FULL-PRECISION COUNTERPART. THE ACCURACY OF ALL COMPARISON METHODS IS SOURCED FROM THEIR ORIGINAL RESEARCH PUBLICATIONS.

Method	Architecture	Weight Bits	TimeStep	Size(M)	Power(mJ)	Acc(%)
XNOR-Net [14]	ResNet18	1	1	-	-	51.2
Bi-Real Net [26]	Bi-Real-18	1	1	-	-	56.4
AdaBin [39]	ResNet18	1	1	-	-	66.4
ReActNet [27]	ReActNet-A	1	1	-	-	69.4
QP-SNN [29]	ResNet-18	8	4	13.28	-	61.36
BESTformer [17]	BESTformer-8-512	1	4	5.57	-	63.46
AGMM [18]	ResNet-18	1	4	-	-	64.67
QSDTransformer [16]	SDTransformer-v2-T	4	4	1.8	2.5	77.5
	SDTransformer-v2-M	4	4	3.9	5.7	78.9
	SDTransformer-v2-L	4	4	6.8	8.7	80.3
Spikformer [38]	Spikformer-8-384	32	4	16.81	7.73	70.24
SDTransformer [40]	SDTransformer-8-384	32	4	16.81	3.90	72.28
Spikingformer-CML [41]	Spikingformer-8-384	32	4	16.81	4.69	74.35
	Spikingformer-8-512	32	4	29.68	7.46	76.54
	Spikingformer-8-768	32	4	66.34	13.68	77.64
	Spikingformer-8-384	1.58	4	1.25(7.44%)	0.33(7.04%)	72.44(-1.91)
Spikingformer-TaWQ	Spikingformer-8-512	1.58	4	2.03(6.84%)	0.43(5.76%)	75.17(-1.37)
	Spikingformer-8-768	1.58	4	4.12(6.21%)	0.63(4.61%)	77.42(-0.22)
QKFormer [42]	QKFormer-10-768	32	4	64.96	38.91	84.22
QKFormer-TaWQ	QKFormer-10-768	1.58	4	4.05(6.23%)	0.99(2.54%)	82.94(-1.28)

TaWQ achieves dual optimization. On the one hand, it implements quantization by converting floating-point weights into 1.58-bit ternary values $\{+1, 0, -1\}$, corresponding to excitatory synapses, asynaptic state, and inhibitory synapses, respectively, further exploiting SNNs' inherent energy efficiency. On the other hand, we innovatively integrate temporal dynamics into the quantization, enabling time-varying synaptic strength in quantized SNNs, which significantly enriches the representational capacity of quantized SNNs.

IV. EXPERIMENTS

We evaluate the proposed quantization method on static datasets (ImageNet [33], CIFAR-10/100 [34]), neuromorphic datasets (CIFAR-10-DVS [35], DVS128-Gesture [36]), and the speech dataset (SHD [37], in Supplementary Material). And we provide the theoretical energy consumption calculation method for TaWQ in the Supplementary Material, along with the detailed process for estimating energy consumption that incorporates the overhead associated with reading/writing weights and feature maps on hardware. Detailed experimental setups are also provided in the Supplementary Material.

A. Results on ImageNet Classification

Comparing with other SNNs. The results on ImageNet are summarized in Table I. All models in the table employ 1-bit activations, except for QSD-Transformer. The Spikingformer-TaWQ models achieve compact model sizes of 1.25M, 2.03M, and 4.12M, while attaining top-1 accuracies of 72.44%, 75.17%, and 77.42%, with corresponding energy consumptions of 0.33mJ, 0.43mJ, and 0.63mJ, respectively. The QKFormer-TaWQ achieves a model size of 4.05M, an accuracy of 82.94%, and an energy consumption of only 0.99mJ. **i) Full-precision SNNs.** The model size of Spikingformer-TaWQ is only 7.44%, 6.84%, and 6.21% of its full-precision

counterpart, Spikingformer-CML, while its energy consumption corresponds to merely 7.04%, 5.76%, and 4.61% of the latter. Notably, the performance degradation is minor, with Spikingformer-TaWQ exhibiting an accuracy reduction of only 0.22% in the "8-768" architecture. Furthermore, in the "8-384" architecture, Spikingformer-TaWQ delivers higher accuracy than both the full-precision Spikformer and SD-Transformer. The model size of the QKFormer-TaWQ is only 6.23% of its full-precision counterpart, QKFormer, with an energy consumption of merely 2.54% and a minor performance degradation of 1.28%. **ii) Quantized SNNs.** Regarding quantized SNNs, QP-SNN, BESTformer, and AGMM exhibit limited performance, whereas our Spikingformer-TaWQ and QKFormer-TaWQ demonstrate a notable advantage in accuracy over these methods. QSDTransformer, which quantizes weights to 4-bit, not only results in larger model sizes than our TaWQ-based models but also consumes significantly more energy. Notably, while QKFormer-TaWQ achieves 2.64% higher accuracy than QSDTransformer, its energy consumption is only 11.38%, and its model size is merely 59.56% of the latter.

Comparing with Binary Neural Networks (BNNs). We select representative ultra-low-bit (1-bit) quantized BNNs as control groups, including XNOR-Net, Bi-Real Net, AdaBin, and ReActNet, which do not mention the calculation of energy consumption. Spikingformer-TaWQ and QKFormer-TaWQ demonstrate substantial performance improvements over these BNNs. Spikingformer-TaWQ with "8-384" architecture outperforms XNOR-Net, Bi-Real Net, AdaBin, and ReActNet by 21.24%, 16.04%, 6.04%, and 3.04% in top-1 accuracy, respectively. And QKformer-TaWQ outperforms them by 31.74%, 26.54%, 16.54%, and 13.54%, respectively.

Attention Map Visualization. We conduct a comparative visualization of attention maps on the validation set for both the full-precision and the TaWQ-quantized model. The atten-

tion maps are computed from the last spiking self-attention module of the network and averaged across all timesteps. As shown in Fig. 4, the attention distributions of the models before and after quantization are largely consistent, focusing on the same key features, such as the eagle’s wings, the bird’s beak, and the cat’s whiskers. This indicates that the TaWQ-quantized model retains feature extraction capabilities largely equivalent to those of the full-precision model.

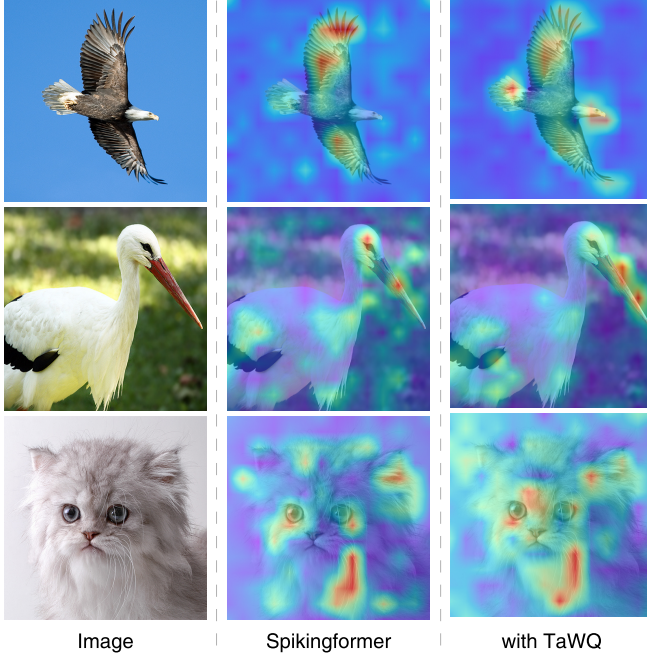


Fig. 4. Attention maps of the full-precision model and the 1.58-bit quantized model with TaWQ. The images are part of ImageNet’s validation set.

Furthermore, we present additional experimental results in the Supplementary Material, including but not limited to the latency and memory footprint of TaWQ-quantized networks, comparisons with Post-training Quantizations (PTQ), and TaWQ-based quantization results for non-Transformer spiking convolutional networks.

B. Results on CIFAR Classification

We quantize Spikformer, Spikingformer, and QKFormer on the CIFAR10 and CIFAR100 datasets, resulting in the quantized models Spikformer-TaWQ, Spikingformer-TaWQ, and QKFormer-TaWQ. The results are shown in Table II, i) **For CIFAR10**, Spikformer-TaWQ and Spikingformer-TaWQ both have a model size of 0.49M ($19.0\times$ smaller), while QKFormer-TaWQ achieves a smaller model size of 0.36M ($18.7\times$ smaller). The accuracy degradation is minor, with drops of only 0.32%, 0.27%, and 0.10%, respectively. Notably, QKFormer-TaWQ, with its compact size, attains an accuracy of 96.08%, outperforming quantized SNNs such as QP-SNN, Q-SNN, and BESTformer by margins of 0.67%, 0.54%, and 0.35%, respectively. Furthermore, it surpasses the full-precision Spikingformer’s accuracy. ii) **For CIFAR100**, the model sizes of Spikformer-TaWQ and Spikingformer-TaWQ are both 0.53M ($17.6\times$ smaller). Their performance shows increases of 0.51% and 0.44%, respectively. QKFormer-TaWQ achieves a model size of 0.39M ($17.3\times$ smaller), with

only a 0.30% accuracy degradation. Notably, Spikingformer-TaWQ attains the highest accuracy of 80.87%, significantly outperforming other quantized SNNs and surpassing all full-precision models except QKFormer. Additionally, a statistical analysis of the firing rates on CIFAR100 is provided in the Supplementary Material.

TABLE II

COMPARISON OF TAWQ’S RESULTS ON CIFAR10 AND CIFAR100. “BITS” DENOTES THE BIT-WIDTH OF WEIGHTS. THE UNIT OF “SIZE” IS “M”. “+TAWQ” DENOTES “XXX-TAWQ”, WHERE “XXX” REFERS TO THE FULL-PRECISION MODEL.

Method	T	Bits	CIFAR10		CIFAR100	
			Size	Acc(%)	Size	Acc(%)
QP-SNN [29]	2	4	3.16	95.41	3.35	75.77
Q-SNN [13]	2	1	1.62	95.54	-	78.82
BESTformer [17]	4	1	1.18	95.73	1.31	79.80
Spikformer [38]	4	32	9.32	95.51	9.32	78.21
+TaWQ	4	1.58	0.49	95.19(-0.32)	0.53	78.72(+0.51)
Spikingformer [41]	4	32	9.32	96.04	9.32	80.37
+TaWQ	4	1.58	0.49	95.77(-0.27)	0.53	80.87(+0.44)
QKFormer [42]	4	32	6.74	96.18	6.74	81.15
+TaWQ	4	1.58	0.36	96.08(-0.10)	0.39	80.85(-0.30)

C. Results on Neuromorphic Classification

Following the static CIFAR datasets protocol, we implement 1.58-bit quantization on Spikformer, Spikingformer, and QKFormer by TaWQ. The results are shown in Table III. In terms of model size, the quantized Spikformer-TaWQ and Spikingformer-TaWQ achieve a compact size of 0.14M, representing merely 5.4% of their full-precision counterparts. QKFormer-TaWQ attains even higher compression with 0.08M, which is only 5.3% of the full-precision counterpart. Accuracy evaluation reveals minor performance degradation. For CIFAR10-DVS, Spikformer-TaWQ and QKFormer-TaWQ show marginal decreases of 0.7% and 0.9% respectively, while Spikingformer-TaWQ demonstrates a 0.3% improvement. For DVS128-Gesture, Spikformer-TaWQ maintains parity with Spikformer, while Spikingformer-TaWQ and QKFormer-TaWQ exhibit minor accuracy reductions of 0.3% and 0.7%, respectively.

TABLE III

COMPARISON OF TAWQ’S RESULTS ON CIFAR10-DVS AND DVS128-GESTURE. THE UNIT OF “SIZE” IS (M), AND “BITS” DENOTES THE BIT-WIDTH OF WEIGHTS. “+TAWQ” DENOTES “XXX-TAWQ”, WHERE “XXX” REFERS TO THE FULL-PRECISION MODEL.

Method	CIFAR10-DVS				DVS128-Gesture			
	Size	T	Bits	Acc(%)	Size	T	Bits	Acc(%)
QP-SNN [29]	1.61	10	8	82.1	-	-	-	-
Q-SNN [13]	-	10	1	81.6	-	16	1	97.9
BESTformer [17]	1.18	16	1	80.8	-	-	-	-
Spikformer [38]	2.57	16	32	80.9	2.57	16	32	98.3
+TaWQ	0.14	16	1.58	80.2(-0.7)	0.14	16	1.58	98.3(-0.0)
Spikingformer [41]	2.57	16	32	81.4	2.57	16	32	98.6
+TaWQ	0.14	16	1.58	81.7(+0.3)	0.14	16	1.58	98.3(-0.3)
QKFormer [42]	1.50	16	32	84.0	1.50	16	32	98.6
+TaWQ	0.08	16	1.58	83.1(-0.9)	0.08	16	1.58	97.9(-0.7)

D. Information Entropy of Quantized Weights

To assess how TaWQ shapes the weight distribution, we employ information entropy of the quantized weights, with

the calculation method detailed in the Supplementary Material. A perfectly balanced ternary weight distribution, equal probabilities for +1, 0, and -1 (i.e., each at $1/3$), yields the theoretical maximum entropy of 1.0986 nats. As summarized in Fig. 5, after training QKFormer-TaWQ on CIFAR100, every layer converges to an almost balanced ternary profile: the mean probabilities across all layers are $P_p = 0.3283$, $P_z = 0.3320$, and $P_n = 0.3396$ (P_p , P_z , and P_n denote the probabilities of +1, 0, and -1 in the weights, respectively). This corresponds to an average entropy of 1.0952 nats, just 0.0034 nats below the optimum (1.0986), indicating that the network exploits the full expressive capacity of the ternary weights. Notably, although the initial quantized weight distribution under $C_{th} = 0.25$ deviates significantly from a uniform distribution, the training procedure drives all the ternary values towards nearly $1/3$, automatically maximizing weight entropy.

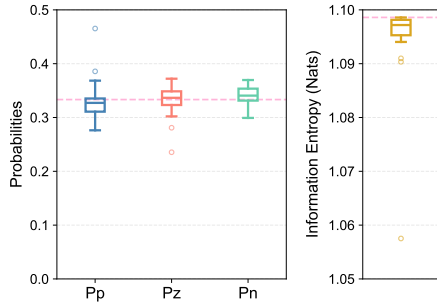


Fig. 5. Information entropy and weight proportion of TaWQ-quantized QKFormer. 'Pp', 'Pz', and 'Pn' represent the probabilities of +1, 0, and -1 in the weight, respectively. The pink dashed line denotes the optimum.

E. Ablation Study of Temporal-adaptive Dynamics

We conduct an ablation study on temporal-adaptive dynamics in TaWQ. After removing it, TaWQ degenerates into the $\mathcal{S}(C_s, C_{th})$ in Eq. (5), yielding Spikformer-WQ through quantization of Spikformer. The accuracy of Spikformer-TaWQ and Spikformer-WQ is shown in Table IV. On the DVS128-Gesture dataset, Spikformer-TaWQ outperforms Spikformer-WQ by a margin of 1.1% in accuracy. Similarly, on CIFAR100, the former exhibited a 0.85% higher accuracy. The proportion of weights in Spikformer-TaWQ more closely approaches the optimal value of $1/3$ compared to Spikformer-WQ. Consequently, TaWQ yields higher information entropy, as illustrated in Fig. 6 and Fig. 7. These results collectively substantiate the necessity of the temporal-adaptive dynamics. Additional ablation studies, including analyses of timesteps and bit-widths, are provided in the Supplementary Material.

TABLE IV
TEMPORAL DYNAMICS ABLATION STUDY RESULTS

Dataset	Method	Acc(%)
DVS128-Gesture	Spikformer-WQ	97.2
	Spikformer-TaWQ	98.3(+1.10)
CIFAR100	Spikformer-WQ	77.87
	Spikformer-TaWQ	78.72(+0.85)

V. CONCLUSION

We propose Temporal-adaptive Weight Quantization (TaWQ), an ultra-low-bit weight quantization method inspired

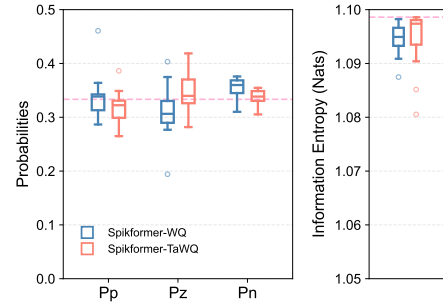


Fig. 6. Information entropy and weight proportion on static image dataset CIFAR100. The pink dashed line denotes the optimum.

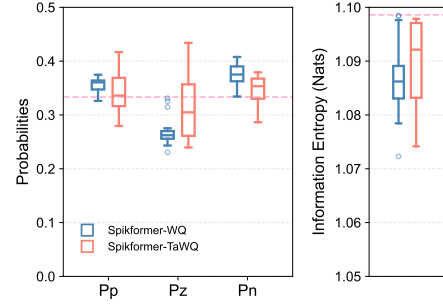


Fig. 7. Information entropy and weight proportion on neuromorphic dataset DVS128-Gesture. The pink dashed line denotes the optimum.

by astrocyte-mediated synaptic modulation. By integrating temporal dynamics into the weight quantization process, TaWQ allocates weights temporal-adaptively. Statistics on the learned weight distributions reveal that TaWQ drives the network towards a near-optimal ternary-weight regime, almost fully exploiting the expressiveness of low-bit weights. Extensive experiments on both static and neuromorphic datasets validate that TaWQ narrows the accuracy gap to full-precision models while delivering significant energy savings. The findings underscore TaWQ's potential to guide future neuromorphic algorithms/devices development that enables low-power SNNs deployment in the real world.

REFERENCES

- [1] M. Wolfgang, "Networks of spiking neurons: the third generation of neural network models," *Neural Netw.*, vol. 10, no. 9, pp. 1659–1671, 1997.
- [2] R. Kaushik, J. Akhilesh and P. Priyadarshini, "Towards Spike-based Machine Intelligence With Neuromorphic Computing," *Nature*, vol. 575, no. 7784, pp. 607–617, 2019.
- [3] M. Yao *et al*, "Scaling Spike-Driven Transformer With Efficient Spike Firing Approximation Training," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 47, no. 4, pp. 2973–2990, 2025.
- [4] H. Suzana, "The human brain in numbers: a linearly scaled-up primate brain," *Front. Neurosci.*, vol. 3, pp. 857, 2009.
- [5] Y. Hu, Q. Zheng, X. Jiang and G. Pan, "Fast-SNN: Fast Spiking Neural Network by Converting Quantized ANN," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 45, no. 12, pp. 14546–14562, 2023.
- [6] M. Yao *et al*, "Attention Spiking Neural Networks," *IEEE Trans. Pattern Anal. Machine Intell.*, vol. 47, no. 4, pp. 9393–9410, 2023.
- [7] A. Filipp *et al*, "Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip," *Nature*, vol. 34, no. 10, pp. 1537–1557, 2015.
- [8] P. Jing *et al*, "Towards artificial general intelligence with hybrid Tianjic chip architecture," *Nature*, vol. 572, no. 7767, pp. 106–111, 2019.
- [9] M. Davies, "Taking neuromorphic computing to the next level with Loihi2," *Intel Labs' Loihi*, vol. 2, no. 1, 2021.
- [10] R. H. Edwards, "The Neurotransmitter Cycle and Quantal Size," *Neuron*, vol. 55, no. 6, pp. 835–858, 2007.

- [11] Y. Yue *et al.*, “Distinct transmission sites within a synapse for strengthening and homeostasis,” *Sci. Adv.*, vol. 11, no. 15, pp. eads5750, 2025.
- [12] F.J. Urbano, E.S. Piedras-Renteria, K. Jun, H. Shin, O.D. Uchitel and R.W. Tsien, “Altered properties of quantal neurotransmitter release at endplates of mice lacking P/Q-type Ca²⁺ channels,” *Proc. Natl. Acad. Sci.*, vol. 100, no. 6, pp. 3491–3496, 2003.
- [13] W. Wei *et al.*, “Q-SNNs: Quantized Spiking Neural Networks,” in *Proc. 32nd ACM Int. Conf. Multimedia*, New York, NY, USA, 2024, pp. 8441–8450.
- [14] M. Rastegari, V. Ordonez, J. Redmon and A. Farhadi, “XNOR-Net: ImageNet Classification Using Binary Convolutional Neural Networks,” in *Proc. Eur. Conf. Comput. Vis.*, Amsterdam, Netherlands, 2016, pp. 525–542. DOI: 10.1007/978-3-319-46493-0_32.
- [15] H. Wang *et al.*, “BitNet: Scaling 1-bit Transformers for Large Language Models,” 2023, *arXiv:2310.11453*.
- [16] X. Qiu *et al.*, “Quantized Spike-driven Transformer,” in *Proc. Int. Conf. Learn. Represent.*, Singapore, 2025.
- [17] H. Cao *et al.*, “Binary Event-Driven Spiking Transformer,” 2025, *arXiv:2501.05904*.
- [18] Y. Liang *et al.*, “Towards Accurate Binary Spiking Neural Networks: Learning with Adaptive Gradient Modulation Mechanism,” in *Proc. AAAI Conf. Artif. Intell.*, Philadelphia, Pennsylvania, USA, 2025, vol. 39, no. 2, pp. 1402–1410.
- [19] G. Perea, M. Navarrete and A. Araque, “Tripartite synapses: astrocytes process and control synaptic information,” *Trends. Neurosci.*, vol. 32, no. 8, pp. 421–431, 2009.
- [20] K. Qian *et al.*, “Revisiting the critical roles of reactive astrocytes in neurodegeneration,” *Mol. Psychiatry*, vol. 28, pp. 2697–2706, 2023.
- [21] K. Seo *et al.*, “Astrocytic inhibition of lateral septal neurons promotes diverse stress responses,” *Nat. Commun.*, vol. 15, no. 10091, 2024.
- [22] J. Lee *et al.*, “Astrocytes phagocytose adult hippocampal synapses for circuit homeostasis,” *Nature*, vol. 590, pp. 612–617, 2021.
- [23] H. Lee, M. A. Wheeler and F. J. Quintana, “Function and therapeutic value of astrocytes in neurological diseases,” *Nat. Rev. Drug Discov.*, vol. 21, pp. 339–358, 2022.
- [24] N. J. Allen, “Astrocyte glypicans 4 and 6 promote formation of excitatory synapses via GluA1 AMPA receptors,” *Nature*, vol. 486, pp. 410–414, 2012.
- [25] S. Zhou, Y. Wu, Z. Ni, X. Zhou, H. Wen and Y. Zou, “DoReFa-Net: Training Low Bitwidth Convolutional Neural Networks with Low Bitwidth Gradients,” 2018, *arXiv:1606.06160*.
- [26] Z. Liu, B. Wu, W. Luo, X. Yang, W. Liu and K. Cheng, “Bi-real net: Enhancing the performance of 1-bit cnns with improved representational capability and advanced training algorithm,” in *Proc. Eur. Conf. Comput. Vis.*, Munich, Germany, 2018, pp. 722–737. DOI: 10.1007/978-3-030-01267-0_44.
- [27] Z. Liu, Z. Shen, M. Savvides and K. Cheng, “Reactnet: Towards precise binary neural network with generalized activation functions,” in *Proc. Eur. Conf. Comput. Vis.*, Glasgow, UK, 2020, pp. 143–159. DOI: 10.1007/978-3-030-58568-6_9.
- [28] S. Ma *et al.*, “The Era of 1-bit LLMs: All Large Language Models are in 1.58 Bits,” 2024, *arXiv:2402.17764*.
- [29] W. Wei *et al.*, “QP-SNN: Quantized and Pruned Spiking Neural Networks,” in *Proc. Int. Conf. Learn. Represent.*, Singapore, 2025.
- [30] A. Flores-Valle, I. Vishniakou and J. D. Seelig, “Dynamics of glia and neurons regulate homeostatic rest, sleep and feeding behavior in *Drosophila*,” *Nat. Neurosci.*, vol. 28, pp. 1226–1240, 2025.
- [31] A. Denizot, M. Arizono, U. V. Nägerl, H. Soula and H. Berry, “Simulation of calcium signaling in fine astrocytic processes: Effect of spatial properties on spontaneous activity,” *PLoS Comput. Biol.*, vol. 15, no. 8, pp. 1–33, 2019.
- [32] G. O. Mizuno *et al.*, “Aberrant Calcium Signaling in Astrocytes Inhibits Neuronal Excitability in a Human Down Syndrome Stem Cell Model,” *Cell Rep.*, vol. 24, no. 2, pp. 355–365, 2018.
- [33] J. Deng, W. Dong, R. Socher, L. Li, K. Li and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Miami, FL, USA, 2009, pp. 248–255.
- [34] A. Krizhevsky, “Learning multiple layers of features from tiny images,” Toronto, Canada, 2009.
- [35] H. Li, H. Liu, X. Ji, G. Li and L. Shi, “Cifar10-dvs: an event-stream dataset for object classification,” *Front. Neurosci.*, vol. 11, 2017.
- [36] A. Amir *et al.*, “A Low Power, Fully Event-Based Gesture Recognition System,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Honolulu, HI, USA, 2017, pp. 7388–7397.
- [37] B. Cramer, Y. Stradmann, J. Schemmel and F. Zenke, “The Heidelberg Spiking Data Sets for the Systematic Evaluation of Spiking Neural Networks,” *IEEE Trans. Neural Netw. Learning Syst.*, vol. 33, no. 7, pp. 2744–2757, 2022.
- [38] Z. Zhou *et al.*, “Spikformer: When Spiking Neural Network Meets Transformer,” in *Proc. Int. Conf. Learn. Represent.*, Kigali, Rwanda, 2023.
- [39] Z. Tu, X. Chen, P. Ren and Y. Wang, “AdaBin: Improving Binary Neural Networks with Adaptive Binary Sets,” in *Proc. Eur. Conf. Comput. Vis.*, Tel Aviv, Israel, 2022, pp. 379–395. DOI: 10.1007/978-3-031-20083-0_23.
- [40] M. Yao *et al.*, “Spike-driven Transformer,” in *Proc. Adv. Neural Inform. Process. Syst.*, New Orleans, LA, USA, 2023, vol. 36.
- [41] C. Zhou *et al.*, “Enhancing the Performance of Transformer-based Spiking Neural Networks by Improved Downsampling with Precise Gradient Backpropagation,” 2023, *arXiv:2305.05954*.
- [42] C. Zhou *et al.*, “QKFormer: Hierarchical Spiking Transformer using Q-K Attention,” in *Proc. Adv. Neural Inform. Process. Syst.*, Vancouver, BC, Canada, 2024, vol. 37, pp. 13074–13098.
- [43] H. Qin *et al.*, “Forward and Backward Information Retention for Accurate Binary Neural Networks,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recog.*, Seattle, WA, USA, 2020, pp. 2247–2256.
- [44] M. Horowitz, “1.1 Computing’s energy problem (and what we can do about it),” in *Proc. IEEE Int. Solid-State Circuits Conf. (ISSCC)*, San Francisco, CA, USA, 2014, pp. 10–14.
- [45] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier and Y. Tian, “Deep residual learning in spiking neural networks,” in *Proc. Adv. Neural Inform. Process. Syst.*, Virtual, 2021, vol. 34, pp. 21056–21069.
- [46] F. Ottati *et al.*, “Astrocytic inhibition of lateral septal neurons promotes diverse stress responses,” *Nat. Commun.*, vol. 13, no. 4, pp. 1015–1025, 2023.
- [47] S. Shen, D. Zhao, G. Shen and Y. Zeng, “TIM: an efficient temporal interaction module for spiking transformer,” in *Proc. IJCAI*, Jeju, South Korea, 2024, pp. 3133–3141.

VI. SUPPLEMENTARY MATERIAL

A. Efficiency of Biological Nervous Systems

We provide quantitative support for our statement about the greater efficiency of biological nervous systems by presenting precise data. The human brain, operating at approximately 20 watts, sustains about 8.6×10^6 neurons and over 1×10^{14} synapses [4]. SNNs, even when deployed on existing low-power neuromorphic chips (such as TrueNorth [7], Tianjic [8], and Loihi2 [9]), demonstrate efficiency levels that remain several orders of magnitude below biological nervous systems. As shown in the Table VI-A, ‘#’ represents a number, and ‘#Neurons (synapses)/W’ denotes the number of neurons or synapses sustained per watt.

TABLE V
QUANTITATIVE ANALYSIS OF ENERGY EFFICIENCY

Platform	Power(W)	# Neurons	#Neurons/W	#Synapses	#Synapses/W
TrueNorth	0.065	1.00×10^6	1.54×10^7	2.56×10^6	3.94×10^7
Tianjic	0.95	4.00×10^4	4.21×10^4	1.00×10^7	1.05×10^7
Loihi2	1	1.00×10^6	1.00×10^6	1.20×10^8	1.20×10^8
Brain	20	8.60×10^{10}	4.30×10^9	1.00×10^{14}	5.00×10^{12}

B. Information Entropy Computation

In TaWQ, the quantized weights are 1.58-bit ternary values $\{+1, 0, -1\}$, and the probability of each value can be formulated as follows.

$$f(w_{tri,q}) = \begin{cases} p_p, & w_{tri,q} = +1, \\ p_z, & w_{tri,q} = 0, \\ p_n, & w_{tri,q} = -1, \end{cases} \quad (15)$$

where $w_{tri,q}$ is a quantized weight, and $p_p + p_z + p_n = 1$. We employ information entropy to quantify the information content carried by quantized weights, with the mathematical expression defined as:

$$\mathcal{H}(w_{tri,q}) = -[p_p \ln(p_p) + p_z \ln(p_z) + p_n \ln(p_n)]. \quad (16)$$

Let $p_p^*, p_z^*, p_n^* = \arg \max_{p_p, p_z, p_n} (\mathcal{H}(w_{tri,q}))$. Analytically, it can be derived that the quantized weights achieve maximum information entropy (1.0986 nats) when $p_p^* = p_z^* = p_n^* = \frac{1}{3}$, satisfying Shannon's entropy maximization principle under a uniform probability distribution.

C. Convergence Analysis of TaWQ Backpropagation

Since the step function is non-differentiable, we use a surrogate gradient during backpropagation. Boundedness of the weight gradient is a necessary condition for network convergence. We demonstrate that each component of Eq. (7) in the Main Manuscript is bounded.

1) $\frac{\partial L}{\partial \mathbf{W}_{tri,q}[t]}$ is computed from the preceding layer's spiking feature maps. Given finite size and the fact that spiking feature maps contain only values of 0 or 1, so $\frac{\partial L}{\partial \mathbf{W}_{tri,q}[t]}$ is bounded.

2) $\frac{\partial \mathbf{W}_{tri,q}[t]}{\partial \mathbf{C}_s[t]} < 1$, as determined by Eq. (8) in the Main Manuscript.

3) $\frac{\partial \mathbf{C}_s[t]}{\partial \mathbf{I}_n} = \lambda$.

4) $\frac{\partial \mathbf{C}_s[t+1]}{\partial \mathbf{C}_s[t]} = 1 - \lambda$ if $\mathbf{W}_{tri,q}[t] = 0$, otherwise $\frac{\partial \mathbf{C}_s[t+1]}{\partial \mathbf{C}_s[t]} = 0$.

5) $\frac{\partial \mathbf{C}_s[t+1]}{\partial \mathbf{W}_{tri,q}[t]} = (1 - \lambda) \times \mathbf{C}_s[t]$ if $\mathbf{W}_{tri,q}[t] < 0$, and $\frac{\partial \mathbf{C}_s[t+1]}{\partial \mathbf{W}_{tri,q}[t]} = (\lambda - 1) \times \mathbf{C}_s[t]$ if $\mathbf{W}_{tri,q}[t] > 0$.

if \mathbf{C}_s is bounded and the time steps T is finite, then the weight gradient remains bounded. \mathbf{C}_s is derived from \mathbf{I} . We initialize \mathbf{I} using the Kaiming Uniform distribution, ensuring its initial values are bounded. During training, we employ gradient clipping to constrain the magnitude of weight updates. Consequently, after a finite number of training steps, \mathbf{I} remains bounded. This approach guarantees the boundedness of \mathbf{C}_s .

Based on the analysis above, we conclude that when the time steps T is finite, the weight gradient remains bounded. This satisfies the necessary condition for network convergence.

D. Consistent Computational Complexity.

The computational complexity of the TaWQ is the same as that of the matrix product using time-invariant weights, the latter is employed in existing quantized SNNs [13], [17]. Assuming each product between matrix $\mathbf{W}_{tri,q}[t]$ and $\mathbf{X}_i[t]$ has a complexity of $O(N)$, which of the TaWQ is $O(TN)$. If it is a general matrix product, the time-invariant weight \mathbf{W} multiplies temporal inputs $\mathbf{X}_i[t]$ across T timesteps, resulting in a complexity of $O(TN)$, which aligns with the TaWQ. Without introducing additional computational complexity, the TaWQ enhances the representational capacity of quantized SNNs by dynamically assigning distinct weights $\mathbf{W}_{tri,q}[t]$ to $\mathbf{X}_i[t]$, thereby fully exploiting temporal information.

E. Multi-bit TaWQ

After the 1.58-bit TaWQ, we conduct further in-depth exploration and develop its multi-bit variant, termed mTaWQ (multi-bit Temporal-adaptive Weight Quantization). The dynamics of mTaWQ are defined as follows:

$$\begin{cases} \mathbf{C}_s[t+1] = \lambda \cdot \mathbf{C}_s[t](1 - |\mathbf{W}_{tri,q}[t]|/n) + (1 - \lambda) \cdot \mathbf{I}_n, \\ \mathbf{W}_{tri,q}[t+1] = \mathcal{Q}(\mathbf{C}_s[t+1], n). \end{cases} \quad (17)$$

Here, n is a positive integer, where $+n$ represents the maximum quantized weight value and $-n$ is the minimum, that is, the quantized weights are constrained to $\{-n, -(n-1), \dots, 0, \dots, +(n-1), +n\}$. When $n = 1$ and $n = 2$, the weights are quantized to 1.58-bit and 2.32-bit. With $n = 4$, the quantization yields 3.17-bit. And for $n = 8$, the weights are quantized to 4.09-bit. $\mathcal{Q}(\mathbf{C}_s, n)$ is the quantization function, and we define it as:

$$\mathcal{Q}(\mathbf{C}_s, n) = \text{round}(\text{clamp}(\mathbf{C}_s, -n, +n)). \quad (18)$$

When $n = 1$, $\mathcal{Q}(\mathbf{C}_s, n)$ is functionally equivalent to $\mathcal{S}(\mathbf{C}_s, C_{th})$ with $C_{th} = 0.5$. And the surrogate gradient $\mathcal{Q}'(\mathbf{C}_s, n) = 1$ if $-n < \mathbf{C}_s < n$, otherwise $\mathcal{Q}'(\mathbf{C}_s, n) = 0$.

The ablation study results of bit-width is presented in Section VI-I6.

F. Datasets

We conduct extensive experiments on static image datasets (ImageNet [33], CIFAR-10/100 [34]), neuromorphic datasets (CIFAR10-DVS [35], DVS128-Gesture [36]) and speech dataset SHD [37] to comprehensively evaluate the proposed TaWQ.

ImageNet is the most widely used large-scale benchmark dataset in image classification, comprising 1,000 classes, with a training set of over 1.2 million images and a validation set of 50,000 images. The CIFAR10 dataset contains 50,000 training images and 10,000 test images, while CIFAR100 shares the same total image count (50,000 for training and 10,000 for test) but differs in class granularity: CIFAR10 has 10 classes, whereas CIFAR100 includes 100 classes. Both CIFAR datasets' resolution is 32x32.

CIFAR10-DVS is a neuromorphic dataset converted by sampling static images through a Dynamic Vision Sensor (DVS) camera, containing 10 classes and 10,000 samples with a 9:1 train-test split ratio. DVS128-Gesture is another neuromorphic dataset for gesture recognition, comprising 1,342 samples across 11 gesture classes, collected from 29 individuals under three distinct illumination conditions.

SHD is an audio-based classification benchmark designed for evaluating spiking neural networks. It consists of spoken digits (0-9) in both German and English, covering 20 distinct classes. The dataset contains 10,000 samples, split into 8156 for training and 2264 for testing. The number of channels is 700.

G. Experimental Setups

The experiments employ LIF neurons with a firing threshold of 1.0, a reset membrane potential of 0, and a membrane time

constant of 2.0. These configurations are consistent with many established works, such as Spikformer [38], and represent commonly adopted settings in the field. For experiments on ImageNet, we employ 8 Ascend 910C NPUs, and other datasets are conducted using 1 Ascend 910C NPU.

Experimental Setup on ImageNet. In this experiment, we employ TaWQ to quantize full-precision models with a batch size of 512 and the AdamW optimizer. The base learning rate (lr_{base}) is set to 6×10^{-4} , and the actual learning rate is calculated as $lr_{\text{base}} \times \text{BatchSize}/256$, resulting in 1.2×10^{-3} , and employ CosineAnnealing Scheduler. Weight decay is 0.05. The training epochs for Spikingformer-TaWQ and QKFormer-TaWQ were set to 310 and 200, respectively, aligning with their full-precision counterparts. Warmup epochs on ImageNet are 5. Additionally, data augmentation techniques are also aligned with their full-precision counterparts.

Experimental Setup on CIFAR. The batch size is set to 32, the timesteps are 4, and the learning rate is set to 2×10^{-3} with AdamW optimizer and CosineAnnealing Scheduler. Weight decay is 0.05. We train models for 400 epochs from scratch, with a warmup for 20 epochs.

Experimental Setup on Neuromorphic Datasets Both batch size and timesteps are uniformly set to 16, with learning rates configured as 5×10^{-3} and 2×10^{-3} for CIFAR10-DVS and DVS128-Gesture, respectively, and training 106 and 200 epochs from scratch. Optimizer is AdamW with Cosine Annealing Scheduler, the weight decay is 0.06.

H. Energy Consumption

1) *Theoretical Energy Consumption:* Comparative analysis of energy consumption in SNNs under the same conditions is critical. Numerous research in the SNN community have adopted energy metrics measured by Horowitz et al. [44] on a 45nm hardware platform, where a single 32-bit multiply-accumulate (MAC) operation consumes $E_{MAC} = 4.6$ pJ, comprising 3.7 pJ for multiplication and $E_{AC} = 0.9$ pJ, where the E_{AC} represents the energy consumption of an accumulate (AC) operation. And for a single 8-bit accumulate, $E_{AC} = 0.03$ pJ.

Estimating the energy consumption of SNNs requires first calculating the synaptic operations (SOPs) of neurons, which can be described as follows.

$$\text{SOPs}^i = \sum_{t=1}^T fr_t^i \times \text{OPs}_t^i, \quad (19)$$

where fr_t^i and OPs_t^i are the firing rate and operations of layer i at timestep t , respectively, T is the timesteps in SNNs. OPs comprise binary operations (BOPs) and floating-point operations (FLOPs). Following ReActNet [27], the OPs is defined as:

$$\text{OPs}_t^i = \text{BOPs}_t^i/64 + \text{FLOPs}_t^i. \quad (20)$$

Notably, the weights quantized by TaWQ are 1.58-bit ternary values $\{+1, 0, -1\}$. We define a synapse ratio sr_t^i for layer i at timestep t , which refers to the proportion of $+1$ (excitatory synapses) and -1 (inhibitory synapses) relative to the total weight count. This allows the conversion of ternary

operations (TOPs) to equivalent binary operations (BOPs) via the relation $\text{BOPs}_t^i = sr_t^i \times \text{TOPs}_t^i$. In the non-quantized layer, $\text{TOPs}_t^i = 0$, leading to $\text{BOPs}_t^i = 0$, and $\text{OPs}_t^i = \text{FLOPs}_t^i$. In the TaWQ-quantized layer, sr_t^i is the actual synapse ratio with $\text{FLOPs}_t^i = 0$ and $\text{OPs}_t^i = \text{TOPs}_t^i/64 \times sr_t^i$.

The TaWQ-quantized SNNs' energy consumption is as follows:

$$E_{\text{quant}} = E_{MAC} \cdot \sum_{l=1}^L \text{FLOPs}_{float}^l + E_{AC} \cdot \sum_{n=1}^N \text{SOPs}^n. \quad (21)$$

Here, N represents the number of TaWQ-quantized layers, and L is the number of layers with floating-point MAC.

2) *Energy Consumption with Residual Integer Spikes:* QKFormer adopts the same shortcuts as SEWResNet [45], which result in the summation of spikes producing non-binary values (integers > 1). In the Main Manuscript, operations involving non-binary values and floating-point weights are treated as multiple floating-point additions followed by summation, which is consistent with the original QKFormer paper. However, non-binary values should properly be regarded as floating-point values, and their operations with floating-point weights should be interpreted as floating-point multiplication rather than addition. This makes the E_{AC} inapplicable and necessitates the use of E_{MAC} instead.

We recalculate the energy consumption for QKFormer. A key difference is that TaWQ-quantized weights take values in $+1, 0, -1$, the computations between > 1 integers arising from shortcuts, and these quantized weights still involve additions. Therefore, the use of E_{AC} in the calculation remains necessary. If we treat > 1 integers as full-precision floating-point values, $E_{AC} = 0.9$ pJ. Conversely, if we consider them as 8-bit integers, E_{AC} becomes 0.03 pJ.

We compute the energy consumption for both scenarios on ImageNet, as shown in the Table VI, where "Power_add" refers to the original calculation method used in the paper, which interprets > 1 integers as the sum of multiple spikes. "Power_int" represents treating > 1 integers as 8-bit integers. Conversely, "Power_fp" denotes treating them as full-precision floating-point values. For non-quantized models (where weights remain full-precision), even when > 1 integers are processed as 8-bit integers, the underlying operations still involve floating-point multiplication. Consequently, only the "Power_fp" outcome is achievable for non-quantized models in such cases.

The table data reveal that, when interpreting > 1 integers as 8-bit integers, the full-precision models consume 32.05mJ and 74.12mJ. In contrast, the TaWQ-based models achieve significantly lower energy consumption, at 0.54mJ and 1.52mJ, representing reductions of 98.32% and 97.95%, respectively. Even when interpreting > 1 integers as full-precision floating-point numbers, the TaWQ-based models still exhibit substantial advantages in energy consumption, consuming only 6.00mJ and 24.37mJ. This translates to reductions of 81.28% and 67.12%, respectively. Thus, these results conclusively demonstrate that TaWQ effectively and substantially reduces the energy consumption of SNNs.

3) *Energy Consumption with Reading/writing Overhead:* Following [46], we have recalculated the energy consumption,

TABLE VI
RECALCULATED ENERGY CONSUMPTION OF QKFORMER AND QKFORMER-TAWQ WITH $T = 4$.

Method	Architecture	Bits	Power_add(mJ)	Power_int(mJ)	Power_fp(mJ)
QKFormer [42]	QKFormer-10-384	32	15.13	-	32.05
	QKFormer-10-768	32	38.91	-	74.12
QKFormer-TaWQ	QKFormer-10-384	1.58	0.45	0.54	6.00
	QKFormer-10-768	1.58	0.99	1.52	24.37

TABLE VII
ENERGY CONSUMPTION RESULTS ON HARDWARE WITH $T = 4$.

Method	Architecture	Bits	Hardware Bits	Power (mJ)	Power_rd (mJ)
Spikingformer [41]	Spikingformer-8-768	32	8	0.64	0.14
Spikingformer-TaWQ		1.58	2	0.36	0.048

with both weights and intermediate states quantized to 8-bit, while spike events are encoded and packed into 8-bit memory words during reading/writing operations (8 spikes per word). We first extend the per-kernel computation described in [46] to multiple kernels, resulting in $N_{rd} = C_o \times C_i \times k_h \times k_w$, where N_{rd} is the number of weights. In the single-timestep scenario, we quantize weights to 1.58-bit via TaWQ while storing each weight with 2 bits. Extending [46]’s spike-encoding methodology, the total reading energy becomes $E_{rd_{tot}} = N_{rd} \times (E_{rd}/4 + E_{rd}/8)$, which is equal to $E_{wr_{tot}}$ for a writing operation. For the input encoding layer (first layer), where both weights and feature maps remain 8-bit, this yields $E_{rd_{tot}} = N_{rd} \times (E_{rd}/1 + E_{rd}/1)$. For non-TaWQ quantized models with 8-bit weights, $E_{rd_{tot}} = N_{rd} \times (E_{rd}/1 + E_{rd}/8)$.

We subsequently extend these single-timestep calculations to T timesteps, where weight-reading operations are executed T times. Finally, we incorporate the spatial dimensions of the feature map by including height H and width W in the calculations. Additionally, energy consumption in the neuron is computed consistently using 8-bit precision. The specific energy consumption results for Spikingformer models on ImageNet are presented in the Table VII.

It can be observed that the energy consumption of the 8-bit Spikingformer is 0.64mJ, while that of Spikingformer-TaWQ is 0.36mJ, which is just 56.25% of the former. We additionally calculated the energy consumption specifically for weight reading (Power_rd) separately. The weight reading energy consumption for Spikingformer is 0.14mJ, while for Spikingformer-TaWQ it is only 34.29% of the former’s consumption, that is 0.048mJ.

In summary, when accounting for both reading/writing energy and internal neuronal computation energy, TaWQ-quantized models demonstrate pronounced energy advantages over their 8-bit counterparts, despite introducing additional temporal dimension during weight quantization.

I. More Experimental results

1) *Statistics of Firing Rates*: We examine the firing rates of Q, K, and V in the last two spiking self-attention (SSA) modules of both QKFormer and QKFormer-TaWQ trained on CIFAR100, as illustrated in Fig. 8. The QKFormer-TaWQ demonstrates a firing rate trend consistent with the QKFormer,

with the highest firing rate in Q, followed by K, and the lowest in V. The Pearson correlation coefficient of the firing rates between QKFormer and QKFormer-TaWQ reaches 0.9733, indicating a strong correlation before and after quantization. Additionally, the mean firing rate after TaWQ quantization is 0.1445, showing a minor deviation from the pre-quantization value of 0.1283. In summary, the TaWQ-quantized model effectively captures the firing rate characteristics of the full-precision model.

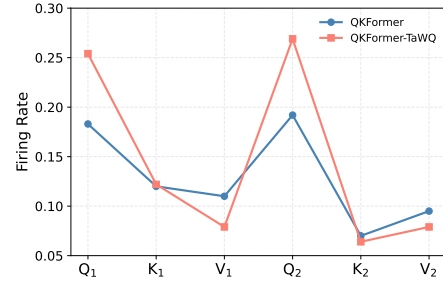


Fig. 8. Firing rates of full-precision and TaWQ-quantized QKFormer.

2) *Latency and Memory Footprint*: We further evaluate the latency and memory footprint of the TaWQ quantization method during inference at a single timestep. While the full-precision and quantized model weights are in 32-bit and 1.58-bit, respectively, the inference is conducted on an Ascend 910 NPU with FP16 floating-point precision due to hardware constraints, meaning both the 32-bit and 1.58-bit weights were converted to 16-bit for execution. The batch size is set to 32. The results in Table VIII indicate that under identical 16-bit conditions, the memory footprint remains the same at 6.88 G, and the latency is nearly identical, with a marginal difference of only 0.2 ms. Significantly, TaWQ-quantized networks exclusively employ addition operations with +1 and -1, therefore, the quantized weights could potentially largely reduce computational demands on hardware specifically designed for low-bit computing.

3) *Comparisons with Post-training Quantizations*: A comparative analysis is conducted between the 1.58-bit TaWQ-quantized models and their 8-bit quantized counterparts implemented using the Post-Training Quantization (PTQ) method on the ImageNet dataset, with $T = 4$. Since Spikingformer and QKFormer only release their largest trained models, we

TABLE VIII
LATENCY AND MEMORY FOOTPRINT OF QKFORMER AND QKFORMER-TAWQ ON A SINGLE TIMESTEP.

Method	Architecture	Bits	NPU Bits	Batch	Latency(ms)	Memory(G)
QKFormer [42]	QKFormer-10-768	32	16	32	201.5	6.88
QKFormer-TaWQ		1.58			201.7	6.88

quantized Spikingformer-8-768 and QKFormer-10-768 using TaWQ for comparison. The results are shown in the Table IX. The results demonstrate that our 1.58-bit Spikingformer-TaWQ achieves 77.42% accuracy, experiencing only a 0.22% accuracy drop compared to the full-precision Spikingformer, whereas the 8-bit Spikingformer-PTQ exhibits an accuracy drop of 2.38%. Meanwhile, QKFormer-10-768-TaWQ achieves 82.94% accuracy, surpassing the PTQ-quantized 8-bit model by 0.96%.

TABLE IX
COMPARISON RESULTS OF TAWQ-QUANTIZED AND PTQ (POST-TRAINING QUANTIZATION) MODELS.

Method	Architecture	Bits	Acc(%)
Spikingformer [41]	Spikingformer-8-768	32	77.64
Spikingformer-TaWQ		1.58	77.42(-0.22)
Spikingformer-PTQ		8	75.26(-2.38)
QKFormer [42]	QKFormer-10-768	32	84.22
QKFormer-TaWQ		1.58	82.94(-1.28)
QKFormer-PTQ		8	81.98(-2.24)

4) *Results of Non-Transformer Spiking Networks:* We train TaWQ-quantized SEWResNet on the ImageNet dataset using the same epochs, learning rate, batch size, and other hyperparameter settings as those specified in [45]. We present the result of the 1.58-bit SEW-ResNet18-TaWQ in Table X. It achieves a final accuracy of 62.06% on ImageNet, surpassing the 8-bit QP-SNN by 0.70%. Compared with the full-precision SEW-ResNet18, this exhibits only a 1.16% accuracy degradation.

TABLE X
RESULTS OF NON-TRANSFORMER STRUCTURE ON IMAGENET.

Method	Architecture	Bits	Acc(%)
XNOR-Net [14]	ResNet-18	1	51.2
Bi-Real Net [26]	Bi-Real-18	1	56.4
QP-SNN [29]	ResNet-18	8	61.36
SEWResNet	SEWResNet-18	32	63.22
SEWResNet-TaWQ		1.58	62.06(-1.16)

5) *Results on SHD classification:* We further validate performance on speech classification tasks using the SHD dataset with a larger timesteps. Input data to the network takes the form $B \times T \times C$, where B is the batch size, T is the temporal dimension, and C is the channel dimension. Following the methodology of [47], the input is processed using spatio-temporal bins, which reduces its dimensionality to $B \times 100 \times 140$. Specifically, $T = 100$ and $C = 140$. We utilize the AdamW optimizer with a Cosine Annealing scheduler. The learning rate is $1e-3$, and the batch size is 32. The experimental results on the SHD dataset are presented in Table XI, where

only minor performance degradation is observed. Specifically, the TaWQ-quantized Spikformer exhibits an accuracy drop of merely 0.35% compared to its full-precision counterpart, while the TaWQ-quantized Spikingformer shows a reduction of only 0.22% in accuracy. The first implementation of Spikformer for the SHD dataset is achieved by [47]. Our analysis reveals that, compared to this prior work, our quantization method maintains performance advantages with a smaller size.

TABLE XI
EXPERIMENTAL RESULTS ON THE SHD. "*" DENOTES THE SELF-REIMPLEMENTED MODEL.

Method	Architecture	Acc(%)
Spikformer [47]	Spikformer-2-256	85.1
TIM [47]		86.3
Spikformer [38]	Spikformer-1-128*	91.16
Spikformer-TaWQ		90.81(-0.35)
Spikingformer [41]	Spikingformer-1-128*	91.74
Spikingformer-TaWQ		91.52(-0.22)

6) *Ablation Study: Ablation Study of Bit-width.* We conduct an ablation study into multi-bit TaWQ (mTaWQ) using the CIFAR100 and QKFormer-mTaWQ, with the results shown in Table XII. The 1.58-bit model is quantized using TaWQ, while the higher-bit model is obtained via mTaWQ. The model size decreases as the bit-width reduces, while the accuracy does not show a significant decline. When quantized to 4.09-bit, the model size is 0.91M, achieving an accuracy of 80.97%. At 1.58-bit, the model size is reduced to 42.86% of the 4.09-bit model, with only a 0.12% decrease in accuracy, indicating that the 1.58-bit weights already encapsulate sufficient information.

TABLE XII
BIT-WIDTH ABLATION STUDY RESULTS.

Method	Bits	T	Size(M)	Acc(%)
mTaWQ	4.09	4	0.91	80.97
	3.17	4	0.72	80.61
	2.32	4	0.55	80.85
TaWQ	1.58	4	0.39	80.85

Ablation Study of Timesteps. The performance of SNNs shows a strong dependence on timestep configurations. We conduct a timestep ablation study on the CIFAR-100 dataset. As shown in Fig. 9, the accuracy of QKFormer exhibits an increasing trend with timestep progression, peaking at $T = 4$. Notably, the accuracy of the quantized QKFormer-TaWQ follows an identical trend to QKFormer, indicating that the performance profile remains consistent despite TaWQ's application. By analyzing the accuracy difference, we observe that QKFormer-TaWQ achieves the closest accuracy with

QKFormer at $T = 3$, with a marginal gap of 0.26%, followed by $T = 4$ with a marginal gap of 0.30%. Overall, selecting $T=4$ achieves high quantization performance while maintaining low degradation.

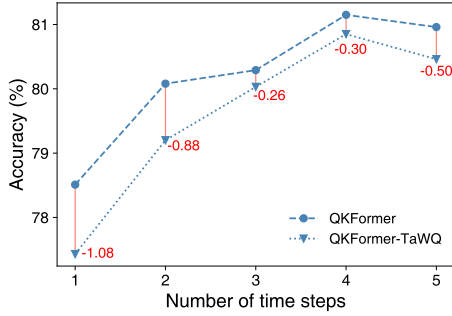


Fig. 9. Timesteps ablation study results of TaWQ-quantized QKFormer.

Ablation Study of Quantization Threshold. The threshold parameter C_{th} in TaWQ affects the performance of quantized SNNs. We conduct C_{th} ablation experiments on the CIFAR-100 dataset using the QKFormer-TaWQ, and the experimental results are shown in Table XIII. We keep the timestep fixed at $T = 4$ and weights quantized to 1.58-bit, only by varying the value of C_{th} , it can be observed that the highest accuracy of 80.85% is achieved at $C_{th} = 0.25$, followed by 80.43% ($C_{th} = 0.15$) and 80.40% ($C_{th} = 0.35$), the accuracy further decreases when $C_{th} = 0.50$, yields 80.34%. These results indicate that $C_{th} = 0.25$ represents the closest optimal value for QKFormer-TaWQ among these C_{th} on the CIFAR-100 benchmark.

TABLE XIII
THE C_{th} ABLATION STUDY RESULTS.

C_{th}	Bits(W-A)	T	Acc(%)
0.15	1.58-1	4	80.43
0.25		4	80.85
0.35		4	80.40
0.50		4	80.34

J. Future work

Our future work will research TaWQ-based quantization in diverse tasks such as detection, segmentation, and language. Additionally, multi-bit quantization variants of TaWQ will be extended to large-scale models. Furthermore, TaWQ-quantized models will be deployed on hardware platforms, including neuromorphic chips and Field Programmable Gate Arrays (FPGAs), to evaluate actual energy consumption and performance under real-world situations.