

Ternary Spike-based Neuromorphic Signal Processing System

Shuai Wang^{a,1}, Dehao Zhang^a, Ammar Belatreche^b, Yichen Xiao^a, Hongyu Qing^a, Wenjie Wei^a, Malu Zhang^{a,*} and Yang Yang^a

^aDepartment of Computer Science and Engineering, University of Electronic Science and Technology of China, Chengdu, 611731, China

^bDepartment of Computer and Information Sciences, Faculty of Engineering and Environment, Northumbria University, Newcastle upon Tyne NE1 8ST, U.K.

ARTICLE INFO

Keywords:

Quantization Spiking Neural Networks
Neural Encoding for Signals
Neuritic Signal Processing
Ternary Spiking Neural Networks
Keyword Spotting and EEG

ABSTRACT

Deep Neural Networks (DNNs) have been successfully implemented across various signal processing fields, resulting in significant enhancements in performance. However, DNNs generally require substantial computational resources, leading to significant economic costs and posing challenges for their deployment on resource-constrained edge devices. In this study, we take advantage of spiking neural networks (SNNs) and quantization technologies to develop an energy-efficient and lightweight neuromorphic signal processing system. Our system is characterized by two principal innovations: a threshold-adaptive encoding (TAE) method and a quantized ternary SNN (QT-SNN). The TAE method can efficiently encode time-varying analog signals into sparse ternary spike trains, thereby reducing energy and memory demands for signal processing. QT-SNN, compatible with ternary spike trains from the TAE method, quantifies both membrane potentials and synaptic weights to reduce memory requirements while maintaining performance. Extensive experiments are conducted on two typical signal-processing tasks: speech and electroencephalogram recognition. The results demonstrate that our neuromorphic signal processing system achieves state-of-the-art (SOTA) performance with a 94% reduced memory requirement. Furthermore, through theoretical energy consumption analysis, our system shows 7.5 \times energy saving compared to other SNN works. The efficiency and efficacy of the proposed system highlight its potential as a promising avenue for energy-efficient signal processing.

1. Introduction

Deep Neural Networks (DNNs) have revolutionized traditional signal processing techniques, achieving higher accuracy in many tasks such as speech recognition (Radford et al., 2023; Liu et al., 2023a; Kim et al., 2023) and electroencephalogram (EEG) analysis (Cai et al., 2019; Su et al., 2022). However, with the development of the Internet of Things and edge computing, the substantial computational and memory requirements of DNNs pose challenges for their direct deployment on edge devices. It limits real-time signal processing and immediate decision-making (Troussard et al., 2023; Giordano et al., 2021; Parekh et al., 2022; Faisal et al., 2019; Valsalan et al., 2020). Thus, devising more lightweight, energy-efficient, and high-performance intelligent signal processing models remains a challenging problem awaiting resolution.

Spiking Neural Networks (SNNs) (Wu et al., 2019; Taherkhani et al., 2020; Bouvier et al., 2019) are inspired by the information transmission mechanisms observed in biological neurons and are considered the third generation of Artificial Neural Networks (ANNs). Spiking neurons compute only upon the arrival of input spikes and remain silent otherwise (Tavanaei et al., 2019). Such event-driven mechanism ensures sparse operations and mitigates extensive floating-point multiplication (MAC) operations in SNNs (Caviglia et al., 2014; Zhang et al., 2021a). Thus, when implemented in hardware, SNNs utilizing accumulate

(AC) operations exhibit a substantially lower power consumption (Orchard et al., 2015) compared to MAC-dependent ANNs. In the past years, the signal processing field has seen many innovative and efficient SNN-based solutions (Cai et al., 2021; Chu et al., 2022; Safa et al., 2021), all achieving satisfactory performance. However, they still face two major bottlenecks.

The first bottleneck lies in the lack of efficient neural encoding methods for signals. Some researchers (Wu et al., 2018; Pan et al., 2020; Xiao et al., 2017) preprocess signals into spectra using Fast Fourier Transform (FFT) and Mel-Frequency Cepstral Coefficients (MFCC), then encode these spectra into spike trains for SNN-based models. However, these preprocessing technologies need massive computing resources, opposing our goal of creating energy-efficient systems. Alternatively, many researchers have adopted the direct encoding method to avoid using preprocessing techniques. This approach directly feeds the raw signal into SNN-based intelligent models, with the model's initial layer as the encoding layer (Weidel & Sheik, 2021; Yang et al., 2022). However, this approach fails to account for signal characteristics and significantly increases MAC operations in the initial model layer. Therefore, devising more efficient neural signal encoding schemes remains a challenge.

In addition, the complexity and memory requirements of SNNs present another significant bottleneck. Recently, an increasing number of high-performance SNNs have been proposed, achieving satisfactory results across many tasks (Zhou et al., 2024; Wang et al., 2023a; Yao et al., 2024; Zhou et al., 2022). However, these models typically exhibit considerable complexity, necessitating extensive computational resources and memory requirements, rendering them unsuitable for deploying on resource-limited edge devices.

*This work was supported by the National Science Foundation of China under Grant 62106038, and in part by the Sichuan Science and Technology Program under Grant 2023YFG0259.

 maluzhang@uestc.edu.cn (M. Zhang)
ORCID(s): 0000-0002-2345-0974 (M. Zhang)

To further exploit the energy efficiency and hardware-friendly advantages of SNNs, many researchers (Hu et al., 2021a; Li et al., 2022; Sulaiman et al., 2020; Castagnetti et al., 2023) have explored quantizing synaptic weights to lower bit-width successfully. However, few works (Yin et al., 2023) focus on the unique membrane potentials within SNNs, leaving substantial room for improvement.

In this research, we present a novel ternary spike-based neuromorphic signal processing system to address the challenges previously mentioned. It primarily comprises two innovative components. First, we propose a threshold-adaptive encoding (TAE) method for effectively encoding raw signals into ternary spike trains and avoiding using high-energy-consuming signal preprocessing techniques. Second, we propose a dual-scale quantization ternary spiking neural network (QT-SNN). It can further process ternary spike signals from our TAE method, using lower bit-width synapse weights and membrane potential. Finally, we validate the performance of our system on two classical signal-processing tasks: keyword recognition and EEG identification. Satisfyingly, compared to other similar efforts, our work achieves state-of-the-art (SOTA) performance with reduced memory usage and lower energy consumption. The principal contributions are summarized as follows:

- We propose an innovative threshold-adaptive encoding (TAE) method. This scheme adaptively adjusts thresholds in response to the time-varying characteristics of raw analog signals. This method permits more sparse and efficient signal encoding into ternary spikes, substantially lowering signal transmission's energy and memory requirements.
- We propose a dual-scaling quantization ternary spiking neural network (QT-SNN). QT-SNN can directly process ternary spike trains from our TAE method, and further quantize synaptic weights and membrane potential in SNNs to lower bit-width, effectively reducing the memory requirements of SNNs with enhanced performance.
- We integrate the TAE method with QT-SNN to develop a ternary spike-based neuromorphic signal processing system. Compared with similar works, our system achieves SOTA performance in both speech and EEG recognition tasks while reducing 94% of memory footprints and 7.5 \times energy consumption.

2. Related Work

In this section, we introduce the latest research findings on neural encoding for signals and spiking neural networks. Simultaneously, we analyze the problems and challenges inherent in existing methods.

2.1. Neural Encoding for Signals

Efficient neural coding methods for signals significantly reduce the bandwidth and memory requirements for information transmission. Furthermore, this process provides

robust data support for subsequent intelligent SNN-based models. Currently, neural encoding for signals can primarily be divided into two types.

Spectrum-based encoding: Given signals' rapid time-variance and high temporal complexity (e.g., speech signals always have 16,000Hz to 48,000Hz sampling rates), directly analyzing and encoding raw analog signals is highly challenging. Consequently, many studies employ MFCC and FFT preprocessing techniques to initially convert analog signals into spectra, and then encode these spectra into spike trains. For example, Xiao et al. (2017) encodes the audio spectrum envelope into spike emission timings through temporal encoding; Pan et al. (2020) considered the masking characteristics of human auditory perception and applied time-frequency masking to process the spectrum, achieving more sparse neural signal encoding method. However, these methods face two significant issues. First, the current signal preprocessing techniques are resource-intensive and lack energy efficiency. Second, these preprocessing methods incur high latency, hindering their application in fast-decision edge devices.

Raw signal-based encoding: To avoid high-resource signal preprocessing techniques, many researchers (Tan et al., 2021; Wu et al., 2019; Taherkhani et al., 2020) are turning to the DC method. It feeds raw analog signals directly into the SNN-based models, utilizing the model's initial layer for encoding. For example, Weidel & Sheik (2021) employed dilated temporal convolution to process raw speech datasets, achieving accuracy in keyword recognition tasks nearly comparable to ANNs. Subsequently, Yang et al. (2022) enhanced accuracy in the same tasks by employing residual convolutional blocks to process raw speech directly. Nonetheless, the DC method leads to extensive MAC operations in the model's primary layer, compromising the energy efficiency of SNNs. Recently, some threshold-based encoding methods for signals have been introduced (Kasabov et al., 2013, 2016; Dupeyroux et al., 2022) to enable real-time, efficient encoding of raw signals into ternary spike sequences, effectively addressing the limitations of the previous methods. However, these methods rely on predetermined thresholds and suffer from significant loss of information during the encoding process. Thus, devising methods for more effective and energy-efficient encoding of signals into spike trains remains a significant challenge.

2.2. Spiking Neural Networks

The event-driven mechanisms ensure SNNs consume minimal energy, offering a significant advantage for energy and compute-constrained edge devices. With the introduction of ANN-SNN (Cao et al., 2015; Stöckl & Maass, 2021; Bu et al., 2022) and STBP (Fang et al., 2021; Zhang et al., 2021b) algorithm, the complexity associated with training high-performance SNNs was significantly reduced. Based on these advanced learning algorithms, (Hu et al., 2021b; Zheng et al., 2021) proposed deep residual SNNs and (Yao et al., 2023b; Zhu et al., 2022) further contributed multi-dimensional spiking attention mechanisms, achieving

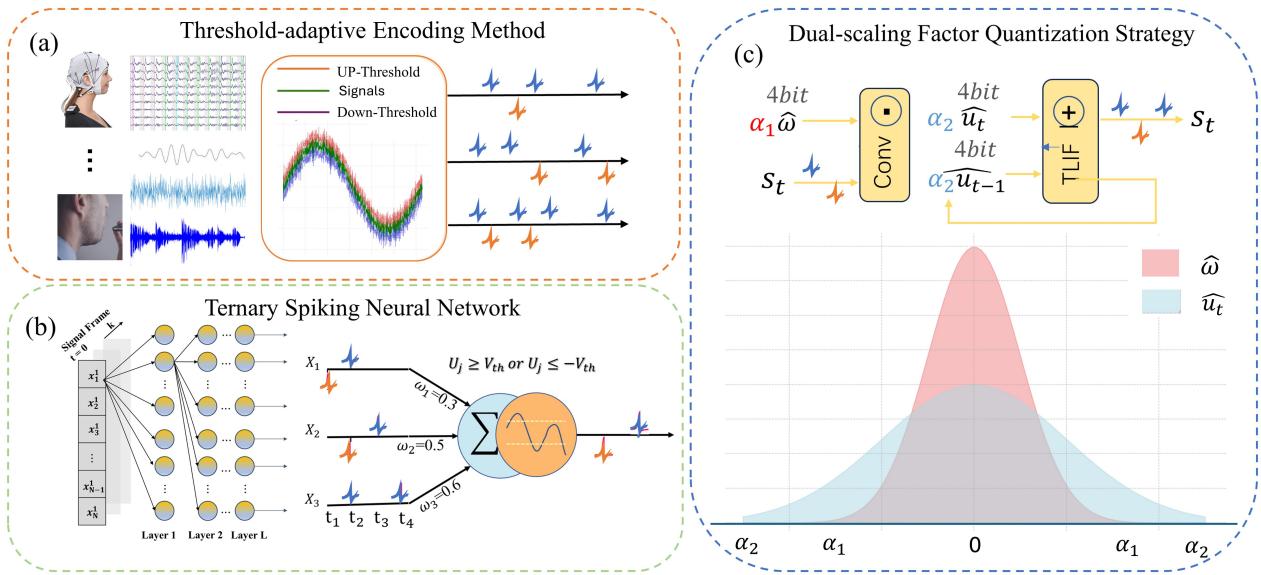


Figure 1: Ternary spike-based neuromorphic signal processing system. (a) Threshold-adaptive Encoding Method. It can efficiently encode analog signal data into ternary spike trains, thereby reducing memory footprints and energy consumption on signal transmission. (b,c) Dual-scaling quantization ternary spiking neural network. It enhances the SNNs' performance and quantizes both synaptic weights and membrane potential to lower bit-width, significantly reducing the network's memory and computational resource requirements.

commendable performance on ImageNet (Deng et al., 2009) and a variety of neuromorphic datasets (Li et al., 2017; Amir et al., 2017). Subsequently, spiking transformers (Zhou et al., 2022; Yao et al., 2024, 2023a) optimized the computation process of spike-based self-attention, further enhancing the performance of SNNs across numerous tasks. Yet, with the incremental enhancement in performance, there is a notable rise in both the model complexity and memory requirements of these networks, which stands in opposition to the edge-friendly objective.

To further reduce SNNs' model size, some researchers try to increase the network's information capacity. Such as PLIF (Fang et al., 2021) utilizes learnable decay to boost neurons' temporal processing ability, enhancing learning and convergence rates; Guo et al. (2023) introduced the ternary spike SNN, which increases the network's information capacity at the neuron scale, achieving enhanced performance with smaller model sizes. However, these methods still employ full-precision weights and membrane potentials, necessitating higher memory requirements. To address this issue, some researchers (Hu et al., 2021a; Li et al., 2022; Sulaiman et al., 2020) have introduced weight quantization into SNNs. For example, ADMM (Deng et al., 2021) optimizes a pre-trained full-precision network for low-precision weight quantization. Chowdhury et al. (2021) utilizes K-means clustering quantization to maintain reasonable accuracy in 5-bit weight SNNs. Subsequently, Yin et al.

(2023) highlights the significance of quantifying membrane potential in diminishing the memory footprint of SNNs, and quantizes both synaptic weights and membrane potentials, further reducing SNNs' memory requirements. However, with the gradual reduction in bit-width for synaptic weights and membrane potentials, the information representation capability of these networks significantly diminishes, leading to notable performance degradation. Therefore, maintaining network performance and significantly reducing the memory and energy requirements of SNNs remains a challenge.

3. Preliminaries

In this section, we present the fundamental principles of threshold-based signal encoding methods. Subsequently, we analyzed traditional binary SNNs' limitations in processing ternary spike trains and introduced the fundamentals and advantages of ternary spike SNNs.

3.1. Threshold-based Encoding for Raw Signals

Although neuromorphic sensors for signals have been developed Chan et al. (2007); Bartolozzi (2018), their widespread implementation remains limited. Consequently, most SNN-based signal-processing tasks still rely on analyzing analog signals captured by traditional sensors. Thus, designing real-time, efficient signal encoding schemes is a major focus. They could transform analog signals from traditional sensors into spike-based signals, significantly

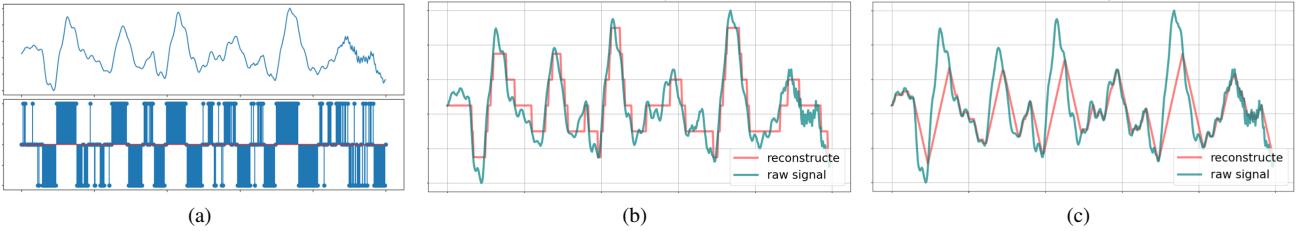


Figure 2: The threshold-based encoding method for raw signals. (a) The threshold-based encoding methods transform raw signals into ternary spike trains consisting of $\{-1, 0, 1\}$. (b) A larger threshold may result in notable reconstruction fluctuations in smooth signal regions. (c) Conversely, a smaller threshold may result in significant peak information loss..

reducing information transmission latency and bandwidth requirements.

Within the current research, three threshold-based encoding algorithms stand out: Threshold-based Representation(TBR), Moving Window(MW), and Step-forward(SF) methods. As depicted in Fig.2(a), the threshold-based encoding methods utilize a fixed threshold to determine the baseline of analog signals, thereby encoding them into ternary spike trains. Among them, SF encoding has achieved significant success in various applications. The SF encoding method calculates the differential of the input signal from a baseline and generates a spike (1 or -1) when the change surpasses a predefined threshold. Notably, the baseline itself is updated according to the spike's polarity. The SF encoding algorithm is described in Algorithm.1.

Algorithm 1 Step-forward encoding (SF)

```

1: Data: input, threshold
2: Result: spikes, init
3:  $L \leftarrow \text{length}(\text{input})$ 
4:  $\text{spikes} \leftarrow \text{zeros}(1, L)$ 
5:  $\text{init}, \text{base} \leftarrow \text{input}(1)$ 
6: for  $i = 2$  to  $L$  do
7:   if  $\text{input}(i) > \text{base} + \text{threshold}$  then
8:      $\text{spikes}(i) \leftarrow 1$ 
9:      $\text{base} \leftarrow \text{base} + \text{threshold}$ 
10:    else if  $\text{input}(i) < \text{base} - \text{threshold}$  then
11:       $\text{spikes}(i) \leftarrow -1$ 
12:       $\text{base} \leftarrow \text{base} - \text{threshold}$ 
13:    end if
14:  end for

```

3.2. Ternary Spiking Neural network

Traditional binary spiking neurons cannot directly process and compute negative spikes in those threshold-based encoding methods. This limits the efficient integration between neural signal encoding and backend signal processing models. Therefore, to accommodate ternary neural encoding and enhance the information capacity of SNNs, we turn to use ternary Leaky Integrate-and-Fire(TLIF) neurons(Guo et al. (2023)). Their neural dynamics can be expressed as follows:

$$u_i^t = \tau u_i^{t-1} \text{Reset}(U_i^{t-1}) + \sum_j w_{ij} o_j^t \quad (1)$$

$$o_i^t = \begin{cases} 1, & \text{if } u_i^t \geq V_{th} \\ -1, & \text{if } u_i^t \leq -V_{th} \\ 0, & \text{otherwise} \end{cases} \quad (2)$$

where τ is the constant leaky factor, u_i^{t+1} is the membrane potential of neuron i at the time step $t + 1$, and $\sum_j w_{ij} o_j^{t+1}$ denotes the pre-synaptic inputs for neuron i . When the membrane potential u_i^{t+1} exceeds the firing threshold V_{th} , the neuron i fires a spike o_i^{t+1} and u_i^{t+1} reset to 0. Eq. 2 describes the firing function and $\text{Reset}(u_i^{t+1}) = (1 - |o_i^{t+1}|)$ is hard reset mechanism. The ternary spike neuron enhances neuronal information capacity while maintaining the benefits of event-driven processing and full-precision floating-point addition.

Firstly, we analyze the increased information capacity of ternary spiking neurons from Shannon's Theorem and information entropy. Information entropy, denoted as $H(X)$, is a metric for the uncertainty or randomness of information in a system, defined by the formula:

$$H(X) = - \sum_{i=1}^n P(x_i) \log_b P(x_i) \quad (3)$$

Where X represents a random variable indicating the state of the system, $P(x_i)$ is the probability of state x_i , n is the number of states, and b is the base of the logarithm, typically 2 for binary systems. For binary SNNs, each neuron can exist in two states (active or inactive). Assuming equal probability for each state ($P(0) = P(1) = 0.5$), the information entropy $H_2 = 1$ bit. In contrast, a ternary spike neuron has three possible states and assumes equal probability for each state ($P(0) = P(1) = P(-1) = \frac{1}{3}$), yields an information entropy $H_3 \approx 1.585$ bits. Comparatively, the ternary spike neuron exhibits an increased information capacity of approximately 0.585 bits per neuron over the binary spike neuron.

Secondly, ternary spike neurons also exhibit characteristics of event-driven computation and full floating-point addition. As depicted in Eq.2, only when the membrane

potential exceeds or falls below a threshold value $\pm V_{\text{th}}$ (-1, 0, and 1), can ternary spike neurons emit spikes; otherwise, they remain silent. During the accumulation of membrane potentials, owing to the ternary input nature of o_j (encompassing 0, 1, and -1), the membrane potential update u_i can be implemented through full floating-point addition, as illustrated in Eq.1. Consequently, the ternary spike neuron model fits ternary neural encoding methods, while retaining the advantages of binary spike SNNs and enhancing the information capacity of SNNs.

4. Method

In this section, we introduce the TAE method and QT-SNN to address the bottlenecks mentioned in SNN-based intelligent signal processing tasks. Specifically, the TAE method can encode analog signals into ternary spike trains with adaptive threshold, exhibiting greater sparsity, and better performance. Meanwhile, QT-SNN tailors to the TAE method, and significantly reduces the energy consumption and memory requirements of backend intelligent signal processing models. Next, we will introduce the details of the TAE method and QT-SNN separately.

4.1. Threshold-Adaptive Encoding

As previously mentioned, those threshold-based signal encoding methods efficiently encode analog signals into ternary spike trains, significantly reducing the bandwidth required for information transmission. Nevertheless, their efficacy largely hinges on the initial *threshold* setting. For example, in the SF algorithm.1, *threshold* is typically set based on extensive experimentation, rather than being dynamically adjusted in response to ongoing signal changes. This static approach may cause some problems. As illustrated in Figs.2(b) and 2(c), a larger threshold can cause notable reconstruction fluctuations in smoother signal regions, and a smaller threshold may result in significant peak information loss. Consequently, a fixed threshold struggles to accommodate the intricacies of time-varying signals fully. To address this limitation, de Gelder (2021) have introduced a group-based SF encoding method. This approach utilizes multiple thresholds to more effectively extract information from both smooth and peak regions of the signal. While this strategy reduces information loss, it also increases computational demands. Therefore, the development of an adaptive threshold mechanism, which dynamically adjusts to the current signal variations, remains a critical challenge in optimizing SF encoding for complex, time-varying signals.

To better address the loss of information due to fixed thresholds, we propose a threshold-adaptive encoding(TAE) method. This scheme adapts the threshold dynamically in response to the current variations in signals, significantly reducing the information loss of encoding in the smoothing and steep phase. The TAE method is shown as Algorithm.2. In our TAE method, *Input* refers to the initial signal input, *Threshold* denotes the initial threshold value, whereas *a* represents the hyperparameter for threshold adaptation, which is typically preset to a value of 1.1. The essence of this

Algorithm 2 Threshold-Adaptive Encoding(TAE)

```

1: Data: input, threshold, a
2: Result: spikes, init
3:  $L \leftarrow \text{length}(\text{input})$ 
4:  $\text{spikes} \leftarrow \text{zeros}(1, L)$ 
5:  $\text{init}, \text{base} \leftarrow \text{input}(1)$ 
6: for  $i = 2$  to  $L$  do
7:   if  $\text{input}(i) \geq \text{base} + \text{threshold}$  then
8:      $\text{spikes}(i) \leftarrow 1$ 
9:      $\text{base} \leftarrow \text{base} + \text{threshold}$ 
10:     $\text{threshold} \leftarrow \text{threshold} \times a$ 
11:   else if  $\text{input}(i) \leq \text{base} - \text{threshold}$  then
12:      $\text{spikes}(i) \leftarrow -1$ 
13:      $\text{base} \leftarrow \text{base} - \text{threshold}$ 
14:      $\text{threshold} \leftarrow \text{threshold} \times a$ 
15:   else
16:      $\text{threshold} \leftarrow \text{threshold}/a$ 
17:   end if
18: end for

```

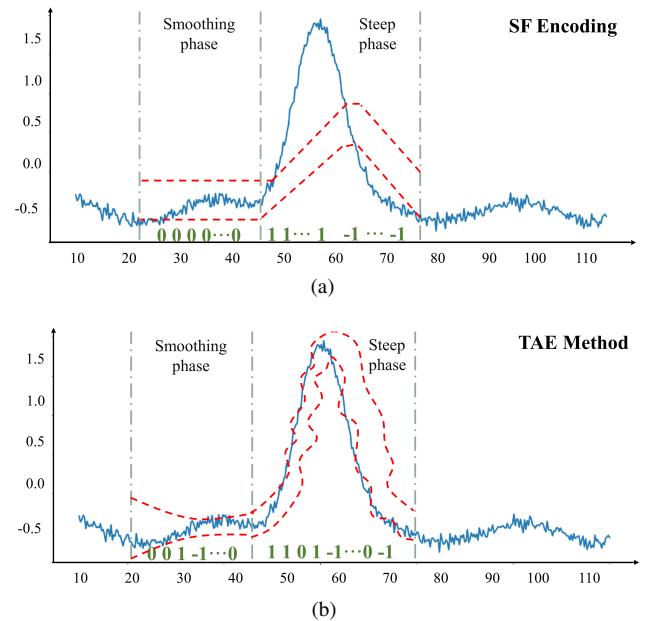


Figure 3: The comparison between SF and TAE methods. (a) In SF coding, the fixed threshold results in a linear increase and decrease of the *base* (indicated in red). (b) The TAE method employs an adaptive thresholding technique, allowing the base to conform to any given curve (red curve).

algorithm lies in its capability to autonomously modulate the threshold size in response to the current dynamic state of the signal. As demonstrated in Fig.3(b), when the signal is in a smoothing phase, TAE gradually decreases the threshold through *a* until the $\text{base} \pm \text{threshold}$ can detect changes in the smoothing signal. Conversely, during peak phases, it increases the threshold by *a* to accommodate the signal's rapid fluctuations. This approach can effectively address the issue of losing smooth information and prominent peaks in threshold-fixed encoding as shown in Fig.3(a),

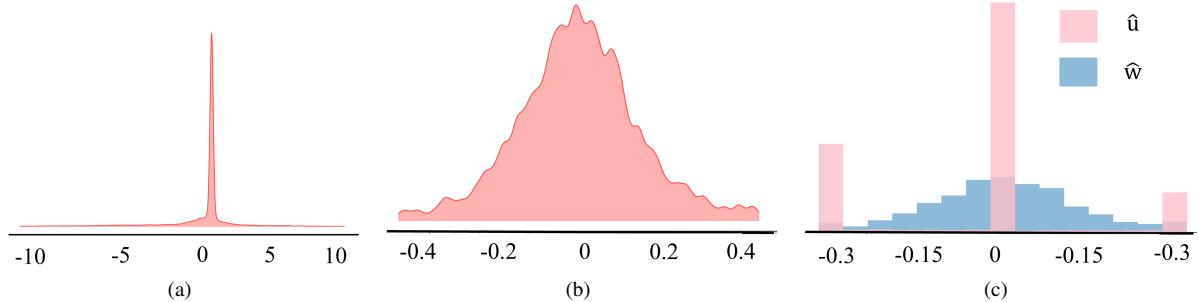


Figure 4: The distribution of weights and membrane potential. (a) In a full-precision ternary spike SNN, the membrane potential and weights follow normal distributions with a mean of zero and differing standard deviations. (b) The MINT method quantizes the membrane potential and weights into 4 bits using nonlinear mapping function \tanh , disrupting the normal distribution characteristics of the membrane potential.

retaining all the advantages of threshold-based neural coding methods. Furthermore, TAE minimizes the occurrence of continuous spikes states in the encoding output, such as $\{-1, -1, \dots, -1\}$ and $\{1, 1, \dots, 1\}$. This results in a lower spiking firing frequency for information transmission. we will demonstrate this point in the experiments of Section 4.1.

4.2. Dual-scaling Factor Quantization Ternary Spiking Neural Networks

Despite ternary spike SNNs capitalize on the energy-efficiency benefits of event-driven processing and can directly process ternary spikes from our TAE method. The membrane potentials and synaptic weights in them remain represented as full-precision floating-point values, as illustrated in Eq.1. This presents two significant issues: firstly, the deployment of models also requires substantial memory footprints. secondly, operations with full-precision floating values are still not energy-efficient. To address these challenges, many researchers Castagnetti et al. (2023); Hwang & Kung (2024) have applied quantization techniques to the synaptic weights to achieve reduced memory consumption and enhanced energy efficiency. Suppose weights are represented by $w \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}}}$, where C_{out} and C_{in} are the number of output and input channels respectively. The quantization of the weights is denoted as:

$$\hat{w} = \Pi_{Q_{\alpha,b}}[w, \alpha] \quad (4)$$

where α is the clipping threshold and the clipping function Π_α clips weights into the interval $[-\alpha, \alpha]$. After clipping, w is projected by $\Pi(\cdot)$ onto the quantization levels \hat{w} . We define $Q(\alpha, b)$ for a set of quantization levels, where b is the bit-width. For uniform quantization, the quantization levels are defined as

$$Q(\alpha, b) = \alpha \times \left\{ 0, \pm \frac{1}{2^{b-1} - 1}, \pm \frac{2}{2^{b-1} - 1}, \dots, \pm 1 \right\}. \quad (5)$$

For every floating-point number, uniform quantization maps it to a b -bit fixed-point representation in $Q_u(\alpha, b)$. α is stored separately as a full-precision floating-point in $Q(\alpha, b)$.

However, few works focus on the unique membrane potential in SNNs. To concurrently quantize membrane potential and synaptic weights within SNNs, further reducing the memory requirements of ternary spike SNNs, we initiate the process by straightforwardly applying Eq.5(UQ) to Eq. 1, employing separate full-precision scaling factors α_1 , α_2 , and α_3 for each integer value. At this stage, the updated equation for membrane voltage can be presented as:

$$\alpha_1 \hat{u}_i^t = \alpha_2 \widehat{\tau u_i^{t-1} \text{Reset}(U_i^{t-1})} + \sum_j \alpha_3 \widehat{w_{ij} o_j^t} \quad (6)$$

However, directly applying Eq.6 for training and inference introduces several MAC operations during the accumulation of membrane potential, compromising the energy efficiency of SNNs. Yin et al. (2023) addressed this by leveraging a nonlinear mapping function to map both u and w to the range $[-1, 1]$, thereby ensuring $\alpha_1 = \alpha_2 = \alpha_3$. This allowed α_1 to be integrated into the threshold for spike emission. Nevertheless, as shown in Fig.4(a) and Fig.4(b), the distributions of u and w indicates a significant discrepancy. Consequently, to better maintain the distribution of membrane potential and synaptic weights, we introduce a dual-scale factor quantization strategy. The membrane potential update formula is presented as follows:

$$\alpha_1 \hat{u}_i^t = \alpha_1 \widehat{\tau u_i^{t-1} \text{Reset}(U_i^{t-1})} + \sum_j \alpha_3 \widehat{w_{ij} o_j^t} \quad (7)$$

In Eq.7, We assume that the membrane potentials across different time steps follow similar normal distributions, thus allowing the use of a uniform scaling factor α_1 . Conversely, synaptic weights are scaled using a distinct factor α_3 . To further enhance the energy efficiency of our QT-SNN in inference, we amalgamate Eq.6 and Eq.2, integrating α_1 and α_3 into the spike emission process. The revised formulation for membrane potential update and spike emission can be described as :

$$\hat{u}_i^t = \widehat{\tau u_i^{t-1} \text{Reset}(U_i^{t-1})} + \sum_j \hat{w}_{ij} o_j^t \quad (8)$$

$$o_i^t = \begin{cases} \frac{a_3}{a_1}, & \text{if } u_i^t \geq \lceil V_{th}/a_1 \rceil \\ -\frac{a_3}{a_1}, & \text{if } u_i^t \leq \lceil -V_{th}/a_1 \rceil \\ 0, & \text{otherwise} \end{cases} \quad (9)$$

During the training phase, as demonstrated in Eq.8, the update of the membrane potential avoids MAC operations. However, during the spike emission process, the intensity of spikes is scaled by the ratio a_3/a_1 , and the threshold is based on V_{th}/a_1 , both of which are parameters subject to learning. To ensure quantization performance and further reduce the energy consumption during the inference phase, we set V_{th}/a_1 as a full-precision learnable parameter. The membrane decay factor τ is set to 2^{-1} , and a_3/a_1 is also defined as a learnable parameter in the form of powers of 2. This configuration allows for the utilization of bit-shift operations for both training and inference phases, with the specific bit-shift formula presented as follows:

$$2^x r = \begin{cases} r \gg x, & \text{if } x > 0 \\ r \ll x, & \text{if } x < 0 \\ r, & \text{otherwise} \end{cases} \quad (10)$$

Subsequently, during the inference phase, the threshold can be calculated using $\lceil V_{th}/a_1 \rceil$, where $\lceil *\rceil$ denotes the ceiling operation; a_3/a_1 can be reintegrated into the membrane voltage accumulation process through bit-shift operations on w , thereby ensuring that spikes are transmitted as $\{-1, 0, 1\}$ throughout the inference process. The specifics of this inference process are described in Algorithm.3.

Our QT-SNN eliminates MAC operations during the inference process, maximally preserving the energy efficiency advantage of SNNs. Additionally, it quantizes both weights and membrane potentials to lower bit-width, reducing the memory required to deploy the model on hardware. During the training phase, treating a_1/a_3 as a learnable spike value effectively enhances the performance of QT-SNN. we will validate these points in the experimental section. Therefore, employing QT-SNN as the backend model for intelligent signal processing tasks markedly improves performance while simultaneously reducing model deployment memory footprints.

5. Experiment

In this section, we systematically validate the performance advantages of our TAE method and QT-SNN through comparative experiments. For TAE method, we assess its spike firing rate and the mean absolute error (MAE) of reconstruction against two other neuromorphic encodings across radar (López-Randulfe et al., 2022) and speech

Algorithm 3 Inference path

```

1: Data:
2: Input spikes  $o_i^t$  to the layer i at time t
3: integer weights  $w_i$  of the layer i
4: leakage factor  $\tau \leftarrow 0.5$ 
5: integer membrane potential  $u_i^{t-1}$  at time t-1
6: integer exponent  $k \leftarrow \log_2 \frac{a_3}{a_1}$ 
7: integer thresholds  $\theta$  of value
8: Result:
9: Output spikes  $o_i^t$  to the layer i at time t
10: integer membrane potential  $u_i^t$  at time t
11:  $X_i^t = \sum_j w_{ij} o_j^t$ 
12:  $H_i^t = X_i^t \gg k + u_i^t \gg 1$ 
13: if  $H_i^t \geq \theta$  then
14:    $o_i^t \leftarrow 1$ 
15:    $u_i^t \leftarrow 0$ 
16: else if  $H_i^t \leq -\theta$  then
17:    $o_i^t \leftarrow 1$ 
18:    $u_i^t \leftarrow 0$ 
19: else
20:    $o_i^t \leftarrow 0$ 
21:    $u_i^t \leftarrow H_i^t$ 
22: end if
```

datasets. Regarding QT-SNN, we examine its performance enhancement over other quantized SNNs on the CIFAR10 datasets. Subsequently, we verify the performance and memory footprints of the proposed ternary spike-based neuromorphic signal processing system on the GSC and EEG datasets. Finally, we analyze the energy efficiency benefits of our system compared to ANNs and other SNNs using DC encoding methods.

5.1. Datasets

GSC Dataset: The Google Speech Commands(GSC) dataset includes 30 short commands for Version 1 (V1) and 35 for Version 2 (V2), recorded by 1,881 and 2,618 speakers, respectively. To make a fair comparison, our experiments are conducted on the 12-class classification and 35-class classification tasks as previous SNN models Yilmaz et al. (2020); Yang et al. (2022); Orchard et al. (2021). While 12-class classification recognizes 12 classes, that include 10 commands: “yes”, “no”, “up”, “down”, “left”, “right”, “on”, “off”, “stop” “go”, and two additional classes: silence, and an unknown class. The unknown class covers the remaining 20 (25) speech commands in the set of 30 (35). The silence class accounting for about 10 % of the total dataset is generated by splicing the noise files in the dataset. Finally, GSC-V1 is split into 56588 training, 7743 validation, and 7835 test utterances, and GSC-V2 is divided into 92843 training, 11003 validation, and 12005 test utterances.

KUL dataset: EEG recordings were obtained from 16 individuals with normal hearing, engaged in the task of concentrating on a specific speaker among two. The auditory stimuli consisted of four tales in Dutch, narrated by three male Flemish speakers. To maintain uniform perceived

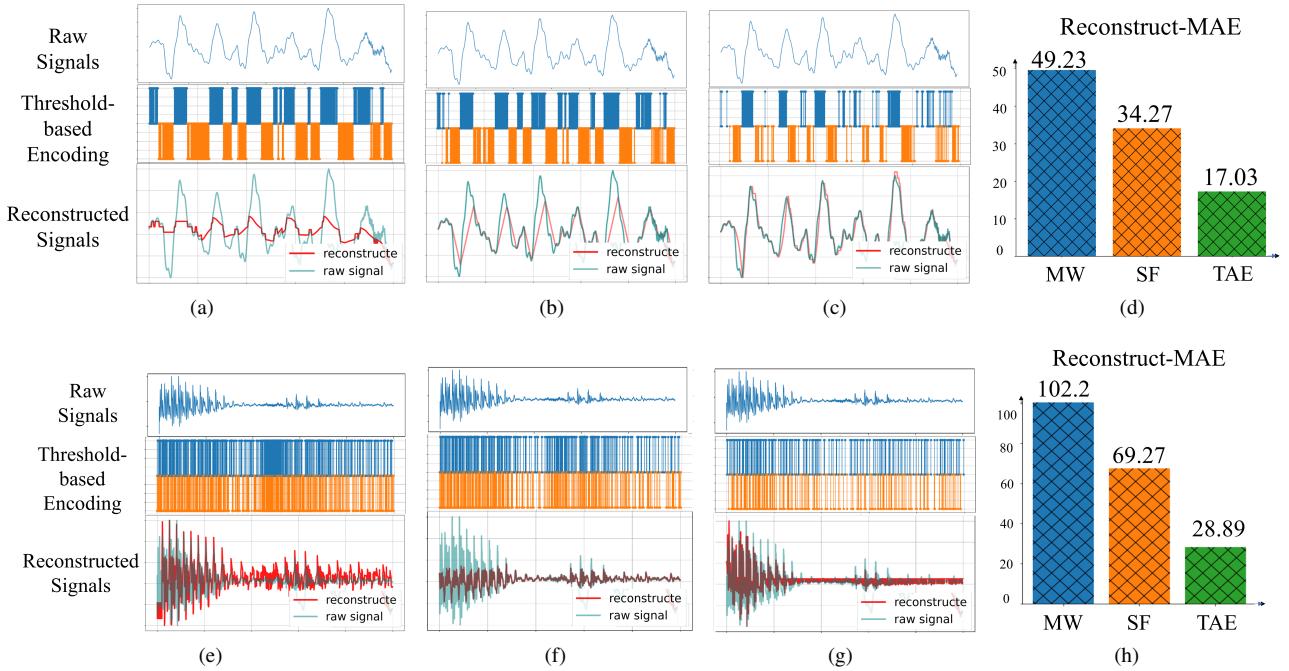


Figure 5: A comparative performance analysis of different encoding methods on radar and GSC datasets. (a-c) depict the encoding and reconstruction capabilities of MW, SF, and our TAE method for a radar signal segment. (e-g) illustrate these encoding strategies applied to a segment of speech signal. (d) and (h) include statistical analysis of the reconstruction-MAE for the overall datasets.

loudness, the intensities of the stimuli were normalized using root mean square (RMS) and presented dichotically (one speaker per ear) or via head-related transfer function filtering, creating an auditory illusion of speech originating from 90° relative to the listener's position. The experiment comprised eight sessions, each lasting six minutes, with the order of conditions randomized for each participant. The EEG data were collected using a 64-channel BioSemi ActiveTwo system at a sampling rate of 8,192 Hz, within a sound-controlled setting. Detailed information on the expanded KUL dataset can be found in the referenced literature (Das et al., 2019; Vandecappelle et al., 2021).

DTU dataset: This dataset comprises EEG recordings from 18 participants with normal hearing, each focusing on one of two speakers in a controlled setting. The auditory stimuli were dialogues between a male and a female speaker. Both native speakers, presented in various acoustic environments. Uniform loudness across the audio streams was achieved through RMS normalization, with the speakers positioned at 0° and 60° azimuths. Each session was designed to last 60 minutes, divided into 512 segments, with both speaker selection and stimuli sequence randomized in each trial. EEG data capture was performed using a 64-channel BioSemi ActiveTwo system, operating at a sampling rate of 512 Hz. Additional details regarding the DTU dataset are available in the cited studies (Fuglsang et al., 2018, 2017).

5.2. Performance of TAE Method

To validate the superiority of our TAE method over alternative threshold-based encoding methods for raw signals,

we conducted further tests focusing on three key metrics: the average mean absolute error (MAE) in signal reconstruction, the average spike firing rate, and the performance of our TAE method. As shown in Fig.5, by counting the average Reconstruct-MAE of the three methods on the GSC and radar () datasets, the TAE method was significantly smaller than the other two methods. The experimental outcomes reveal that the TAE method more closely approximates the original signals and more reliably preserves their critical features. Additionally, as depicted in Fig.6, by conducting

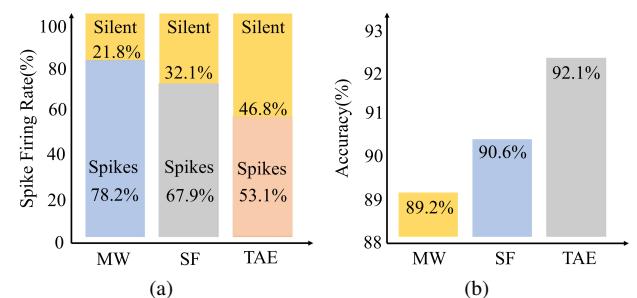


Figure 6: comparing TAE method against other encoding methods on performance and spike firing frequency. (a) Comparing respective accuracies on the GSC dataset among TAE, MW, and SF under a uniform backend classification model. (b) The spike firing rates of TAE, SF, and MW on the GSC dataset.

tests using the same network architecture, we test the average spike firing rates and performance across these

Table 1

Comparison of QT-SNN with other methods

Method	fp32	Precision (W/U)			
		8/8	4/4	2/2	1/32
ResNet-19					
MINT	91.29	91.36	91.45	90.79	-
(Ours)QT-SNN	94.59	94.31	93.99	93.71	-
VGG-16					
TC-SNN	92.68	-	-	-	91.51
MINT	91.15	90.72	90.65	90.56	-
CBP-QSNN	91.79	-	-	-	90.93
(Ours)QT-SNN	94.27	93.88	93.56	93.28	-
VGG-9					
MINT	88.03	87.48	87.37	87.47	-
(Ours)QT-SNN	92.11	91.58	91.12	90.71	-
7Conv+2FC					
CBP-QSNN	89.73	-	-	-	89.01
(Ours)QT-SNN	91.27	90.88	90.90	90.75	-

three encoding schemes on the GSC dataset. The results demonstrate a lower peak firing rate and higher performance for the TAE method, indicating its ability to represent analog signals with better representation and more sparse spikes. Consequently, our TAE method can encode analog signals into ternary spikes more effectively at lower firing rate, significantly reducing the bandwidth requirements for signal transmission.

5.3. Performance of QT-SNN

To assess the performance of our QT-SNN, we compare it against other quantization methods on CIFAR10 with different network structures. As demonstrated in Table.1, In the context of full precision for synaptic weights and membrane potentials, our QT-SNN demonstrates superior performance compared to other methods(MINT(Yin et al. (2023)), TC-SNN(Zhou et al. (2021)), CBP-QSNN(Yoo & Jeong (2023))). This is attributed to the fact that ternary spike neurons, alongside learnable parameters a_1/a_3 , can improve the information capacity, thereby enhancing the performance of QT-SNN. Moreover, our QT-SNN exhibits no significant performance degradation at reduced bit widths for w and u , underscoring the superiority of our quantization approach. Additionally, as illustrated in Fig.7, we analyzed the distribution of membrane potentials and weights across all layers at low bit-width levels, aligning with normal distributions of varying standard deviations. Notably, the proportion of zero values in membrane potentials and weights increases as the bit width decreases, further reducing the computational energy consumption of our QT-SNN.

5.4. Performance of Ternary Spike-based Neuromorphic Signal Processing System

By integrating the strengths of the TAE method and QT-SNN, we design a more lightweight and energy-efficient ternary neuromorphic signal processing system. To validate

the performance and memory efficiency of our system, we establish numerous experiments across two classical signal processing tasks: keyword spotting and EEG recognition. As shown in Tables.2 and 3, our system maintains SOTA performance with both membrane potentials and weights quantized to lower bit-width, outperforming other SNN-based signal processing solutions (Weidel & Sheik, 2021; Zhang et al., 2023; Yin et al., 2021; Stewart et al., 2023; Wang et al., 2023b; Orchard et al., 2021). Regarding memory efficiency, we established ablation experiments to confirm that both the TAE method and QT-SNN contribute to reducing the system's memory footprint. As shown in Fig.8, In scenarios with a batch size of 1, QT-SNN can achieve approximately 90% reduction in memory usage, with negligible differences observed between DC and our TAE methods. However, for larger batch sizes, TAE encoding achieved up to 87% memory reduction compared to the DC method. Consequently, although our TAE method exhibits a marginal performance decrement relative to the DC method, it substantially reduces the system's memory requirements, rendering it more suitable for deployment on resource-limited edge devices. Moreover, we will detail the computational energy consumption advantages of our system through a theoretical energy consumption analysis in the subsequent section.

5.5. Energy Efficiency

To further validate the energy efficiency of our neuromorphic signal processing system, we calculated the theoretical energy consumption of Synaptic Operations (SynOps) for TAE+QT-SNN, DC+SNN, and the same structure ANN networks. According to the standards established in the field of neuromorphic computing(Sengupta et al. (2019); Liu et al. (2023b,c)), the total SynOps required for the ANNs is defined as follows:

$$EN_{(ANNs)} = MAC \times \sum_{l=1}^L f_{in}^l N_l, \quad (11)$$

where f_{in}^l is the number of fan-in connections for the neurons in layer l , N_l is the number of neurons in layer l , and L is the total number of layers in the network.

The DC method, Widely used in SNNs, employs the original signal as the initial layer's input, repeating it across temporal steps. Thus, the total synaptic operation energy consumption for the DC+SNN is defined as:

$$EN_{(DC+SNN)} = MAC \times \sum_{t=1}^T \sum_{j=1}^{N_1} f_{in,j}^1 + AC \times \sum_{t=1}^T \sum_{l=2}^L \sum_{j=1}^{N_l} f_{out,j}^l o_j^l[t], \quad (12)$$

where the left half of the formula represents the synaptic energy consumption of the DC encoding layer, and the right half represents the synaptic operations energy consumption of the subsequent spiking layers. Here, $f_{out,j}^l$ is the number of fan-out connections from neuron j in layer l , T is the analog

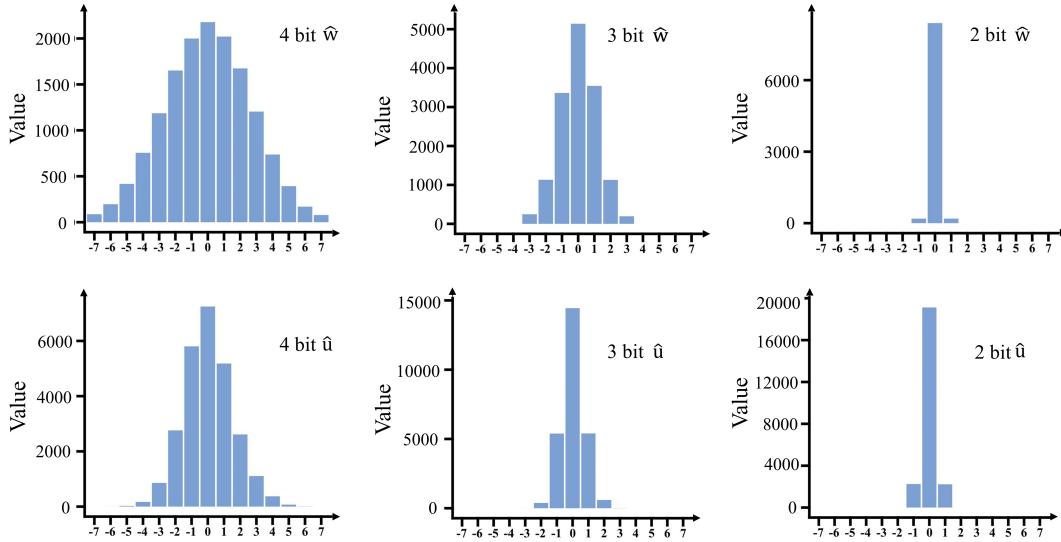


Figure 7: In our QT-SNN, when weights and membrane potentials are quantized to 4, 3, and 2 bits, their distributions adhere to normal distributions with varying standard deviations. Satisfactorily, as the bit-width decreases, the proportion of zero-valued weights (silent) and membrane potentials significantly increases, further enhancing the energy efficiency of the SNN.

Table 2

Comparison of model performance on GSC datasets for 12 and 35 classification.

Model	Network	Encoding	Memory(MB)	Precision(W/U)	Acc(%)
Google Speech Commands Dataset Version 2 (12)					
ST-Attention-SNN Wang et al. (2023b)	SNN	DC	2170	32/32	95.1
SLAYER-RF-CNN Orchard et al. (2021)	SNN	DC	888	32/32	91.4
SpikGRUDampfhofer et al. (2023)	SNN	DC	1587	32/32	92.9
(Our) SNN-KWS	SNN	DC	863	16/16	94.5
(Our) SNN-KWS	SNN	TAE	72	16/16	93.6
(Our) SNN-KWS	SNN	DC	538	4/4	93.9
(Our) SNN-KWS	SNN	TAE	63	4/4	93.2
(Our) SNN-KWS	SNN	DC	527	2/2	93.4
(Our) SNN-KWS	SNN	TAE	59	2/2	92.9
Google Speech Commands Dataset Version 2 (35)					
WaveSence Weidel & Sheik (2021)	SNN	DC	N/A	32/32	79.5
LSTMs-SNN Zhang et al. (2023)	SNN	DC	N/A	32/32	91.5
SRNN+ALIF Yin et al. (2021)	SNN	DC	3290	32/32	92.5
Speech2Spikes Stewart et al. (2023)	SNN	DC	2780	32/32	89.5
(Our) SNN-KWS	SNN	DC	871	16/16	92.8
(Our) SNN-KWS	SNN	TAE	74	16/16	92.3
(Our) SNN-KWS	SNN	DC	853	4/4	92.3
(Our) SNN-KWS	SNN	TAE	68	4/4	91.9
(Our) SNN-KWS	SNN	DC	849	2/2	92.1
(Our) SNN-KWS	SNN	TAE	60	2/2	91.8

time window, and $o_j^l[t]$ is the spike occurrence of neuron j at time t . Regrettably, the extensive use of MAC operations in the first layer significantly limits the energy efficiency of SNNs.

Our system achieves all ternary spike-based information transmission and computation, effectively avoiding MAC operations. Its' total synaptic operation energy consumption

is defined as:

$$En_{(TAE+QT-SNN)} = AC \times \sum_{t=1}^T \sum_{l=1}^L \sum_{j=1}^{N_l} f_{out,j}^l o_j^l[t]. \quad (13)$$

Here, we evaluate the hardware energy costs of our system(TAE+QT-SNN), DC+SNN, and ANN on a speech recognition task using the same network architecture. Our network architecture consists of 4 D-conv1d layers, 2

Table 3

Comparison of model performance on KUL and DTU datasets.

Dataset	Model	Network	Precision	Decision Windows(second)						
				0.1	0.2	0.5	1	2	5	10
KUL	De Cheveigné et al. (2018)	ANN	32	50.9	53.6	55.7	60.2	63.5	67.4	75.9
	Vandecappelle et al. (2021)	ANN	32	74.3	78.2	80.6	84.1	85.7	86.9	87.9
	Cai et al. (2021)	ANN	32	80.8	84.3	87.2	90.1	91.4	92.6	93.9
	(Our) TAE+QT-SNN	SNN	4	91.2	91.7	91.9	92.4	93.6	94.0	94.3
DTU	De Cheveigné et al. (2018)	ANN	32	-	-	-	53.4	57.7	61.9	70.1
	Vandecappelle et al. (2021)	ANN	32	56.7	58.4	61.7	63.6	65.2	67.4	67.8
	Cai et al. (2021)	ANN	32	65.7	68.1	70.8	71.9	73.7	76.1	75.8
	(Our) TAE+QT-SNN	SNN	4	66.8	68.6	70.4	72.2	73.6	76.8	76.6

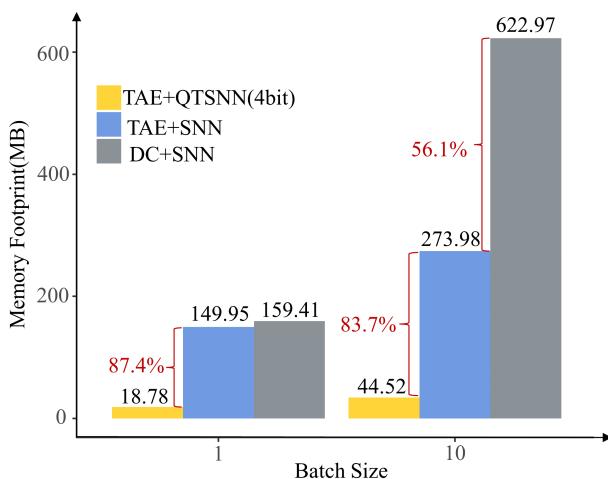


Figure 8: Both TAE and QT-SNN contribute to memory reduction across varying batch sizes. With a batch size of 1, QT-SNN achieves approximately an 87.4% reduction in memory usage. At a batch size of 10, TAE and QT-SNN reduce memory by 83.7% and 56.1%, respectively. Therefore, Both TAE and QT-SNN significantly contribute to reducing our system's memory footprints.

bottleneck blocks, and 1 FC layer. For SNN-based models, we conduct inference using 4 time steps. The average spike sparsity for TAE+QT-SNN is 10.42%, while for DC+SNN, it is 18.27% (excluding the first layer). (Guo et al., 2023; Hu et al., 2021b; Horowitz, 2014; Rueckauer et al., 2017) indicates a MAC operation requires 4.6pJ and an AC operation requires 0.9pJ, we compile the hardware energy costs of the three models as demonstrated in Table.4. The result reveals that our neuromorphic signal processing system offers an approximate 7.5× and 12.25× energy saving compared to similar DC-based SNN model and ANN-based model with the same model structures. Our system significantly enhances the energy efficiency and hardware-friendliness of intelligent signal processing models.

Table 4
Comparison of hardware energy consumption for different models

Method	MAC	AC	Timestep	Energy
Our	0M	15.19M	4	5.70uJ
DC+SNN	1.86M	13.33M	4	42.99uJ
ANNs	15.19	0M	1	69.87uJ

6. Conclusion

This study introduces a ternary spike-based neuromorphic signal processing system to achieve lightweight, energy-efficient, and high-performance intelligent signal processing models for resource-constrained edge devices. This approach effectively addresses the challenges of neural signal encoding and model complexity inherent in existing SNN-based solutions. Our system incorporates two innovative components: Firstly, we introduce the TAE method, which more efficiently encodes raw analog signals into ternary spike trains, facilitating more sparse information transmission. Secondly, our QT-SNN model complements our TAE method, achieving direct processing of ternary spike signal information. Moreover, it quantizes both synaptic weights and membrane potentials to lower bit-width, significantly reducing the memory footprint and enhancing the hardware friendliness of intelligent signal processing models. Extensive experimental evidence demonstrates that our system, in comparison to other similar SNN-based works, achieves superior performance with greater lightweight and energy-saving characteristics. Future work will focus on deploying our system on neuromorphic chips, promising to offer a novel perspective for edge signal processing.

CRediT authorship contribution statement

Shuai Wang: Conceptualization, Methodology, Experimentation, Writing - Original draft preparation. **Dehao Zhang:** Methodology, Investigation. **Ammar Belatreche:** Methodology, Investigation, Experimentation. **Yichen Xiao:** Methodology, Investigation. **Hongyu Qing:** Methodology, Investigation, Experimentation. **Wenjie Wei:** Methodology,

Investigation, Experimentation. **Malu Zhang**: Conceptualization, Supervision, Funding acquisition.. **Yang Yang**: Supervision, Investigation.

References

- Amir, A., Taba, B., Berg, D., Melano, T., McKinstry, J., Di Nolfo, C., Nayak, T., Andreopoulos, A., Garreau, G., Mendoza, M. et al. (2017). A low power, fully event-based gesture recognition system. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 7243–7252).
- Bartolozzi, C. (2018). Neuromorphic circuits impart a sense of touch. *Science*, 360, 966–967.
- Bouvier, M., Valentian, A., Mesquida, T., Rummens, F., Reyboz, M., Vianello, E., & Beigne, E. (2019). Spiking neural networks hardware implementations and challenges: A survey. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 15, 1–35.
- Bu, T., Ding, J., Yu, Z., & Huang, T. (2022). Optimized potential initialization for low-latency spiking neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence* (pp. 11–20). volume 36. Doi: 10.1609/aaai.v36i1.19874.
- Cai, S., Chen, Y., Huang, S., Wu, Y., Zheng, H., Li, X., & Xie, L. (2019). Svm-based classification of semg signals for upper-limb self-rehabilitation training. *Frontiers in neurorobotics*, 13, 31.
- Cai, S., Su, E., Xie, L., & Li, H. (2021). Eeg-based auditory attention detection via frequency and channel neural attention. *IEEE Transactions on Human-Machine Systems*, 52, 256–266.
- Cao, Y., Chen, Y., & Khosla, D. (2015). Spiking deep convolutional neural networks for energy-efficient object recognition. *International Journal of Computer Vision*, 113, 54–66. Doi: 10.1007/s11263-014-0788-3 .
- Castagnetti, A., Pegatoquet, A., & Miramond, B. (2023). Trainable quantization for speedy spiking neural networks. *Frontiers in Neuroscience*, 17, 1154241.
- Caviglia, S., Valle, M., & Bartolozzi, C. (2014). Asynchronous, event-driven readout of poset devices for tactile sensing. In *2014 IEEE International Symposium on Circuits and Systems (ISCAS)* (pp. 2648–2651). IEEE. Doi: 10.1109/iscas.2014.6865717 .
- Chan, V., Liu, S.-C., & van Schaik, A. (2007). Aer ear: A matched silicon cochlea pair with address event representation interface. *IEEE Transactions on Circuits and Systems I: Regular Papers*, 54, 48–59.
- Chowdhury, S. S., Garg, I., & Roy, K. (2021). Spatio-temporal pruning and quantization for low-latency spiking neural networks. In *2021 International Joint Conference on Neural Networks (IJCNN)* (pp. 1–9). IEEE.
- Chu, H., Yan, Y., Gan, L., Jia, H., Qian, L., Huan, Y., Zheng, L., & Zou, Z. (2022). A neuromorphic processing system with spike-driven snn processor for wearable ecg classification. *IEEE Transactions on Biomedical Circuits and Systems*, 16, 511–523.
- Dampfhofer, M., Mesquida, T., Hardy, E., Valentian, A., & Anghel, L. (2023). Leveraging sparsity with spiking recurrent neural networks for energy-efficient keyword spotting. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)* (pp. 1–5). IEEE.
- Das, N., Francart, T., & Bertrand, A. (2019). Auditory attention detection dataset kuleuvan. *Zenodo*.
- De Cheveigné, A., Wong, D. D., Di Liberto, G. M., Hjortkær, J., Slaney, M., & Lalor, E. (2018). Decoding the auditory brain with canonical component analysis. *NeuroImage*, 172, 206–216.
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., & Fei-Fei, L. (2009). Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition* (pp. 248–255). Ieee.
- Deng, L., Wu, Y., Hu, Y., Liang, L., Li, G., Hu, X., Ding, Y., Li, P., & Xie, Y. (2021). Comprehensive snn compression using admm optimization and activity regularization. *IEEE transactions on neural networks and learning systems*, 34, 2791–2805.
- Dupeyroux, J., Stroobants, S., & De Croon, G. C. (2022). A toolbox for neuromorphic perception in robotics. In *2022 8th International Conference on Event-Based Control, Communication, and Signal Processing (EBC CSP)* (pp. 1–7). IEEE.
- Faisal, A., Kamruzzaman, M., Yigitcanlar, T., & Currie, G. (2019). Understanding autonomous vehicles. *Journal of transport and land use*, 12, 45–72.
- Fang, W., Yu, Z., Chen, Y., Masquelier, T., Huang, T., & Tian, Y. (2021). Incorporating learnable membrane time constant to enhance learning of spiking neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision* (pp. 2661–2671). Doi: 10.1109/iccv48922.2021.00266 .
- Fuglsang, S. A., Dau, T., & Hjortkær, J. (2017). Noise-robust cortical tracking of attended speech in real-world acoustic scenes. *NeuroImage*, 156, 435–444.
- Fuglsang, S. A., Wong, D., & Hjortkær, J. (2018). Eeg and audio dataset for auditory attention decoding. *Zenodo*.
- de Gelder, L. (2021). Population step forward encoding algorithm: Improving the signal encoding accuracy and efficiency of spike encoding algorithms, .
- Giordano, C., Brennan, M., Mohamed, B., Rashidi, P., Modave, F., & Tighe, P. (2021). Assessing artificial intelligence for clinical decision-making. *Frontiers in digital health*, 3, 645232.
- Guo, Y., Chen, Y., Liu, X., Peng, W., Zhang, Y., Huang, X., & Ma, Z. (2023). Ternary spike: Learning ternary spikes for spiking neural networks. *arXiv preprint arXiv:2312.06372*, .
- Horowitz, M. (2014). 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. Doi: 10.1109/isscc.2014.6757323 .
- Hu, S., Qiao, G., Chen, T., Yu, Q., Liu, Y., & Rong, L. (2021a). Quantized stdp-based online-learning spiking neural network. *Neural Computing and Applications*, 33, 12317–12332.
- Hu, Y., Tang, H., & Pan, G. (2021b). Spiking deep residual networks. *IEEE Transactions on Neural Networks and Learning Systems*, 34, 5200–5205.
- Hwang, S., & Kung, J. (2024). One-spike snn: Single-spike phase coding with base manipulation for ann-to-snn conversion loss minimization. *arXiv preprint arXiv:2403.08786*, .
- Kasabov, N., Dhoble, K., Nuntalid, N., & Indiveri, G. (2013). Dynamic evolving spiking neural networks for on-line spatio-and spectro-temporal pattern recognition. *Neural Networks*, 41, 188–201.
- Kasabov, N., Scott, N. M., Tu, E., Marks, S., Sengupta, N., Capecci, E., Othman, M., Doborjeh, M. G., Murli, N., Hartono, R. et al. (2016). Evolving spatio-temporal data machines based on the neucube neuromorphic framework: Design methodology and selected applications. *Neural Networks*, 78, 1–14.
- Kim, K., Wu, F., Peng, Y., Pan, J., Sridhar, P., Han, K. J., & Watanabe, S. (2023). E-branchformer: Branchformer with enhanced merging for speech recognition. In *2022 IEEE Spoken Language Technology Workshop (SLT)* (pp. 84–91). IEEE.
- Li, C., Ma, L., & Furber, S. (2022). Quantization framework for fast spiking neural networks. *Frontiers in Neuroscience*, 16, 918793.
- Li, H., Liu, H., Ji, X., Li, G., & Shi, L. (2017). Cifar10-dvs: an event-stream dataset for object classification. *Frontiers in neuroscience*, 11, 244131.
- Liu, A. H., Hsu, W.-N., Auli, M., & Baevski, A. (2023a). Towards end-to-end unsupervised speech recognition. In *2022 IEEE Spoken Language Technology Workshop (SLT)* (pp. 221–228). IEEE.
- Liu, H., Chen, Y., Zeng, Z., Zhang, M., & Qu, H. (2023b). A low power and low latency fpga-based spiking neural network accelerator. In *2023 International Joint Conference on Neural Networks (IJCNN)* (pp. 1–8). IEEE.
- Liu, Y., Chen, Z., Wang, Z., Zhao, W., He, W., Zhu, J., Wang, O., Zhang, N., Jia, T., Ma, Y. et al. (2023c). Aa 22nm 0.43 pj/sop sparsity-aware in-memory neuromorphic computing system with hybrid spiking and artificial neural network and configurable topology. In *2023 IEEE Custom Integrated Circuits Conference (CICC)* (pp. 1–2). IEEE. Doi: 10.1109/cicc57935.2023.10121315 .

- López-Randulfe, J., Reeb, N., Karimi, N., Liu, C., Gonzalez, H. A., Dietrich, R., Vogginger, B., Mayr, C., & Knoll, A. (2022). Time-coded spiking fourier transform in neuromorphic hardware. *IEEE Transactions on Computers*, *71*, 2792–2802.
- Orchard, G., Frady, E. P., Rubin, D. B. D., Sanborn, S., Shrestha, S. B., Sommer, F. T., & Davies, M. (2021). Efficient neuromorphic signal processing with loihi 2. In *2021 IEEE Workshop on Signal Processing Systems (SiPS)* (pp. 254–259). IEEE. Doi: 10.1109/sips52927.2021.00053 .
- Orchard, G., Jayawant, A., Cohen, G. K., & Thakor, N. (2015). Converting static image datasets to spiking neuromorphic datasets using saccades. *Frontiers in neuroscience*, *9*, 437. Doi: 10.3389/fnins.2015.00437 .
- Pan, Z., Chua, Y., Wu, J., Zhang, M., Li, H., & Ambikairajah, E. (2020). An efficient and perceptually motivated auditory neural encoding and decoding algorithm for spiking neural networks. *Frontiers in neuroscience*, *13*, 1420. Doi: 10.3389/fnins.2019.01420 .
- Parekh, D., Poddar, N., Rajpurkar, A., Chahal, M., Kumar, N., Joshi, G. P., & Cho, W. (2022). A review on autonomous vehicles: Progress, methods and challenges. *Electronics*, *11*, 2162.
- Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., & Sutskever, I. (2023). Robust speech recognition via large-scale weak supervision. In *International Conference on Machine Learning* (pp. 28492–28518). PMLR.
- Rueckauer, B., Lungu, I.-A., Hu, Y., Pfeiffer, M., & Liu, S.-C. (2017). Conversion of continuous-valued deep networks to efficient event-driven networks for image classification. *Frontiers in neuroscience*, *11*, 682. Doi: 10.3389/fnins.2017.00682 .
- Safa, A., Corradi, F., Keuninckx, L., Ocket, I., Bourdoux, A., Catthoor, F., & Gielen, G. G. (2021). Improving the accuracy of spiking neural networks for radar gesture recognition through preprocessing. *IEEE Transactions on Neural Networks and Learning Systems*, *34*, 2869–2881.
- Sengupta, A., Ye, Y., Wang, R., Liu, C., & Roy, K. (2019). Going deeper in spiking neural networks: Vgg and residual architectures. *Frontiers in neuroscience*, *13*, 95. Doi: 10.3389/fnins.2019.00095 .
- Stewart, K. M., Shea, T., Pacik-Nelson, N., Gallo, E., & Danilescu, A. (2023). Speech2spikes: Efficient audio encoding pipeline for real-time neuromorphic systems. In *Proceedings of the 2023 Annual Neuro-Inspired Computational Elements Conference* (pp. 71–78). Doi: 10.1145/3584954.3584995 .
- Stöckl, C., & Maass, W. (2021). Optimized spiking neurons can classify images with high accuracy through temporal coding with two spikes. *Nature Machine Intelligence*, *3*, 230–238.
- Su, E., Cai, S., Xie, L., Li, H., & Schultz, T. (2022). Stanet: A spatiotemporal attention network for decoding auditory spatial attention from eeg. *IEEE Transactions on Biomedical Engineering*, *69*, 2233–2242.
- Sulaiman, M. B. G., Juang, K.-C., & Lu, C.-C. (2020). Weight quantization in spiking neural network for hardware implementation. In *2020 IEEE International Conference on Consumer Electronics-Taiwan (ICCE-Taiwan)* (pp. 1–2). IEEE.
- Taherkhani, A., Belatreche, A., Li, Y., Cosma, G., Maguire, L. P., & McGinnity, T. M. (2020). A review of learning in biologically plausible spiking neural networks. *Neural Networks*, *122*, 253–272.
- Tan, P.-Y., Wu, C.-W., & Lu, J.-M. (2021). An improved stbp for training high-accuracy and low-spike-count spiking neural networks. In *2021 Design, Automation & Test in Europe Conference & Exhibition (DATE)* (pp. 575–580). IEEE.
- Tavaneai, A., Ghodrati, M., Kheradpisheh, S. R., Masquelier, T., & Maida, A. (2019). Deep learning in spiking neural networks. *Neural networks*, *111*, 47–63.
- Troussard, C., Dufrechou, L., Tremblin, P.-A., & Eustache, Y. (2023). Real-time monitoring of coastal & offshore construction noise for immediate decision making. In *Abu Dhabi International Petroleum Exhibition and Conference* (p. D021S045R002). SPE.
- Valsalan, P., Baomar, T. A. B., & Baabood, A. H. O. (2020). Iot based health monitoring system. *Journal of critical reviews*, *7*, 739–743.
- Vandecappelle, S., Deckers, L., Das, N., Ansari, A. H., Bertrand, A., & Francart, T. (2021). Eeg-based detection of the locus of auditory attention with convolutional neural networks. *Elife*, *10*, e56481.
- Wang, Q., Zhang, T., Han, M., Wang, Y., Zhang, D., & Xu, B. (2023a). Complex dynamic neurons improved spiking transformer network for efficient automatic speech recognition. In *Proceedings of the AAAI Conference on Artificial Intelligence* (pp. 102–109). volume 37.
- Wang, Y., Shi, K., Lu, C., Liu, Y., Zhang, M., & Qu, H. (2023b). Spatial-temporal self-attention for asynchronous spiking neural networks. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI-23, Edith Elkind, Ed* (pp. 3085–3093). volume 8. Doi: 10.24963/ijcai.2023/344 .
- Weidel, P., & Sheik, S. (2021). Wavesense: Efficient temporal convolutions with spiking neural networks for keyword spotting. *arXiv preprint arXiv:2111.01456*, . Doi: 10.48550/arXiv.2111.01456 .
- Wu, J., Chua, Y., Zhang, M., Li, H., & Tan, K. C. (2018). A spiking neural network framework for robust sound classification. *Frontiers in neuroscience*, *12*, 836. Doi: 10.3389/fnins.2018.00836 .
- Wu, Y., Deng, L., Li, G., Zhu, J., Xie, Y., & Shi, L. (2019). Direct training for spiking neural networks: Faster, larger, better. In *Proceedings of the AAAI conference on artificial intelligence* (pp. 1311–1318). volume 33.
- Xiao, R., Yan, R., Tang, H., & Tan, K. C. (2017). A spiking neural network model for sound recognition. In *Cognitive Systems and Signal Processing: Third International Conference, ICCSIP 2016, Beijing, China, November 19–23, 2016, Revised Selected Papers 3* (pp. 584–594). Springer.
- Yang, Q., Liu, Q., & Li, H. (2022). Deep residual spiking neural network for keyword spotting in low-resource settings. *Proc. Interspeech 2022*, (pp. 3023–3027). Doi: 10.21437/interspeech.2022-107 .
- Yao, M., Hu, J., Hu, T., Xu, Y., Zhou, Z., Tian, Y., Bo, X., & Li, G. (2023a). Spike-driven transformer v2: Meta spiking neural network architecture inspiring the design of next-generation neuromorphic chips. In *The Twelfth International Conference on Learning Representations*.
- Yao, M., Hu, J., Zhou, Z., Yuan, L., Tian, Y., Xu, B., & Li, G. (2024). Spike-driven transformer. *Advances in Neural Information Processing Systems*, *36*.
- Yao, M., Zhao, G., Zhang, H., Hu, Y., Deng, L., Tian, Y., Xu, B., & Li, G. (2023b). Attention spiking neural networks. *IEEE transactions on pattern analysis and machine intelligence*, .
- Yilmaz, E., Gevrek, O. B., Wu, J., Chen, Y., Meng, X., & Li, H. (2020). Deep convolutional spiking neural networks for keyword spotting. In *Proceedings of INTERSPEECH* (pp. 2557–2561). Doi: 10.21437/interspeech.2020-1230 .
- Yin, B., Corradi, F., & Bohté, S. M. (2021). Accurate and efficient time-domain classification with adaptive spiking recurrent neural networks. *Nature Machine Intelligence*, *3*, 905–913. Doi: 10.1101/2021.03.22.436372 .
- Yin, R., Li, Y., Moitra, A., & Panda, P. (2023). Mint: Multiplier-less integer quantization for spiking neural networks. *arXiv preprint arXiv:2305.09850*, .
- Yoo, D., & Jeong, D. S. (2023). Cbp-qsnn: Spiking neural networks quantized using constrained backpropagation. *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, *13*, 1137–1146. doi:10.1109/JETCAS.2023.3328911 .
- Zhang, A., Li, X., Gao, Y., & Niu, Y. (2021a). Event-driven intrinsic plasticity for spiking convolutional neural networks. *IEEE Transactions on Neural Networks and Learning Systems*, *33*, 1986–1995.
- Zhang, M., Wang, J., Wu, J., Belatreche, A., Amornpaisanon, B., Zhang, Z., Miriyala, V. P. K., Qu, H., Chua, Y., Carlson, T. E. et al. (2021b). Rectified linear postsynaptic potential function for backpropagation in deep spiking neural networks. *IEEE transactions on neural networks and learning systems*, *33*, 1947–1958.
- Zhang, S., Yang, Q., Ma, C., Wu, J., Li, H., & Tan, K. C. (2023). Long short-term memory with two-compartment spiking neuron. *arXiv preprint arXiv:2307.07231*, . Doi: 10.48550/arXiv.2307.07231 .
- Zheng, H., Wu, Y., Deng, L., Hu, Y., & Li, G. (2021). Going deeper with directly-trained larger spiking neural networks. In *Proceedings of the AAAI conference on artificial intelligence* (pp. 11062–11070). volume 35.
- Zhou, S., Li, X., Chen, Y., Chandrasekaran, S., & Sanyal, A. (2021). Temporal-coded deep spiking neural network with easy training and robust performance. In *35th AAAI Conference on Artificial Intelligence*,

AAAI 2021 35th AAAI Conference on Artificial Intelligence, AAAI 2021 (pp. 11143–11151). Association for the Advancement of Artificial Intelligence.

Zhou, Z., Che, K., Fang, W., Tian, K., Zhu, Y., Yan, S., Tian, Y., & Yuan, L. (2024). Spikformer v2: Join the high accuracy club on imagenet with an snn ticket. *arXiv preprint arXiv:2401.02020*, .

Zhou, Z., Zhu, Y., He, C., Wang, Y., Yan, S., Tian, Y., & Yuan, L. (2022). Spikformer: When spiking neural network meets transformer. *arXiv preprint arXiv:2209.15425*, .

Zhu, R.-J., Zhao, Q., Zhang, T., Deng, H., Duan, Y., Zhang, M., & Deng, L.-J. (2022). Tcja-snn: Temporal-channel joint attention for spiking neural networks. *arXiv preprint arXiv:2206.10177*, .