

DS-ATGO: Dual-Stage Synergistic Learning via Forward Adaptive Threshold and Backward Gradient Optimization for Spiking Neural Networks

Jiaqiang Jiang^{1,2}, Wenfeng Xu^{1,2}, Jing Fan^{1,2}, Rui Yan^{1,2*}

¹College of Computer Science and Technology, Zhejiang University of Technology, Hangzhou 310023

²Zhejiang Key Laboratory of Visual Information Intelligent Processing, Hangzhou 310023

{jqjiang, wfxu, fanjing, ryan}@zjut.edu.cn

Abstract

Brain-inspired spiking neural networks (SNNs) are recognized as a promising avenue for achieving efficient, low-energy neuromorphic computing. Direct training of SNNs typically relies on surrogate gradient (SG) learning to estimate derivatives of non-differentiable spiking activity. However, during training, the distribution of neuronal membrane potentials varies across timesteps and progressively deviates toward both sides of the firing threshold. When the firing threshold and SG remain fixed, this may lead to imbalanced spike firing and diminished gradient signals, preventing SNNs from performing well. To address these issues, we propose a novel dual-stage synergistic learning algorithm that achieves forward adaptive thresholding and backward dynamic SG. In forward propagation, we adaptively adjust thresholds based on the distribution of membrane potential dynamics (MPD) at each timestep, which enriches neuronal diversity and effectively balances firing rates across timesteps and layers. In backward propagation, drawing from the underlying association between MPD, threshold, and SG, we dynamically optimize SG to enhance gradient estimation through spatio-temporal alignment, effectively mitigating gradient information loss. Experimental results demonstrate that our method achieves significant performance improvements. Moreover, it allows neurons to fire stable proportions of spikes at each timestep and increases the proportion of neurons that obtain gradients in deeper layers.

Code — <https://github.com/jqjiang1999/DS-ATGO>

Introduction

As a new paradigm that combines biological plausibility with computational efficiency, spiking neural networks (SNNs) have recently attracted widespread attention (Maass 1997). Unlike artificial neural networks (ANNs), which work with continuous activation and real-valued encoding (Ju et al. 2025b,a; Zhang et al. 2025; Feng et al. 2025; Zhang, Yuan, and Pan 2024), SNNs operate with asynchronous, discrete spiking signals. This spike train-based way renders SNNs highly promising for spatio-temporal (ST) information processing and energy-efficient computation compared to ANNs (Roy, Jaiswal, and Panda 2019).

*Corresponding author: Rui Yan (ryan@zjut.edu.cn)

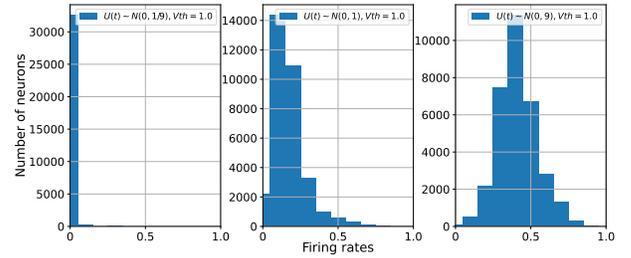


Figure 1: The distributions of firing rates at differ variances of membrane potential when $V_{th} = 1.0$ (Zheng et al. 2021).

Despite these advantages, the non-differentiable nature of spike activity hinders gradient backpropagation, making it more difficult to directly train high-performance SNNs than ANNs. To overcome this issue, surrogate gradient (SG) learning (Wu et al. 2018) estimates the gradients of output signals by introducing continuous smooth functions, which reconstruct a complete propagation path of ST gradients.

However, as spikes propagate through layers, the distribution of membrane potentials shifts and may fall into inappropriate areas (Guo et al. 2022b; Liu et al. 2025) (Fig. 2). To our best knowledge, most existing SG-based methods for directly training SNNs use a fixed threshold and SG. Thus, membrane potential dynamics (MPD) may present two challenges that lead to training difficulties.

The first challenge is the imbalanced spike firing of SNNs with a fixed threshold. The neuron fires a spike only when its membrane potential exceeds the threshold V_{th} . In Fig. 1, when membrane potentials are too small compared to the threshold, excessive sparse firing in neurons will lead to the “*spike vanishing problem*”. When membrane potentials are too large, neurons will be over-firing, reducing the ability to represent the input pattern differentially. To maintain the stability of information flow in SNN, we need to balance the threshold and membrane potentials to keep neurons in a moderately active state. Neuroscience has observed in various brain regions that the thresholds of biological neurons are not constant but exhibit variability between and within neurons (Azouz and Gray 2000; Farries, Kita, and Wilson 2010). This phenomenon, known as threshold plasticity, can be regarded as an adaptation to membrane potentials, which plays an essential role in maintaining neuronal firing home-

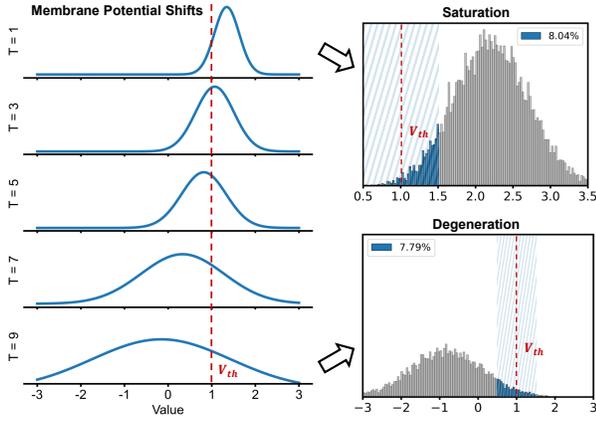


Figure 2: The distribution of membrane potentials deviating from the threshold in a vanilla SNN with ten timesteps. When almost all the membrane potentials of neurons are beyond V_{th} , called saturation. Conversely, called degeneration.

ostasis (Fontaine, Peña, and Brette 2014). To model threshold plasticity, (Wang, Cheng, and Lim 2022; Rathi and Roy 2023; Sun et al. 2024; Hasssan, Meng, and Seo 2024) directly set the threshold as a learnable parameter and optimize it dynamically via gradients. Although this effectively enhances neuronal firing levels, its optimization relies solely on the gradient of thresholds itself, leaving uncertainty in maintaining the balance of firing rates in SNNs.

The other challenge is the diminished gradient information of SNNs with a fixed SG. The neuron obtains a gradient only when its membrane potential falls within the gradient-available interval of SG. In Fig. 2, when membrane potentials is too small or too large compared to the threshold, the limited gradient-available interval of fixed SG causes the membrane potentials of many neurons to fall into the areas where the approximate derivatives are zero or a tiny value. In this case, only a few neurons contribute gradients, leading to the “*gradient vanishing problem*”. (Lian et al. 2023) dynamically adjusted the SG based on changes in the MPD distribution. It ignores the dynamic nature of evolving MPD in the temporal dimension, which may exacerbate the inaccurate estimation of gradients. Then, (Jiang et al. 2025; Liu et al. 2025) optimized the SG of different timesteps in a temporal-aligned manner. In SG learning, the gradient-available interval of SG is symmetrically centered around the threshold. However, these methods overlook the association between MPD, threshold, and SG, making SG fail to accurately capture the dynamic deviations of membrane potentials relative to thresholds. Therefore, the neuronal threshold can be regarded as an idealized tunable parameter that regulates spike firing and affects gradient optimization.

In this paper, we propose a novel dual-stage synergistic learning via forward adaptive threshold and backward gradient optimization for SNNs, named DS-ATGO. By utilizing the distribution characteristics of MPD at each timestep, we adaptively adjust thresholds to enhance the ability of neurons to encode dynamic information, achieving a linear response to input signals. Building on the observation that

adaptive thresholds reflect shifts in membrane potentials, we further establish a correlation between SG and MPD via thresholds. Based on that, we dynamically adjust SG in a threshold-driven manner to align with evolving MPD, contributing smoother spatio-temporal gradients and maintaining the optimal gradient domain across timesteps. The overall framework of DS-ATGO is illustrated in Fig. 3. The main contributions of our work are summarized as follows:

- We propose an adaptive threshold mechanism that adapts the firing threshold in the temporal dimension by leveraging the MPD distribution. It balances neuronal spike firing across layers, enhancing the cross-layer information transmission capability of SNNs.
- We propose a threshold-driven gradient optimization method that dynamically adjusts SG at each timestep by associating SG with evolving MPD via adaptive thresholds, effectively enhancing ST gradient information.
- The experimental results on both static and neuromorphic datasets demonstrate that DS-ATGO achieves outstanding performance with low latency without additional inference overhead.

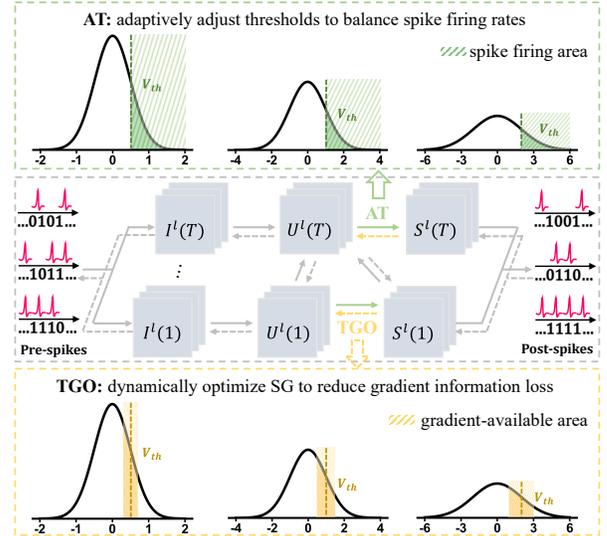


Figure 3: The overall framework of DS-ATGO. Internal dynamics of LIF neurons in a layer (gray). In forward propagation, the adaptive threshold (AT) mechanism promotes neurons to generate stable firing rates under different MPD distributions (green). In backward propagation, the threshold-driven SG optimization (TGO) method dynamically scales SG to respond to evolving MPD (yellow).

Related Works

Spiking Threshold Plasticity

Spiking threshold plasticity is a biological mechanism by which neurons dynamically adjust their firing thresholds according to historical activity patterns. To express this characteristic, (Ai et al. 2025) dynamically adjusted the threshold based on the repetitive or high-frequency discharge re-

sponses of neurons, effectively preventing excessive activation of neurons. (Fu et al. 2025) proposed a spatio-temporal threshold adjustment strategy, which enhances the spike coding capacity of SNNs by coupling with neural dynamics. In addition, some studies on SNN training have introduced learning dynamics into spike processes (Wang, Cheng, and Lim 2022; Rathi and Roy 2023). (Sun et al. 2024) proposed a synapse-threshold synergistic learning, which allows stable signal transmission through appropriate firing rates by simultaneously training weights and thresholds. Considering the weights and threshold same landscape makes the learning sub-optimal, (Hasssan, Meng, and Seo 2024) achieved individually adaptive optimization of layer-wise thresholds and weights by introducing a separate gradient path.

Gradient Optimization

The concept of gradient optimization in SNN aims to promote the efficient propagation of error signals in the ST domain by adjusting SG (Fang et al. 2021a). (Guo et al. 2022a) approximated the spike activity gradient by a continuously differentiable evolving asymptotic function, bridging the gap between pseudo and natural derivatives. (Che et al. 2022) proposed a differentiable gradient search method to optimize SG locally. (Lian et al. 2023) unlocked the width limitation of SG based on changes in MPD distributions, increasing the proportion of neurons that obtained gradients. (Wang et al. 2023) adaptively learned the accurate gradients of the loss landscape in SNN by fusing the learnable relaxation degree into a prototype network with random spike noise, effectively eliminating the smoothing error of SG learning. (Wang, Cheng, and Lim 2025) proposed a parametric SG strategy to control and calibrate SG, which mitigates the degradation caused by improper selection of SG. (Jiang et al. 2025; Liu et al. 2025) promotes balanced optimization by enhancing the gradients in a temporal-aligned manner.

Motivation: Overall, although existing methods achieve threshold or SG adjustment, they ignore the intrinsic association among membrane potential, threshold, and SG, which have limitations in jointly solving the two problems caused by membrane potential shifts. This motivates us to achieve the synergistic optimization of thresholds and SG to maintain balanced firing rates and stable gradient flows in SNNs.

Preliminaries

Spiking Neural Model

In this work, we use the Leaky Integrate-and-Fire (LIF) model in SNNs. The dynamics of an iterative LIF neuron (Wu et al. 2019) can be modeled by

$$I_i^l(t) = \sum_{j=1}^{N(l-1)} w_{ij}^l S_j^{l-1}(t), \quad (1)$$

$$U_i^l(t) = \tau U_i^l(t-1) \odot (1 - S_i^l(t-1)) + I_i^l(t), \quad (2)$$

$$S_i^l(t) = \Theta(U_i^l(t)) = \begin{cases} 1, & U_i^l(t) \geq V_{th} \\ 0, & \text{otherwise} \end{cases} \quad (3)$$

where the upper l , subscripts i and t denote the l -th layer, i -th neuron, and t -th timestep, respectively. $N(l-1)$ denotes

the number of neurons in the $(l-1)$ -th layer. w_{ij}^l denotes the synapse weight from the j -th neuron in the $(l-1)$ -th layer to the i -th neuron in the l -th layer. I , U , and S denote the pre-synaptic input, the membrane potential, and the output spike of neurons. V_{th} is the firing threshold. τ is the membrane time constant that indicates how quickly the potential decays over time, affecting the duration of neuronal excitation.

Surrogate Gradient Function

In Eq. 3, the firing function $\Theta(\cdot)$ of SNNs is a Heaviside function. The derivative of output signals $\frac{\partial S}{\partial U}$ is a Dirac function that tends to infinity at the threshold V_{th} and zeros otherwise. Here, we employ the rectangular function (Wu et al. 2018) to estimate the gradients of $\frac{\partial S}{\partial U}$, which is defined as

$$\frac{\partial S_i^l(t)}{\partial U_i^l(t)} \approx h(U_i^l(t)) = \frac{1}{\kappa} \text{sign}(|U_i^l(t) - V_{th}| \leq \frac{\kappa}{2}), \quad (4)$$

where hyperparameter κ controls the width of $h(\cdot)$ to ensure that it integrates to 1, normally set to 1. The gradient $\frac{1}{\kappa}$ is available when the membrane potential $U_i^l(t)$ falls within the interval $[V_{th} - \frac{\kappa}{2}, V_{th} + \frac{\kappa}{2}]$.

Methods

Adaptive Threshold Mechanism

It is difficult for SNNs with a fixed threshold to adapt to variations in input strength, making the network prone to silence or excessive firing when processing ST dynamic data. (Zheng et al. 2021) proposed a threshold-dependent batch normalization that normalized the pre-synaptic input $I(t)$ to $N(0, V_{th}^2)$ instead of $N(0, 1)$. Although this method balances pre-synaptic inputs and thresholds to increase neuronal spike activity, the intrinsic dynamics of neurons still hinder it from maintaining the stability of firing rates across the entire timestep. That motivates us to explore how to adaptively adjust the firing threshold to keep neurons homeostatic, enhancing the dynamic responsiveness of SNNs.

Recent studies have revealed that the threshold of biological neurons exhibits a positive correlation with membrane potentials (Pena and Konishi 2002; Azouz and Gray 2003; Ding et al. 2022). Based on the insight that firing rates are directly regulated by membrane potential and threshold, we can adjust the threshold according to the evolving MPD, promoting neurons to fire appropriate spikes in SNNs. To achieve this, the key lies in balancing the proportion of membrane potentials that exceed the threshold. **Theorem 1** shows that when the threshold is set as the sum of expectation and standard deviation of MPD distributions ($V_{th} = \mu + \sigma$), the proportion of membrane potentials exceeding the threshold remains relatively constant under different distributions. Establishing a positive correlation between thresholds and membrane potentials retains the adaptive properties of biological neurons while ensuring a stable firing rate in SNNs.

Theorem 1. For $U(t)$ that satisfies a normal distribution $N(\mu, \sigma^2)$, the probability that a random variable $U_i(t)$ exceeds $\mu + \sigma$ is relatively constant and given by $P(U_i(t) > \mu + \sigma) = 1 - \Phi(1)$, where $\Phi(\cdot)$ denotes the cumulative distribution function of standard normal distributions.

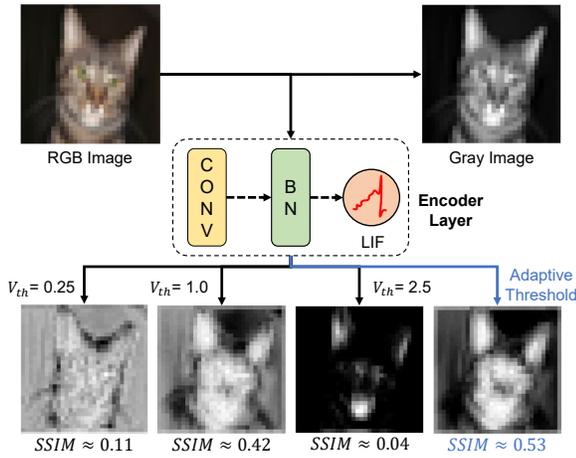


Figure 4: The structural similarity between the gray image and encoded images at different thresholds.

Proof. The proof of Theorem 1 is given in **Appendix A**. \square

To maintain neurons in a moderately active state, we first need to analyze the detailed distributions of membrane potentials. (Zheng et al. 2021) derived a high degree of similarity between the distribution of pre-synaptic inputs and membrane potentials. (Lian et al. 2023; Jiang et al. 2025) extended this theorem that for a given pre-synaptic inputs $I(t) \sim N(0, V_{th}^2)$, the distribution of membrane potentials satisfies $U(t) \sim N(0, (1 + \tau^2)V_{th}^2)$. Given the approximate nature of MPD distributions derived by (Lian et al. 2023; Jiang et al. 2025), we can also directly use the true MPD distribution during training. Moreover, to trade-off energy efficiency and performance, we introduce a factor f_c to control spiking firing rates. We will discuss these two distributions and the effect of factor f_c in the ablation study. Finally, the adaptive threshold mechanism can be formulated as

$$\begin{aligned} \Delta V_{th}^l(t)_n &\approx f_c * \sqrt{(1 + \tau_n^2)} V_{th}, \#Estimated & (5) \\ &= f_c * (\mathbb{E}(U^l(t)_n) + \sqrt{\text{VAR}(U^l(t)_n)}), \#True \end{aligned}$$

where the subscript n denotes the n -th mini-batch, τ is layer-wise learnable (Fang et al. 2021b). In this way, all neurons in different layers will have distinct thresholds across different timesteps. Fig. 4 shows that using the proposed adaptive threshold in the encoding layer produces better quality images than the fixed threshold.

Considering that different mini-batches in the inference stage may cause fluctuations in MPD distributions that affect the threshold, inspired by BN (Ioffe and Szegedy 2015), we also use the moving average of thresholds during training for inference to stabilize neuron outputs. It can be described as

$$\Delta V_{th}^l(t) = m * \Delta V_{th}^l(t)_n + (1 - m) * \Delta V_{th}^l(t), \quad (6)$$

where m is the momentum coefficient, which we set to 0.1.

Threshold-driven Gradient Optimization

In SNNs, adequate gradient information is crucial for achieving efficient learning, as it directly governs weight

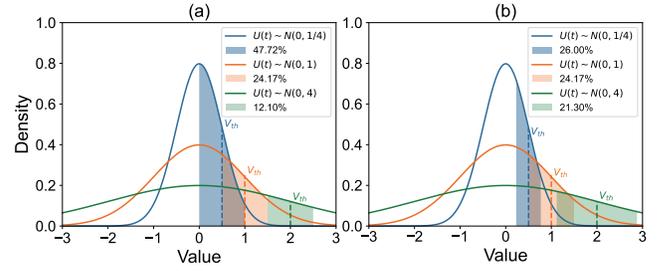


Figure 5: (a) The proportion of membrane potentials that fall into the gradient-available interval with fixed SG ($k = 1$) under different distributions when using adaptive threshold. (b) The proportion of gradient-available when using the adaptive threshold and threshold-driven gradient optimization.

updates and model convergence. When SG remains fixed, membrane potential shifts will cause decreased gradients during backpropagation. The adaptive threshold mechanism adjusts the gradient-available interval of fixed SG to $[\Delta V_{th} - \frac{1}{2}, \Delta V_{th} + \frac{1}{2}]$ instead of $[V_{th} - \frac{1}{2}, V_{th} + \frac{1}{2}]$, but the above issue remains due to the inherent mismatch between fixed SG and evolving MPD. Specifically, in Fig. 5, when the variance of MPD distributions becomes smaller (blue), a large number of neuronal membrane potentials will fall within the gradient-available interval for gradient computation, enlarging the approximation error with the natural gradient. Conversely, when the variance of MPD distributions becomes larger (green), only a few neurons will obtain gradients, and most neurons will fall into the saturation area with zero gradients, resulting in vanishing gradient.

To optimize SG learning, we need to modify SG accordingly to better respond to evolving MPD. (Lian et al. 2023; Jiang et al. 2025) indicated an association between SG and MPD. Given that the adaptive threshold reflects changes in MPD distributions (Eq. 5) and that the threshold determines the location of SG. Then, we propose a threshold-driven gradient optimization method that dynamically adjusts SG width (Fig. 5), ensuring that the gradient-available interval aligns with MPD. For the above two distinct cases, we design two rules to calculate the SG width, as follows:

- $\Delta V_{th} < V_{th}$: When the adaptive threshold ΔV_{th} is less than the initial threshold V_{th} , which indicates the MPD distribution tends to concentrate, thus we need to reduce the SG width to decrease the proportion of neurons fall into the gradient-available interval, suppressing the cumulative enlargement of gradient approximation errors.
- $\Delta V_{th} \geq V_{th}$: When the adaptive threshold ΔV_{th} is greater than the initial threshold V_{th} , which indicates the MPD distribution tends to disperse, thus we need to expand the SG width to increase the proportion of neurons falling into the gradient-available interval, mitigating the loss of gradient information.

In summary, the threshold-driven gradient optimization method can be mathematically written as

$$k = \begin{cases} (1 - \tanh(V_{th} - \Delta V_{th})) * k, & \Delta V_{th} < V_{th} \\ (1 + \tanh(\Delta V_{th} - V_{th})) * k, & \Delta V_{th} \geq V_{th} \end{cases} \quad (7)$$

The Overall Procedure of SNNs

In the output layer, we only accumulate the membrane potential of output neurons without leakage and firing (Rathi and Roy 2023), which can be described by

$$U_i^{l^o} = \frac{1}{T} \sum_{t=1}^T \sum_{j=1}^{N(l^o-1)} w_{ij}^{l^o} S_j^{l^o-1}(t), i \in \{1, 2, \dots, c\} \quad (8)$$

where l^o and c denote the output layer and the number of classes, respectively. As $U_i^l(t)$ not only contributes to $S_i^l(t)$ but also governs $U_i^l(t+1)$, it can be derived by

$$\frac{\partial L}{\partial U_i^l(t)} = \frac{\partial L}{\partial S_i^l(t)} \frac{\partial S_i^l(t)}{\partial U_i^l(t)} + \frac{\partial L}{\partial U_i^l(t+1)} \frac{\partial U_i^l(t+1)}{\partial U_i^l(t)}, \quad (9)$$

$$\begin{aligned} \frac{\partial L}{\partial S_i^l(t)} &= \frac{\partial L}{\partial U_i^l(t+1)} \frac{\partial U_i^l(t+1)}{\partial S_i^l(t)} \\ &+ \sum_{j=1}^{N(l+1)} \frac{\partial L}{\partial U_j^{l+1}(t)} \frac{\partial U_j^{l+1}(t)}{\partial S_i^l(t)}. \end{aligned} \quad (10)$$

where $\frac{\partial S_i^l(t)}{\partial U_i^l(t)}$ are approximated by Eq. 4. Then, the gradient of synaptic weights w_{ij}^l can be derived by the chain rule:

$$\begin{aligned} \frac{\partial L}{\partial w_{ij}^l} &= \sum_{t=1}^T \frac{\partial L}{\partial U_i^l(t)} \frac{\partial U_i^l(t)}{\partial I_i^l(t)} \frac{\partial I_i^l(t)}{\partial w_{ij}^l} \\ &= \sum_{t=1}^T \frac{\partial L}{\partial U_i^l(t)} \sum_{j=1}^{N(l-1)} S_j^{l-1}(t). \end{aligned} \quad (11)$$

Moreover, the pseudocode of the overall procedure is briefed in **Appendix B**.

Experiments

Ablation Study

The performance improvement of our model benefits from the adaptive threshold mechanism (AT) and threshold-driven gradient optimization method (TGO). We conducted experiments to evaluate their contributions. In Fig. 6, applying AT (w/ AT) outperforms Vanilla-SNN on two datasets. This indicates that adaptively adjusting thresholds in response to evolving MPD enables flexible control over the timing and intensity of neuron activation, enhancing the SNN’s ability to model complex dynamic information. Notably, applying TGO (w/ TGO) also achieves significant improvements over Vanilla-SNN, even surpassing AT. The reason is that, unlike AT, which focuses on spike modulation in information encoding, TGO synchronously adjusts SG to accurately capture the dynamic deviation of membrane potentials relative to thresholds. It effectively maintains the balance of gradient signals in the temporal dimension, enhancing the efficiency and correctness of parameter updates. Furthermore, combining the two techniques (w/ AT+TGO) enables SNNs to balance firing rates for improved encoding efficiency and optimize SG through temporal alignment for gradient enhancement, further boosting performance. It achieved remarkable

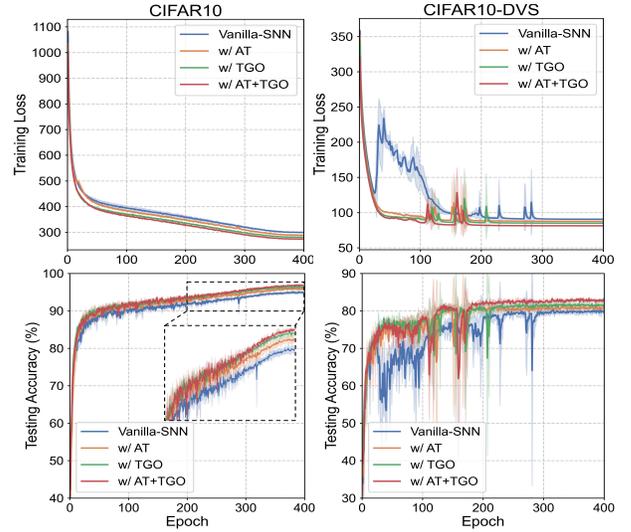


Figure 6: Comparison of training loss and test accuracy.

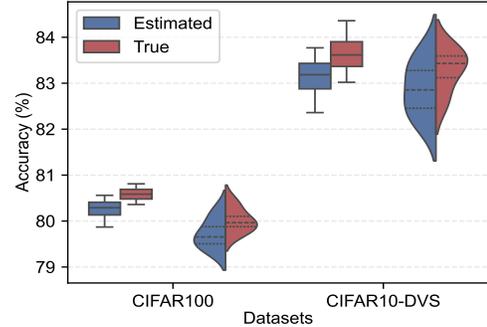


Figure 7: Comparison of estimated with true MPD distribution and w/ vs. w/o moving average (box plot vs. violin plot).

improvements of 1.71% and 2.30% on the CIFAR10 and CIFAR10-DVS datasets, respectively.

Moreover, Fig. 7 illustrates that DS-ATGO driven by the true MPD distribution outperforms the theoretical distribution derived by (Lian et al. 2023), which indicates that the theorem derivation involves certain approximation errors. Then, we evaluated the effectiveness of the moving average rule. When applying this rule in AT, the accuracy of both datasets surpassed the baseline without the rule. Especially for CIFAR10-DVS trained with mini-batches, the threshold moving average more effectively adapts to MPD fluctuations, exhibiting higher performance improvements.

Performance Evaluation

As listed in Table 1, we compare the classification accuracy of our method with other advanced methods. DS-ATGO performed well on all four datasets, achieving an accuracy of 96.91%/80.59% with low latency on CIFAR10/100, 83.70% and 68.86% accuracy on CIFAR10-DVS and ImageNet, respectively. Methods that only adjust thresholds (e.g., LT-SNN) or only optimize SG (e.g., DeepTAGE) both perform worse than ours, indicating that the two challenges arising

Dataset	Method	AT	GO	Architecture	Timestep	Accuracy(%)
CIFAR10	STAtten+ (Lee et al. 2025)	✗	✗	Spikformer-4-384	4	94.36
	Distillation-based SNN (Yu et al. 2025a)	✗	✗	ResNet-18	6 / 4	95.96 / 95.57
	STL-SNN (Sun et al. 2024)	✓	✗	8-layers SNN	8	92.42
	LT-SNN (Hasssan, Meng, and Seo 2024)	✓	✗	Spikformer-4-256	4	95.19
	DeepTAGE (Liu et al. 2025)	✗	✓	ResNet-18	4	95.86
	MPD-AGL (Jiang et al. 2025)	✗	✓	ResNet-19	4 / 2	96.35 / 96.18
	AT-LIF&ASG-S (Ai et al. 2025)	✓	✓	ResNet-18	6 / 4	95.47 / 95.30
Ours	✓	✓	ResNet-19	2	96.91±0.12	
CIFAR100	TMC (Yan et al. 2025)	✗	✗	ResNet-19	6 / 4 / 2	78.05 / 77.52 / 76.35
	SNN-ViT (Wang et al. 2025)	✗	✗	VGG-16	4	80.01
	LT-SNN (Hasssan, Meng, and Seo 2024)	✓	✗	ResNet-19	6 / 2	74.82 / 72.78
	ALSF (Fu et al. 2025)	✓	✓	ResNet-19	2	78.73
	DeepTAGE (Liu et al. 2025)	✗	✓	ResNet-18	4	78.80
	MPD-AGL (Jiang et al. 2025)	✗	✓	ResNet-19	4 / 2	79.72 / 78.84
	AT-LIF&ASG-S (Ai et al. 2025)	✓	✓	ResNet-18	6 / 4	76.44 / 76.42
Ours	✓	✓	ResNet-19	2	80.59±0.17	
CIFAR10-DVS	FSTA-SNN (Yu et al. 2025b)	✗	✗	ResNet-20	10	81.50
	QP-SNN (Wei et al. 2025)	✗	✗	VGG-SNN	10	82.10
	STL-SNN (Sun et al. 2024)	✓	✗	7-layers SNN	20	77.30
	LT-SNN (Hasssan, Meng, and Seo 2024)	✓	✗	VGG-9	30	80.07
	DeepTAGE (Liu et al. 2025)	✗	✓	VGG-11	10	81.23
	MPD-AGL (Jiang et al. 2025)	✗	✓	VGG-SNN	10	82.50
	AT-LIF&ASG-S (Ai et al. 2025)	✓	✓	VGG-SNN	10	78.70
Ours	✓	✓	VGG-SNN	10	83.70±0.41	
ImageNet	AGMM (Liang et al. 2025)	✗	✗	ResNet-18	4	64.67
	ReverB (Guo et al. 2025)	✗	✗	ResNet-18	4	66.22
	MTT (Du et al. 2025)	✗	✗	ResNet-34	4	67.54
	LTF&ALSF (Fu et al. 2025)	✓	✗	SEW ResNet-18	4	63.67
	DeepTAGE (Liu et al. 2025)	✗	✓	ResNet-18	4	68.52
	Ours	✓	✓	ResNet-18	4	68.86±0.25

Table 1: Comparison of classification accuracy. **AT** and **GO** denote adaptive threshold and gradient optimization, respectively.

from membrane potential shifts limit performance. Enhancing SNNs requires considering both adaptive regulation of neuronal firing thresholds and dynamic optimization of gradient learning. In particular, as analyzed in ablation studies, most methods related to gradient optimization are superior to methods for threshold adaptation. Although AT-LIF&ASG-S also achieves threshold and SG adjustment, it ignores the correlation between them, lacking synergistic learning. In summary, our method combines these two techniques to develop a two-stage synergistic learning that achieves higher accuracy with low latency, demonstrating weak-to-strong generalization and computational efficiency.

Spike Firing Rates

To validate whether our method can effectively balance firing rates, we conducted experiments on the CIFAR10 dataset. In Fig. 8, Vanilla-SNN exhibits a low firing rate (average 10.66%) with large fluctuations between layers, which hinders adequate encoding of input information. DIET-SNN (Rathi and Roy 2023) and LTMD (Wang, Cheng, and Lim 2022) methods substantially increase the firing level of neurons by adapting thresholds, avoiding extensive information loss. However, their firing rates across layers exhibit even greater oscillations, and the high firing rates certainly increase energy consumption. Instead, our firing rates remain within a small fluctuation range of nearly $15 \pm 1.62\%$. This stability stems from DS-ATGO’s adaptive adjustment of

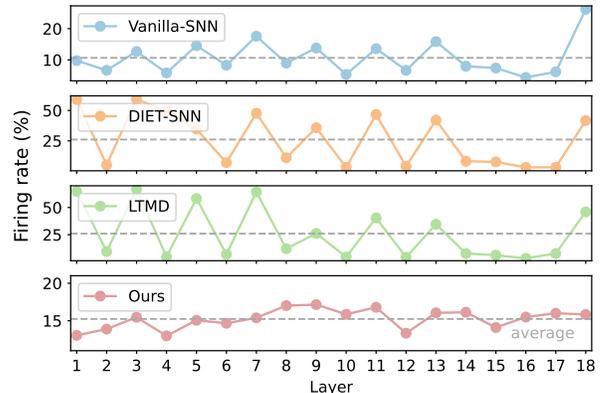


Figure 8: Comparison of firing rates for each layer in ResNet-19, where grey dotted lines denote average values.

thresholds based on the MPD distribution at each timestep, which helps neurons maintain moderate activation. Moreover, the firing rates of the last layer in the other three methods exhibit a steep increase. This may be due to the accumulated deviation of membrane potentials from thresholds over time, resulting in a slow decrease in firing rates. To compensate for the diminishing useful features in the hidden layers, output layer neurons are forced to fire more spikes by dynamically lowering thresholds or enhancing weights in an

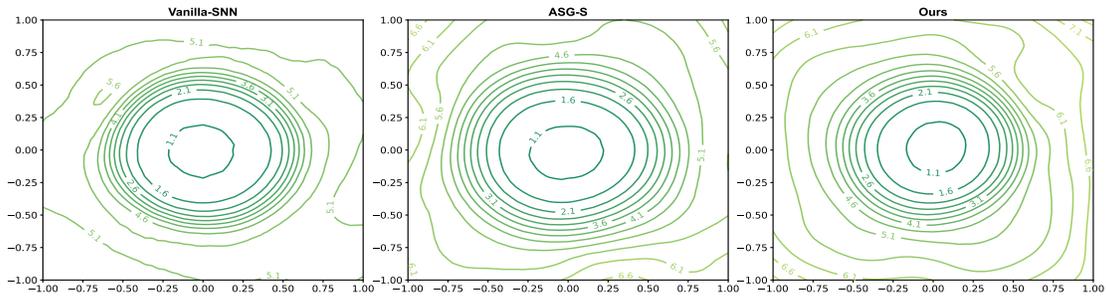


Figure 9: The 2D loss landscape of ResNet-19 trained with different methods on the CIFAR100 dataset.

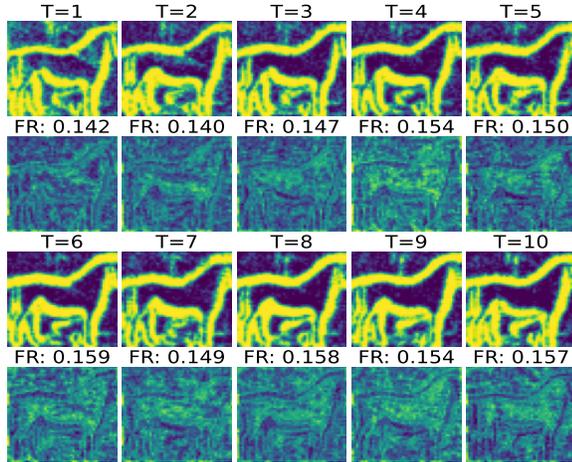


Figure 10: Case study of firing rates (FR) at each timestep.

attempt to reconstruct complete feature representations. Although this temporarily enhances the excitability of output neurons, it causes neurons to over-respond to non-sharp signals, which is not conducive to recognition. Then, we presented the original image (rows 1, 3) and feature maps (rows 2, 4) on the CIFAR10-DVS dataset. In Fig. 10, our method effectively maintains the firing rate stable over timesteps and extracts good feature maps.

Proportion of Gradient Available

To investigate whether our method can effectively mitigate the gradient vanishing problem, we conducted experiments on the CIFAR100 dataset. In Fig. 11, the fixed SG of Vanilla-SNN cannot dynamically match the membrane potential shifts, resulting in only a few neurons in each layer falling within the gradient-available interval. As the number of layers deepens, this exacerbates the loss of gradient information, causing the proportion of neurons that obtain gradients in the last three layers to drop to an average of less than 15.13%. Although ASG-S (Ai et al. 2025) slightly increases the proportion by adjusting SG, it still fails to accurately capture MPD that evolves with timesteps due to the lack of temporal dependence. By comparison, our method increases the gradient-available rate of each layer in ResNet-19 to over 38.53%, and even maintains 36.54% in the deeper layers.

Then, we visualized the 2D loss landscapes of these meth-

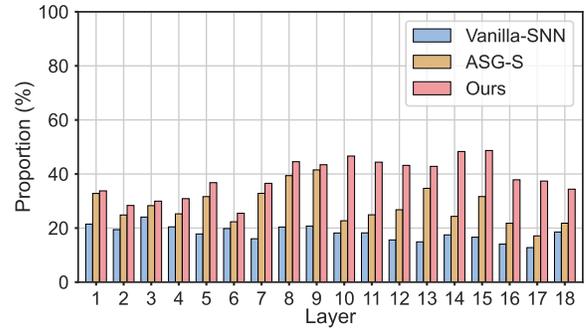


Figure 11: The proportion of neurons falling within the gradient-available interval of each layer in ResNet-19.

ods. In Fig. 9, the loss surface of Vanilla-SNN exhibits two local minima, with obvious protrusions in the contour lines around the center point, reflecting oscillations in weight updates caused by gradient loss. Although the loss landscape of ASG-S contains only one local minimum, its contour lines form an elliptical region. In contrast, the loss surface of our method exhibits a uniformly continuous downward trend in the parameter space through synergistic optimization, generating a sparser and flatter loss landscape.

Conclusion

In this paper, we present a new perspective on the performance limitations of SNNs by analyzing shifts in membrane potential distributions and deviations from the threshold. Fixed thresholds and SG are inherently mismatched with evolving MPD across timesteps, leading to imbalanced firing and diminished gradient problems. Then, we propose the DS-ATGO learning algorithm, which responds to evolving MPD across both temporal and spatial dimensions via forward adaptive thresholding and backward SG optimization. This dual-stage synergistic learning framework aims to help SNNs break through the performance bottlenecks of single-path techniques. Experimental results and analyses demonstrate that DS-ATGO outperforms methods that only adjust thresholds or optimize SG. By balancing the firing rate and enhancing the gradient signal across timesteps, our method improves information encoding efficiency while optimizing loss landscapes, potentially advancing the application of SNNs to more complex tasks and wider scenarios.

Acknowledgements

This work was supported by the National Key R&D Program of China (Grant No. 2025YFG0100700) and the National Natural Science Foundation of China (Grant No. 62276235).

References

- Ai, Q.; Yang, Y.; Cai, M.; Chen, K.; Liu, Q.; and Ma, L. 2025. A cross-layer residual spiking neural network with adaptive threshold leaky integrate-and-fire neuron and learnable surrogate gradient. *Knowledge-Based Systems*, 319: 113575.
- Azouz, R.; and Gray, C. M. 2000. Dynamic spike threshold reveals a mechanism for synaptic coincidence detection in cortical neurons in vivo. *Proceedings of the National Academy of Sciences*, 97(14): 8110–8115.
- Azouz, R.; and Gray, C. M. 2003. Adaptive coincidence detection and dynamic gain control in visual cortical neurons in vivo. *Neuron*, 37(3): 513–523.
- Che, K.; Leng, L.; Zhang, K.; Zhang, J.; Meng, Q.; Cheng, J.; Guo, Q.; and Liao, J. 2022. Differentiable hierarchical and surrogate gradient search for spiking neural networks. *Advances in Neural Information Processing Systems*, 35: 24975–24990.
- Cubuk, E. D.; Zoph, B.; Mane, D.; Vasudevan, V.; and Le, Q. V. 2019. AutoAugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 113–123.
- Deng, J.; Dong, W.; Socher, R.; Li, L.-J.; Li, K.; and Fei-Fei, L. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, 248–255. Ieee.
- Deng, S.; Li, Y.; Zhang, S.; and Gu, S. 2022. Temporal efficient training of spiking neural network via gradient reweighting. In *International Conference on Learning Representations*.
- Ding, J.; Dong, B.; Heide, F.; Ding, Y.; Zhou, Y.; Yin, B.; and Yang, X. 2022. Biologically inspired dynamic thresholds for spiking neural networks. *Advances in neural information processing systems*, 35: 6090–6103.
- Du, K.; Wu, Y.; Deng, S.; and Gu, S. 2025. Temporal flexibility in spiking neural networks: towards generalization across time steps and deployment friendliness. In *The Thirteenth International Conference on Learning Representations*.
- Fang, W.; Yu, Z.; Chen, Y.; Huang, T.; Masquelier, T.; and Tian, Y. 2021a. Deep residual learning in spiking neural networks. *Advances in Neural Information Processing Systems*, 34: 21056–21069.
- Fang, W.; Yu, Z.; Chen, Y.; Masquelier, T.; Huang, T.; and Tian, Y. 2021b. Incorporating learnable membrane time constant to enhance learning of spiking neural networks. In *Proceedings of the IEEE/CVF international conference on computer vision*, 2661–2671.
- Farries, M. A.; Kita, H.; and Wilson, C. J. 2010. Dynamic spike threshold and zero membrane slope conductance shape the response of subthalamic neurons to cortical input. *Journal of Neuroscience*, 30(39): 13180–13191.
- Feng, X.; Zhang, H.; Zhang, Y.; Zhang, L. Y.; and Pan, S. 2025. BiMark: Unbiased Multilayer Watermarking for Large Language Models. In *Forty-second International Conference on Machine Learning*.
- Fontaine, B.; Peña, J. L.; and Brette, R. 2014. Spike-threshold adaptation predicted by membrane potential dynamics in vivo. *PLoS computational biology*, 10(4): e1003560.
- Fu, J.; Gou, S.; Wang, P.; Jiao, L.; Guo, Z.; Li, J.; and Liu, R. 2025. Adaptation and learning of spatio-temporal thresholds in spiking neural networks. *Neurocomputing*, 130423.
- Guo, Y.; Chen, Y.; Zhang, L.; Liu, X.; Wang, Y.; Huang, X.; and Ma, Z. 2022a. IM-Loss: Information maximization loss for spiking neural networks. *Advances in Neural Information Processing Systems*, 35: 156–166.
- Guo, Y.; Tong, X.; Chen, Y.; Zhang, L.; Liu, X.; Ma, Z.; and Huang, X. 2022b. RecDis-SNN: Rectifying membrane potential distribution for directly training spiking neural networks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 326–335.
- Guo, Y.; Zhang, Y.; Jie, Z.; Liu, X.; Tong, X.; Chen, Y.; Peng, W.; and Ma, Z. 2025. ReverB-SNN: Reversing bit of the weight and activation for spiking neural networks. In *Forty-second International Conference on Machine Learning*.
- Hasssan, A.; Meng, J.; and Seo, J.-S. 2024. Spiking neural network with learnable threshold for event-based classification and object detection. In *2024 International Joint Conference on Neural Networks (IJCNN)*, 1–8.
- 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)*, 10–14. IEEE.
- Hu, Y.; Deng, L.; Wu, Y.; Yao, M.; and Li, G. 2024. Advancing spiking neural networks toward deep residual learning. *IEEE transactions on neural networks and learning systems*, 36(2): 2353–2367.
- Ioffe, S.; and Szegedy, C. 2015. Batch normalization: accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, 448–456. pmlr.
- Jiang, J.; Wang, L.; Jiang, R.; Fan, J.; and Yan, R. 2025. Adaptive gradient learning for spiking neural networks by exploiting membrane potential dynamics. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*.
- Ju, J.; Zheng, Y.; Koh, H. Y.; and Pan, S. 2025a. Uni-MRL: Unified MultiModal Molecular Representation Learning with Large Language Models and Graph Neural Networks. In *Pacific-Asia Conference on Knowledge Discovery and Data Mining*, 275–287. Springer.
- Ju, J.; Zheng, Y.; Koh, H. Y.; Wang, C.; and Pan, S. 2025b. M2LLM: Multi-view Molecular Representation Learning with Large Language Models. In *IJCAI-25*, 7437–7445. International Joint Conferences on Artificial Intelligence Organization.

- Krizhevsky, A.; Hinton, G.; et al. 2009. Learning multiple layers of features from tiny images.
- Lee, D.; Li, Y.; Kim, Y.; Xiao, S.; and Panda, P. 2025. Spiking transformer with spatial-temporal attention. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, 13948–13958.
- Li, H.; Liu, H.; Ji, X.; Li, G.; and Shi, L. 2017. CIFAR10-DVS: An event-stream dataset for object classification. *Frontiers in neuroscience*, 11: 244131.
- Li, Y.; Kim, Y.; Park, H.; Geller, T.; and Panda, P. 2022. Neuromorphic data augmentation for training spiking neural networks. In *European Conference on Computer Vision*, 631–649. Springer.
- Lian, S.; Shen, J.; Liu, Q.; Wang, Z.; Yan, R.; and Tang, H. 2023. Learnable surrogate gradient for direct training spiking neural networks. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence*, 3002–3010.
- Liang, Y.; Wei, W.; Belatreche, A.; Cao, H.; Zhou, Z.; Wang, S.; Zhang, M.; and Yang, Y. 2025. Towards accurate binary spiking neural networks: learning with adaptive gradient modulation mechanism. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 1402–1410.
- Liu, W.; Yang, L.; Zhao, M.; Wang, S.; Gao, J.; Li, W.; Li, B.; and Hu, W. 2025. DeepTAGE: Deep temporal-aligned gradient enhancement for optimizing spiking neural networks. In *The Thirteenth International Conference on Learning Representations*.
- Maass, W. 1997. Networks of spiking neurons: the third generation of neural network models. *Neural networks*, 10(9): 1659–1671.
- Pena, J. L.; and Konishi, M. 2002. From postsynaptic potentials to spikes in the genesis of auditory spatial receptive fields. *Journal of Neuroscience*, 22(13): 5652–5658.
- Rathi, N.; and Roy, K. 2023. DIET-SNN: A low-latency spiking neural network with direct input encoding and leakage and threshold optimization. *IEEE Transactions on Neural Networks and Learning Systems*, 34(6): 3174–3182.
- Roy, K.; Jaiswal, A.; and Panda, P. 2019. Towards spike-based machine intelligence with neuromorphic computing. *Nature*, 575(7784): 607–617.
- Samadzadeh, A.; Far, F. S. T.; Javadi, A.; Nickabadi, A.; and Chehreghani, M. H. 2023. Convolutional spiking neural networks for spatio-temporal feature extraction. *Neural Processing Letters*, 55(6): 6979–6995.
- Sun, H.; Cai, W.; Yang, B.; Cui, Y.; Xia, Y.; Yao, D.; and Guo, D. 2024. A synapse-threshold synergistic learning approach for spiking neural networks. *IEEE Transactions on Cognitive and Developmental Systems*, 16(2): 544–558.
- Wang, S.; Cheng, T. H.; and Lim, M.-H. 2022. LTMD: Learning improvement of spiking neural networks with learnable thresholding neurons and moderate dropout. *Advances in Neural Information Processing Systems*, 35: 28350–28362.
- Wang, S.; Cheng, T. H.; and Lim, M.-H. 2025. Potential distribution adjustment and parametric surrogate gradient in spiking neural networks. *Neurocomputing*, 620: 129189.
- Wang, S.; Zhang, M.; Zhang, D.; Belatreche, A.; Xiao, Y.; Liang, Y.; Shan, Y.; Sun, Q.; Zhang, E.; and Yang, Y. 2025. Spiking vision transformer with saccadic attention. In *The Thirteenth International Conference on Learning Representations*.
- Wang, Z.; Jiang, R.; Lian, S.; Yan, R.; and Tang, H. 2023. Adaptive smoothing gradient learning for spiking neural networks. In *International conference on machine learning*, 35798–35816. PMLR.
- Wei, W.; Zhang, M.; Zhou, Z.; Belatreche, A.; Shan, Y.; Liang, Y.; Cao, H.; Zhang, J.; and Yang, Y. 2025. QP-SNN: Quantized and pruned spiking neural networks. In *The Thirteenth International Conference on Learning Representations*.
- Wu, Y.; Deng, L.; Li, G.; Zhu, J.; and Shi, L. 2018. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in neuroscience*, 12: 331.
- Wu, Y.; Deng, L.; Li, G.; Zhu, J.; Xie, Y.; and Shi, L. 2019. Direct training for spiking neural networks: faster, larger, better. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, 1311–1318.
- Yan, J.; Wang, C.; Ma, D.; Tang, H.; Zheng, Q.; and Pan, G. 2025. Training high performance spiking neural network by temporal model calibration. In *Forty-second International Conference on Machine Learning*.
- Yu, C.; Zhao, X.; Liu, L.; Yang, S.; Wang, G.; Li, E.; and Wang, A. 2025a. Efficient logit-based knowledge distillation of deep spiking neural networks for full-range timestep deployment. In *Forty-second International Conference on Machine Learning*.
- Yu, K.; Zhang, T.; Wang, H.; and Xu, Q. 2025b. FSTA-SNN: Frequency-based spatial-temporal attention module for spiking neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, 22227–22235.
- Zhang, H.; Wu, B.; Yang, X.; Yuan, X.; Liu, X.; and Yi, X. 2025. Dynamic graph unlearning: a general and efficient post-processing method via gradient transformation. In *Proceedings of the ACM on Web Conference 2025*, 931–944.
- Zhang, H.; Yuan, X.; and Pan, S. 2024. Unraveling privacy risks of individual fairness in graph neural networks. In *2024 IEEE 40th International Conference on Data Engineering (ICDE)*, 1712–1725. IEEE.
- Zheng, H.; Wu, Y.; Deng, L.; Hu, Y.; and Li, G. 2021. Going deeper with directly-trained larger spiking neural networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, 11062–11070.

Appendix of DS-ATGO: Dual-Stage Synergistic Learning via Forward Adaptive Threshold and Backward Gradient Optimization for Spiking Neural Networks

A Proofs of Theorems

Theorem 1. For $U(t)$ that satisfies a normal distribution $N(\mu, \sigma^2)$, the probability that a random variable $U_i(t)$ exceeds $\mu + \sigma$ is relatively constant and given by $P(U_i(t) > \mu + \sigma) = 1 - \Phi(1)$, where $\Phi(\cdot)$ denotes the cumulative distribution function of standard normal distributions.

Proof. The probability density function of the normal distribution $U(t) \sim N(\mu, \sigma^2)$ can be expressed by

$$f(U_i(t)) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(U_i(t)-\mu)^2}{2\sigma^2}}, U_i(t) \in (-\infty, +\infty). \quad (12)$$

To calculate the probability $P(U_i(t) > \mu + \sigma)$ that $U_i(t)$ is greater than $\mu + \sigma$, which is

$$P(U_i(t) > \mu + \sigma) = \int_{\mu+\sigma}^{+\infty} \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(U_i(t)-\mu)^2}{2\sigma^2}} dU_i(t). \quad (13)$$

Let $U(t)$ be normalized to $Z(t) = \frac{U(t)-\mu}{\sigma}$, we can obtain the standard normal distribution $Z(t) \sim N(0, 1)$. Then, it transforms the calculation $P(U_i(t) > \mu + \sigma)$ into $P(Z_i(t) > 1)$, as follows:

$$P(Z_i(t) > 1) = \int_1^{+\infty} \frac{1}{\sqrt{2\pi}} e^{-\frac{Z_i(t)^2}{2}} dZ_i(t). \quad (14)$$

Based on the properties of standard normal distributions, $P(Z_i(t) > 1) = 1 - P(Z_i(t) \leq 1) = 1 - \Phi(1)$, where $\Phi(\cdot)$ denotes the cumulative distribution function. $\Phi(1)$ is a relatively constant value that does not depend on μ and σ . \square

B Overall procedure of DS-ATGO algorithm

C Supplementary Experiments

Energy Efficiency and Ablation Study of Firing Control f_c

We show the spike firing rates of different methods for each ResNet-19 layer in the *Spiking Firing Rates* Section. Then, we estimate the energy consumption of ANN and various SNN methods ($T = 2$). Following (Rathi and Roy 2023), the number of operations in SNN is specified as $FR^l \times T \times OP_{ANN}^l$, where FR^l denotes the firing rate of the l -layer, T denotes the timesteps and OP_{ANN}^l denotes the number of operations at the l -th layer in the iso-architecture ANN. It is worth mentioning that, as the original images are directly fed into SNN for encoding and the membrane potentials of the output layer are used for prediction in our model, the number of operations for these two layers is specified as the number of operations in the corresponding layers of the iso-architecture ANN and scaled by T . (Horowitz 2014)

Algorithm 1: The overall procedure of SNNs in one epoch with DS-ATGO algorithm.

Input: Initial SNN model and hyperparameters; total train and test iteration N_{train}, N_{test} ; class label vector: Y .

Output: The well-trained SNN and classification accuracy.

Training:

```

1: for  $n = 1$  to  $N_{train}$  do
2:   for  $l = 1$  to  $L - 1$  do
3:     Compute  $I^l, U^l$  // (1) and (2)
4:     for  $t = 1$  to  $T$  do
5:       Calculate  $\Delta V_{th}^l(t)_n, \Delta V_{th}^l(t)$  // (5) and (6)
6:       Compute  $S^l(t) \leftarrow U^l(t) \geq \Delta V_{th}^l(t)_n$  // (3)
7:       Calculate the width of SG  $k^l(t)$  // (10)
8:     end for
9:   end for
10:  Compute the output  $O_n^L$  and loss  $L_{CE}$  // (11)
11:  Compute the gradient w.r.t.  $w$  // (13-15)
12:  Update  $w \leftarrow w - \eta \frac{\partial L}{\partial w}$ , where  $\eta$  is learning rate
13: end for

```

Inference:

```

1: for  $n = 1$  to  $N_{test}$  do
2:   for  $l = 1$  to  $L - 1$  do
3:     Compute  $I^l$  and  $U^l$  // (1) and (2)
4:     Compute  $S^l \leftarrow U^l \geq \Delta V_{th}^l$  // (3)
5:   end for
6:   Compute  $O^L = \frac{1}{T} \sum_{t=1}^T (w^L \otimes S^{L-1})$  // (11)
7:   Compute accuracy from outputs  $O^L$  and labels  $Y$ 
8: end for

```

Methods	AC	MAC	Energy	Train	Test
ANN	0M	2285.35M	10.51mJ	38s	2s
Vanilla-SNN	443.14M	7.08M	0.43mJ	67s	4s
DIET-SNN	1084.02M	7.08M	1.01mJ	69s	4s
LTMD	1047.98M	7.08M	0.98mJ	68s	4s
Ours ($f_c = 1.0$)	690.84M	7.08M	0.65mJ	74s	4s
Ours ($f_c = 1.2$)	493.32M	7.08M	0.48mJ	74s	4s

Table 2: The energy consumption and running time of ANN and various SNN methods ($T = 2$) on the CIFAR10 dataset.

measured in 45-nm CMOS technology that an AC operation costs $0.9pJ$ and a MAC operation costs $4.6pJ$.

To investigate the impact of the firing rate control factor f_c on the trade-off between performance and energy consumption, we present a comparison of performance and energy consumption under different values of f_c , as shown in Fig. 12. As listed in Table 2, the energy consumption of our

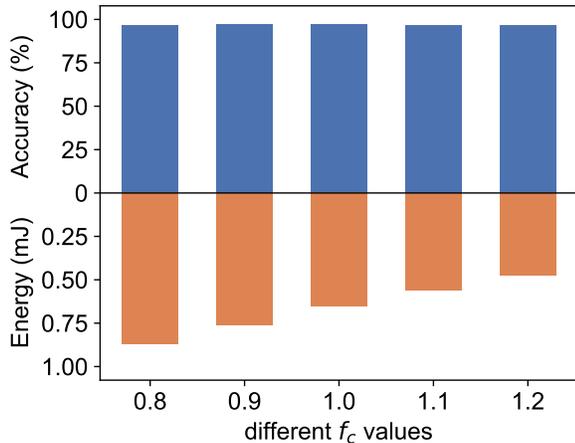


Figure 12: The trade-off between performance and energy consumption with different values of f_c .

method is approximately $0.65mJ$ when the firing control factor $f_c = 1.2$, much lower than the $10.51mJ$ of ANN (around $16\times$). When $f_c = 1.2$, our energy consumption is comparable to Vanilla-SNN, while performance still improves by 1.06%. Moreover, we recorded the running time of these methods. As we need to compute the thresholds and widths of SG for each timestep in each layer, the training time of our method is roughly 74s per epoch, which is 7 seconds slower than Vanilla-SNN. However, the inference time of our method is the same as that of the Vanilla-SNN, meaning that it does not introduce additional inference overhead.

Ablation Study of Momentum Coefficient m

In the *Ablation Study* Section, we demonstrate that applying the moving average rule in the adaptive threshold mechanism can effectively improve the inference performance. Here, we further investigated the impact of the momentum coefficient on the moving average rule. As shown in Fig. 13 (bar chart), different values of the momentum coefficient exhibit negligible influence on the accuracy of the inference phase, but all outperform the baseline without the moving average rule. To explore the underlying reason for this phenomenon, we also present a comparison of the moving average thresholds (average of two timesteps) used for inference, obtained after training with different momentum coefficients. Fig. 13 (green line chart) clearly shows that the moving average thresholds computed under different momentum factors exhibit a high degree of consistency. Finally, we plot the changes in the thresholds (Eq. 5) of each layer in ResNet-19 during training (Fig. 14), where the momentum coefficient and f_c are set to 0.1 and 1.0, respectively.

Comparison of t-SNE

In this section, we perform a t-SNE visualization on the feature vectors extracted from the last convolutional layer of ResNet-19. As shown in Fig. 15, our method clearly distinguishes feature information of different categories. In particular, it can effectively separate similar categories such as ‘truck’ and ‘automobile’, while there is only partial overlap

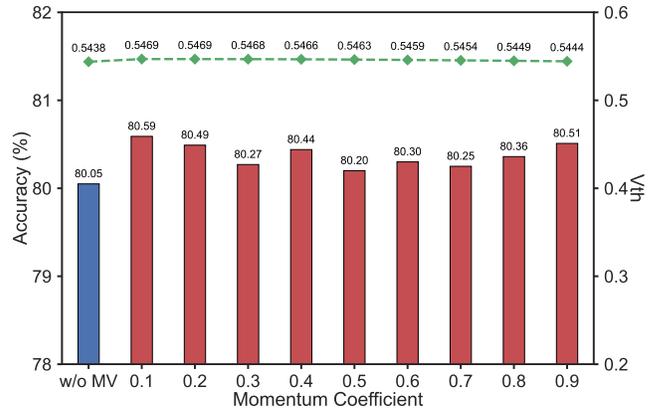


Figure 13: The impact of momentum coefficients in the moving average (MV) rule (bar chart) and a comparison of the moving average thresholds used for inference, obtained after training on the CIFAR100 dataset with different momentum coefficients (line chart).

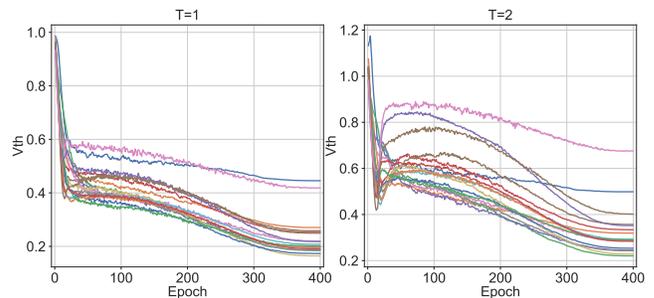


Figure 14: The changes in the thresholds (Eq. 5) of each layer in ResNet-19 during training on the CIFAR100 dataset.

between the categories ‘dog’ and ‘cat’. However, Vanilla-SNN performs poorly in distinguishing features between multiple animal categories (such as ‘dog’, ‘cat’, ‘deer’, and ‘bird’). This indicates that our method has stronger feature extraction capabilities than Vanilla-SNN and effectively captures the crucial information of different categories, which is beneficial for classification.

D Implementation Details

Environment and Hyperparameter Settings

All experiments are performed on a workstation equipped with Ubuntu 20.04.5 LTS, an AMD Ryzen Threadripper 3960X CPU running at 2.20 GHz, 128 GB RAM, and four NVIDIA GeForce RTX 4090 GPUs with 24GB DRAM. The code is implemented in the PyTorch framework with Python 3.9, and the weights are initialized using the default Kaiming initialization method in PyTorch 1.12.1.

Table 3 lists the hyperparameters used in our work. SGD optimizer with an initial learning rate $\eta = 0.1$, 0.9 momentum, and weight decay $1e-4$ is used in all four datasets. All experiments used the CosineAnnealingLR scheduler to adjust η , which will cosine decay to 0 over epochs. The classification accuracy on the CIFAR10/100 and CIFAR10-

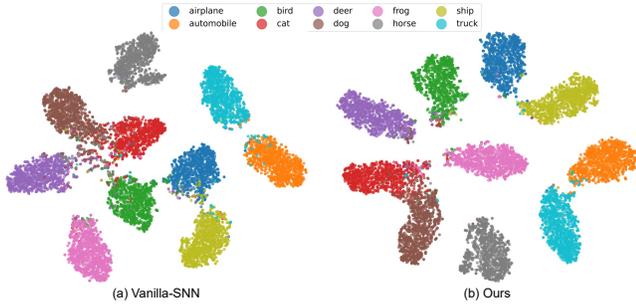


Figure 15: t-SNE visualization of feature vectors using ResNet-19 on the CIFAR10 dataset.

Hyperparameters	CIFAR10	CIFAR100	ImageNet	CIFAR10-DVS
V_{th}	1.0	1.0	1.0	1.0
τ	0.2	0.2	0.2	0.2
Epoch	400	400	400	400
Batch Size	100	100	32	50
Optimizer	SGD	SGD	SGD	SGD
Weight Decay	$1e^{-4}$	$1e^{-4}$	$1e^{-4}$	$1e^{-4}$
η	0.1	0.1	0.1	0.1

Table 3: Hyperparameter Settings.

DVS datasets are obtained by running with five random seeds, whereas that on ImageNet dataset is obtained using three random seeds. For the CIFAR10/100 and ImageNet datasets, the network loss L is calculated using the cross-entropy function with label smoothing (Hu et al. 2024). For the CIFAR10-DVS dataset, the network loss L is calculated using TET loss (Deng et al. 2022).

Datasets and Preprocessing

CIFAR-10: The CIFAR-10 dataset (Krizhevsky, Hinton et al. 2009) consists of 60,000 RGB static images across 10 classes, each with a 32×32 pixels resolution. These images are split into 50,000 for training and 10,000 for testing. In data preprocessing, we normalized the dataset by subtracting the global mean value of pixel intensity and dividing by the standard variance of RGB channels. Each image was randomly cropped to 32×32 pixels after padding 4 pixels and randomly horizontally flipped with 0.5 probability. AutoAugment (Cubuk et al. 2019) was used for data augmentation.

CIFAR-100: The CIFAR-100 dataset (Krizhevsky, Hinton et al. 2009) also contains 60,000 RGB static images with a resolution of 32×32 pixels in 100 classes, which are split into 50,000 training images and 10,000 test images. We adopt the same preprocessing and data augmentation strategy to the CIFAR-100 dataset as the CIFAR-10 dataset.

CIFAR10-DVS: The CIFAR10-DVS dataset (Li et al. 2017) is converted from 10,000 CIFAR10 images and is the most challenging mainstream neuromorphic dataset. It consists of 10 classes, each with 1,000 samples and a resolution of 128×128 pixels, but was not split into training and test sets. Following previous studies (Zheng et al. 2021; Deng et al. 2022; Samadzadeh et al. 2023), we also used 90% of the samples in each class for training and the other 10%

for testing. In data preprocessing, we reduced the temporal resolution by segmenting the event stream into 10 temporal blocks, accumulating spike events within each block, and resizing the spatial resolution to 48×48 (Li et al. 2022; Lian et al. 2023). Horizontal flipping with a random probability greater than 0.5 and rolling with 5 offsets were applied for data augmentation (Deng et al. 2022).

ImageNet: The ImageNet-1k dataset (Deng et al. 2009) is more challenging than static CIFAR family datasets. It consists of 1250k RGB static images for training and 50k RGB static images for testing across 1000 classes. In data preprocessing, we normalized the dataset by subtracting the global mean value of pixel intensity and dividing by the standard variance of RGB channels. Each image was randomly cropped to 224×224 pixels and randomly horizontally flipped with 0.5 probability. AutoAugment was used for data augmentation.

Network Architectures

We adopt the widely-used ResNet-18/19 (Zheng et al. 2021; Fang et al. 2021a; Hu et al. 2024) and VGGSNN (Deng et al. 2022) network structures. The details of the network architecture are listed in Table 4 and Table 5, respectively. xCy represents a convolutional layer with output channels equal to x , $kernel\ size = y$, $stride\ set = 1$, and $padding = 1$. xFC represents a fully-connected layer with output features equal to x . $2AP$ represents the average-pooling layer with $kernel\ size = 2$ and $stride = 2$. z is the number of classes.

conv1	128C3
block1	$\begin{pmatrix} 128C3 \\ 128C3 \end{pmatrix} \times 3$
block2	$\begin{pmatrix} 256C3 \\ 256C3 \end{pmatrix}^* \times 3$
block3	$\begin{pmatrix} 512C3 \\ 512C3 \end{pmatrix}^* \times 3$
ResNet-18	average pool, zFC
ResNet-19	average pool, 256-d FC zFC

* means the first basic block in the series performs downsampling directly with convolution kernels and a stride of 2.

Table 4: ResNet-18/19 structures.

	64C3-LIF
	128C3-LIF-2AP
VGGSNN	256C3-LIF-256C3-LIF-2AP
	512C3-LIF-512C3-LIF-2AP
	512C3-LIF-512C3-LIF-2AP - zFC

Table 5: VGGSNN structures.