# Revisiting Direct Encoding: Learnable Temporal Dynamics for Static Image Spiking Neural Networks

Huaxu He[a,b,*]

[a]*School of Computer and Information Engineering, Henan University, Kaifeng, 475004, Henan, China*
[b] *Henan Key Laboratory of Big Data Analysis and Processing, Henan University, Kaifeng, 475004, Henan, China*

**Abstract**

Handling static images that lack inherent temporal dynamics remains a fundamental challenge for spiking neural networks (SNNs). In directly trained SNNs, static inputs are typically repeated across time steps, causing the temporal dimension to collapse into a rate like representation and preventing meaningful temporal modeling. This work revisits the reported performance gap between direct and rate based encodings and shows that it primarily stems from convolutional learnability and surrogate gradient formulations rather than the encoding schemes themselves. To illustrate this mechanism level clarification, we introduce a minimal learnable temporal encoding that adds adaptive phase shifts to induce meaningful temporal variation from static inputs.

*Keywords:* Spiking neural network, Rate coding,Surrogate gradient, Learnable temporal encoding

## 1. Introduction

Spiking neural networks (SNNs), with event driven sparsity and spatiotemporal processing, offer brain inspired and energy efficient computation [1, 2, 3, 4, 5]. However, in static image tasks, existing SNNs still rely on approximate schemes such as rate coding, leaving their temporal modeling capability underutilized. Achieving meaningful temporal dynamics under low time-step constraints remains challenging [6].

---

*Corresponding author: Huaxu He (email: 104753230975@henu.edu.cn)

Within current SNN research, training is broadly divided into ANN-to-SNN conversion and direct training [7, 8, 9, 10]. Conversion methods approximate ANN activations with firing rates and typically require many time steps [11]. while direct training allows end-to-end learning with just a few steps [12]. Among direct methods, direct encoding is the most representative, enabling efficient inference at minimal time steps.

We revisit the widely adopted direct encoding mechanism in directly trained SNNs. Although direct encoding enables efficient inference at minimal time steps—demonstrating a clear advantage in computational efficiency—it effectively repeats identical inputs across time steps, causing the temporal dimension to lose its genuine dynamic evolution [13].

Specifically, in directly trained SNNs, static images are typically repeated along the temporal dimension and fed into leaky integrate-and-fire (LIF) neurons to construct a formal time sequence [14]. However, since the inputs at all time steps are identical, the temporal dimension functions merely as a statistical averaging channel rather than a true temporal modeling mechanism. For LIF neurons under constant input current, the spiking behavior becomes simple and stable, implying that the overall network effectively performs rate-based coding [15], A minimal illustrative example of this equivalence is provided in the Appendix B.
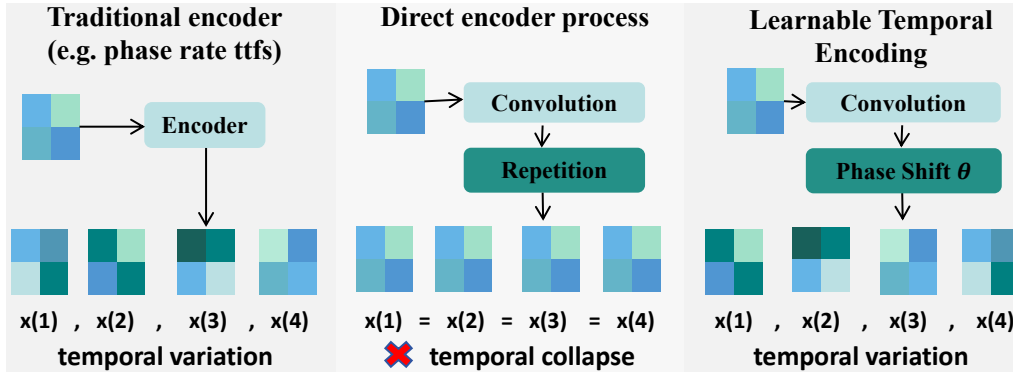


Figure 1: Comparison of traditional encoders, direct encoding, and the proposed learnable temporal encoding. Traditional encoders produce time-varying inputs via fixed rules, while direct encoding replicates static inputs, causing temporal collapse. The proposed method restores temporal variation through learnable phase shifts. Here, $x$ represents the input feature maps at each time step.

To further analyze the performance discrepancy between direct and rate coding reported in previous studies, we show that the difference mainly

arises from the encoding target: rate coding is applied directly to raw inputs, whereas direct encoding operates on static feature maps extracted by convolutional layers. Empirically, the two exhibit highly similar performance in static tasks. The previously reported performance gap is thus largely due to implementation details such as feature selection and surrogate gradient configurations rather than intrinsic differences in their coding principles.

Neurophysiological studies indicate that biological neurons encode information not only through firing rates but also through temporal features like spike timing. Motivated by this, we introduce a learnable temporal encoding that adds trainable phase parameters at the input stage. These parameters induce phase shifts across time steps, forming distinguishable temporal patterns while preserving the structural simplicity of direct encoding. This mechanism equips the network with temporal correlations for learning discriminative representations, as illustrated in Figure 1. The main contributions of this work are summarized as follows:

- We identify the functional equivalence between direct and rate encoding in static image SNNs and show that their reported performance gap mainly results from convolutional learnability and surrogate gradient.

- We introduce a minimal learnable temporal encoding for static inputs, using trainable phase parameters to induce temporal variation with minimal structural overhead.

- We validate this design on object recognition and detection tasks using the CIFAR-10, CIFAR-100, and VOC datasets, showing improved performance under low time-step constraints.

## 2. Related Work

### 2.1. SNN Encoding Schemes

The encoding problem in SNNs can be categorized into two types:encoding for static images in directly trained SNNs [16], and feature encoding for intermediate layers in ANN-to-SNN conversion.This work focuses on the former.

In directly trained SNNs, the predominant strategy is direct encoding, which replicates static feature maps along the temporal dimension and feeds them into LIF neurons, producing spike sequences that formally exhibit temporal structure. This approach is simple, stable during training, and performs well under very few time steps.

3

Other commonly used schemes include rate encoding, time-to-first-spike (TTFS), and phase encoding. Rate encoding conveys signal magnitude through firing rates, whereas TTFS and phase encoding rely on spike timing, with TTFS enforcing a single spike per neuron [17, 18, 19].

We show that direct encoding is functionally equivalent to rate encoding, with the key difference being the encoding target: rate encoding operates directly on raw images, whereas direct encoding applies to static feature maps produced by convolutional layers. Based on this observation, we further introduce phase encoding and TTFS encoding to enhance the temporal representation capability of SNNs.

## 2.2. Spiking Neurons

The LIF neuron model with soft reset, widely used in current SNNs [20, 21], It can be described as follows:

$$H[t] = \mathcal{L}V[t-1] + (1-\mathcal{L})I[t], \quad 0 < \mathcal{L} < 1 \tag{1}$$

$$S[t] = \Theta\left(H[t] - V_{\text{th}}\right) \tag{2}$$

$$V[t] = H[t] - S[t] \times V_{\text{th}} \tag{3}$$

Here, $\mathcal{L}$ denotes the membrane potential decay constant, $I[t]$ is the input current at time step $t$, and the firing threshold is $V_{\text{th}} = 1$. A spike $S[t]$ is emitted when the membrane potential $H[t]$ exceeds $V_{\text{th}}$. $\Theta(x)$ represents the Heaviside step function, which equals 1 when $x \geq 0$ and 0 otherwise. After firing, the membrane potential $V[t]$ is reset by subtracting $S[t] \times V_{\text{th}}$; if no spike occurs, $V[t] = H[t]$.

To better process phase dependent information, we modify the standard LIF neuron as follows:

$$H[t] = \mathcal{L}_t^{\text{learn}}V[t-1] + \beta_t^{\text{learn}}I[t] \tag{4}$$

where $\mathcal{L}_t^{\text{learn}}$ represents a learnable membrane potential decay constant at time step $t$, and $\beta_t^{\text{learn}}$ denotes a learnable weight that controls the contribution of the input current at each time step.

Since the input current can be expressed as $I[t] = W^{(n-1)}S^{(n-1)}[t]$, the modified neuron enables the network to adaptively rescale the spike inputs $S^{(n-1)}[t]$ through the learnable parameters $\beta_t^{\text{learn}}$. Here, $n$ denotes the layer index and $W$ represents the convolutional or fully connected weights. The learnable parameters are shared among all neurons within each channel.

4

## 3. Method

This section first introduces how to optimize rate coding through learnable convolutional layers and surrogate gradients, enabling comparable performance to direct encoding under very few time steps. Subsequently, a learnable temporal encoding mechanism is proposed to enhance the network's ability to model temporal dependencies.

### 3.1. Learnability of Convolutional Layers and Surrogate Gradient

In conventional SNN encoding schemes such as rate coding, TTFS, and phase coding a relatively large number of time steps are typically required to faithfully represent all input information (e.g., every pixel in an image). This not only escalates training costs and energy consumption but also inevitably introduces redundant spike activities lacking discriminative power at the input stage, thereby undermining the efficiency of temporal feature representation.

In contrast, direct encoding achieves competitive performance with extremely few time steps by exploiting the **learnability of preceding convolutional layers**, rather than relying on intrinsic temporal dynamics. These convolutional layers extract discriminative spatial features prior to the spiking transformation.

Motivated by this observation, we prepend convolutional layers to traditional encoding modules, enabling the network to adaptively form task effective spike representations and substantially narrow the performance gap with direct encoding.

However, directly applying binary operations (e.g., rounding or Heaviside step functions) to the convolutional outputs leads to gradient blockage, thereby preventing effective parameter updates in the preceding convolutional layers during backpropagation. To address this issue, a surrogate gradient is employed to smoothly approximate the non differentiable spiking activation, enabling end-to-end optimization [22]. This design significantly enhances feature learning under low time step constraints and provides a solid foundation for subsequent temporal modeling.

### 3.2. Learnable Temporal Encoding

Building upon the analysis of the temporal degeneration in direct encoding, we propose a learnable temporal encoding mechanism that introduces

5

trainable temporal dependencies under static input conditions. This mechanism transforms the convolutional feature maps into spike sequences with controllable temporal variations. (See Figure 2 for an illustration of the mechanism).

Formally, given an input feature map $X$, it is replicated along the temporal axis to form $\{X_t\}_{t=1}^{T}$. At each time step, the spike activation is computed as:

$$s_t = \Theta(x_t - \theta_t), \tag{5}$$

where $\Theta(\cdot)$ is the Heaviside step function.which is approximated by a surrogate gradient during backpropagation.To emulate temporal progression, the firing threshold is adaptively decayed across time steps according to the following.

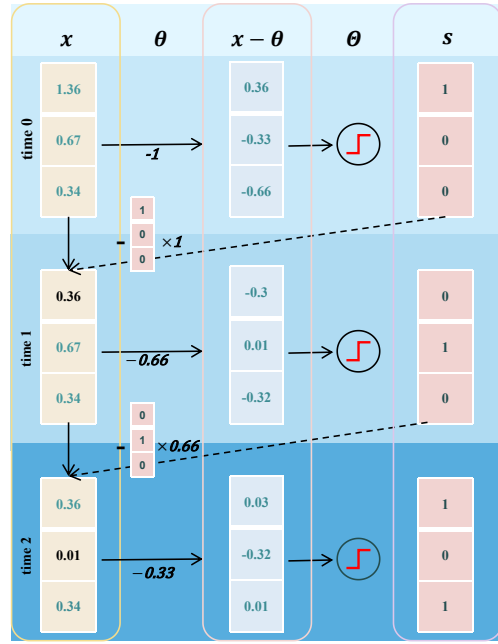$$\theta_{t+1} = \theta_t \times \sigma(a_t), \tag{6}$$



Figure 2: Example of the learnable temporal encoding (phase encoding). Values $\theta_0 = 1$, $\theta_1 = 0.66$, and $\theta_2 = 0.33$ are illustrative; actual $\theta_t$ are learnable.

The adaptively decaying threshold $\theta_t$ (controlled by a learnable, channel-wise shared $a_t$) introduces a time varying activation boundary, enabling the network to learn temporal patterns in spike generation under different encoding rules. After spike generation, different update rules lead to three classical temporal encoding forms:

1)Phase Encoding.

$$x_{t+1} = x_t - s_t \times \theta_t \tag{7}$$

Where each spike subtracts the threshold term from the input, producing progressively phase shifted feature sequences.

*2)TTFS.*

$$x_{t+1} = x_t - s_t \times x_t \tag{8}$$

Which enforces that each neuron fires at most once by setting its input to zero after the first spike.

*3)Rate Encoding.*

When no temporal update is applied, the input features are simply repeated across time, corresponding to conventional rate coding.

This unified framework jointly describes these three spike generation paradigms and enables end-to-end optimization via learnable thresholds and surrogate gradients. It preserves structural simplicity while providing explicit temporal structures, allowing static inputs to be represented as temporally dependent spike patterns.

## 4. Experiments

In this section, we evaluate four encoding strategies on both image classification and object detection tasks. For classification, we use the high performance SNN architecture QKFormer-3-256 [23]. For object detection, we develop a spiking version of YOLOv5n [24], using SEW-ResNet blocks for residual connections to ensure fair comparison across encoding schemes [25]. All experiments are conducted under an extremely low time-step regime ($T = 4$) to evaluate efficiency in ultra short temporal windows. Additional model configurations are provided in the Appendix C.

*4.1. Ablation Study*

Using TTFS as a representative case, we observe that most of its performance gain arises simply from replacing raw image encoding with convolutional features. This improvement is largely attributable to the substantial

Table 1: Ablation study on CIFAR-10 using TTFS. LT = Learnable Threshold $\theta_t$ ,LC = Learnability of Convolution Layer,SG=Surrogate Gradient.

| Method | LT | LC | SG | Acc@1 |
|---|---|---|---|---|
| TTFS | | | | 83.22 |
| TTFS+LT | ✓ | | | 83.44 |
| TTFS+LT+LC | ✓ | ✓ | | 94.32 |
| TTFS+LT+LC+SG | ✓ | ✓ | ✓ | 95.23 |

channel expansion in the first convolutional layer, which compensates for the shallow temporal depth (T=4).

In contrast, surrogate gradients are ineffective when applied directly to raw images, as no learnable parameters precede spike generation. When combined with a learnable convolution, surrogate gradients offer only a modest additional gain ($\approx$1 percentage point), as shown in Table 1.

For reference, rate and phase encoding achieve 86.67% and 87.77% top-1 accuracy under the same four step setting, further confirming that TTFS's performance gain mainly comes from convolutional feature extraction.

## 4.2. Validation Experiments

The results show that all four encoding methods achieve comparable performance on image classification tasks, where temporal modeling plays a limited role under static inputs.In contrast, the performance differences become more evident on the detection task. Compared with image classification, detection poses a more complex prediction setting, which makes the model more sensitive to the availability of temporal variation. As a result, the proposed phase-based encoding shows a clear advantage, while direct and rate encoding remain limited due to the absence of meaningful temporal dynamics. The results across datasets are summarized in Table 2.

## 4.3. Discussion

To illustrate how channel redundancy affects the role of surrogate gradients, we consider a setting where the first convolutional layer is restricted to three output channels under a TTFS like firing scheme. In this case, performance increases from 83.22% to 86.43% when a learnable convolution is

Table 2: Performance comparison of four encoding schemes on CIFAR-10, CIFAR-100, and VOC. For fairness, all methods share the same convolutional front-end, consisting of the learnability of the convolution layer and surrogate gradient, as direct encoding inherently relies on convolutional learnability; learnable thresholds are also included for consistency.

| Encoding | CIFAR-10 Acc@1 | CIFAR-100 Acc@1 | VOC mAP@0.5 |
|---|---|---|---|
| Direct | 95.60 | 77.90 | 0.452 |
| Rate+Conv | 95.62 | 78.08 | 0.444 |
| Phase+Conv | **95.64** | **78.56** | **0.467** |
| TTFS+Conv | 95.26 | 77.93 | 0.460 |

used, but increases to 93.71% when surrogate gradients are used in conjunction with the learnable convolution.

This example demonstrates that the benefit of learnable convolution relies on sufficient channel expansion in the first layer. When the channel capacity is limited, the influence of surrogate gradients becomes significantly more pronounced, revealing their essential contribution beyond the encoding scheme.

Moreover, different encoding schemes exhibit distinct behaviors across classification and detection tasks. We argue that this discrepancy likely does not stem from the encoding schemes themselves, but rather from whether the network actually requires additional temporal depth for effective representation. Using direct encoding as an example, reducing the time steps in QKFormer from four to two yields only a marginal drop on CIFAR-10 (95.6%→94.89%), yet causes a noticeably larger decline on VOC detection (mAP 0.452→0.365).

This suggests that modern SNN backbones can solve coarse grained classification with minimal reliance on temporal structure, whereas fine-grained localization tasks remain more sensitive to reductions in temporal integration.

## 5. Conclusion and Future Work

This work re-examines direct encoding in static-image SNNs, demonstrating its functional equivalence to rate coding and showing that the reported performance gaps mainly arise from the learnability of convolutional layers and the use of surrogate gradients. To address the lack of temporal dynamics, we introduce a minimal learnable temporal encoding that adds adaptive phase shifts to generate meaningful temporal variation.

Experiments verify that this temporal encoding restores useful time structure and improves performance on tasks that rely on temporal processing. While this study provides a unified perspective on three mainstream encoding paradigms, their practical advantages may vary across scenarios for instance, TTFS based methods may be preferable for ultra low power applications, whereas rate based approaches may offer greater robustness in noisy environments.

Overall, our findings highlight the necessity of explicit temporal structure when static inputs contain no inherent time variation. Future work will

explore broader integration of temporal design principles with other components of SNNs to enhance learning under static-input conditions.

## Appendix A. Experimental Results with LIF Neurons

We compare standard LIF neurons with learnable LIF neurons to investigate the impact of neuron model variations on encoding efficacy. The baseline ANN (YOLOv5n) achieves an mAP@.5 of 0.668 on VOC. As shown in Table A.1, standard LIF neurons show minimal performance differences across encoding strategies.

Table A.1: Experimental results with LIF neurons.

| Encoding | CIFAR-10 Acc@1 | CIFAR-100 Acc@1 | VOC mAP@.5 |
|---|---|---|---|
| Direct Encoding | 95.37 | 78.35 | 0.451 |
| Rate Encoding | 95.46 | 78.08 | 0.442 |
| Phase Encoding | 95.56 | 77.71 | 0.451 |
| TTFS | 95.23 | 77.82 | 0.439 |

## Appendix B. Concrete Example of Direct Encoding to Rate Coding Equivalence

To intuitively illustrate the rate-coding equivalence of direct encoding under static inputs, we consider the LIF neuron model with $T = 4$ time steps, membrane decay constant $\tau = 0.5$, and soft reset (membrane potential subtracts 1 upon firing). When input current $X$ is constant across time steps, it uniquely maps to a firing frequency, forming rate coding. Only seven distinct firing patterns are possible, as listed in Table B.1.

Furthermore, we conduct a **temporal shuffling experiment** to assess reliance on temporal ordering by randomly permuting spike trains across time steps while preserving the total number of spikes per neuron. Minimal performance degradation suggests that the network relies on firing rates rather than precise temporal patterns.

On CIFAR-10, accuracy drops slightly from 95.37% to 95.01% after shuffling, confirming that information is encoded via spike rates.

Table B.1: Firing patterns for constant input current $X$ over four time steps.

| Firing Pattern | Boundary Range |
| --- | --- |
| 0000 | $X < 1.066$ |
| 0001 | $1.066 \leq X < 1.142$ |
| 0010 | $1.142 \leq X < 1.333$ |
| 0101 | $1.333 \leq X < 1.714$ |
| 0110 | $1.714 \leq X < 1.866$ |
| 0111 | $1.866 \leq X < 2$ |
| 1111 | $2 \leq X$ |

## Appendix C. Experimental Settings

To ensure reproducibility, we adopt the following unified experimental settings: time steps $T = 4$; LIF neuron membrane decay constant $\tau = 0.5$; surrogate gradient: Sigmoid; pretrained weights: none; initial threshold values: 1. For CIFAR tasks, we use QKFormer-3-256 with parameters consistent with the original implementation. For object detection, we employ Spiking-YOLOv5n with 300 training epochs and batch size 24. Random seeds are set to 35, 1000, and 0 (three runs). To balance computational cost and repeatability, detection experiments use VOC images resized to $320 \times 320$. This setting does not alter the relative performance trends across encoding schemes, affecting only absolute accuracy levels. All other hyperparameters follow the public source code. Experiments are conducted using PyTorch on a single-GPU environment.

## References

[1] Y. Cao, Y. Chen, D. Khosla, Spiking deep convolutional neural networks for energy-efficient object recognition, Int. J. Comput. Vision 113 (1) (2015) 54–66. doi:10.1007/s11263-014-0788-3.
URL https://doi.org/10.1007/s11263-014-0788-3

[2] W. Maass, Networks of spiking neurons: the third generation of neural network models, Neural Networks 10 (9) (1997) 1659–1671.

[3] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain, Y. Liao, C.-K. Lin, A. Lines,

R. Liu, D. Mathaikutty, S. McCoy, A. Paul, J. Tse, G. Venkataramanan, Y.-H. Weng, A. Wild, Y. Yang, H. Wang, Loihi: A neuromorphic many-core processor with on-chip learning, IEEE Micro 38 (1) (2018) 82–99. doi:10.1109/MM.2018.112130359.

[4] S. K. Esser, P. A. Merolla, J. V. Arthur, A. S. Cassidy, R. Appuswamy, A. Andreopoulos, D. J. Berg, J. L. Mckinstry, T. Melano, D. R. Barch, Convolutional networks for fast, energy-efficient neuromorphic computing, Proc Natl Acad Sci U S A 113 (41) (2016) 11441–11446.

[5] K. Roy, A. Jaiswal, P. Panda, Towards spike-based machine intelligence with neuromorphic computing, Nature 575 (7784) (2019) 607–617.

[6] Y. Guo, X. Huang, Z. Ma, Direct learning-based deep spiking neural networks: a review, Frontiers in Neuroscience (2023).

[7] B. Han, G. Srinivasan, K. Roy, Rmp-snn: Residual membrane potential neuron for enabling deeper high-accuracy and low-latency spiking neural network, IEEE (2020).

[8] Y. Hu, H. Tang, G. Pan, Spiking deep residual network (2018).

[9] X. Chen, Q. Yang, J. Wu, H. Li, K. C. Tan, A hybrid neural coding approach for pattern recognition with spiking neural networks, IEEE transactions on pattern analysis and machine intelligence 46 (5) (2023) 3064–3078.

[10] S. R. Kheradpisheh, M. Ganjtabesh, S. J. Thorpe, T. Masquelier, Stdp-based spiking deep convolutional neural networks for object recognition, Neural Networks the Official Journal of the International Neural Network Society 99 (2017) 56.

[11] T. Wang, Q. Shen, X. Li, Y. Zhang, Z. Wang, C. Yan, Precise spiking neurons for fitting any activation function in ann-to-snn conversion, Applied Intelligence 55 (6) (2025) 463.

[12] Y. Hu, L. Deng, Y. Wu, M. Yao, G. Li, Advancing spiking neural networks toward deep residual learning, IEEE Transactions on Neural Networks and Learning Systems 36 (2) (2024) 2353–2367, publisher Copyright: IEEE. doi:10.1109/TNNLS.2024.3355393.

[13] Y. Wu, L. Deng, G. Li, J. Zhu, L. Shi, Direct training for spiking neural networks: Faster, larger, better, 2019, pp. 1311–1318.

[14] E. Hunsberger, C. Eliasmith, Spiking deep networks with lif neurons, arXiv preprint arXiv:1510.08829 (2015).

[15] E. M. Izhikevich, Which model to use for cortical spiking neurons?, IEEE Transactions on Neural Networks 15 (5) (2004) 1063–1070.

[16] N. Rathi, K. Roy, 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) (2021) 3174–3182.

[17] C. Shoushun, A. Bermak, Arbitrated time-to-first spike cmos image sensor with on-chip histogram equalization, IEEE Transactions on Very Large Scale Integration Systems 15 (2007) 346–357.

[18] J. Kim, H. Kim, S. Huh, J. Lee, K. Choi, Deep neural networks with weighted spikes, Neurocomputing 311 (2018) 373–386.

[19] Y. Kim, H. Park, A. Moitra, A. Bhattacharjee, Y. Venkatesha, P. Panda, Rate coding or direct coding: Which one is better for accurate, robust, and energy-efficient spiking neural networks?, in: ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), IEEE, 2022, pp. 71–75.

[20] Neurons, Maass, Single neurons , populations , plasticity (2002).

[21] A. N. Burkitt, A review of the integrate-and-fire neuron model: I. homogeneous synaptic input, Biological cybernetics 95 (1) (2006) 1–19.

[22] E. O. Neftci, H. Mostafa, F. Zenke, Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks, IEEE Signal Processing Magazine 36 (6) (2019) 51–63.

[23] C. Zhou, H. Zhang, Z. Zhou, L. Yu, L. Huang, X. Fan, L. Yuan, Z. Ma, H. Zhou, Y. Tian, Qkformer: Hierarchical spiking transformer using q-k attention (2024).

[24] G. Jocher, ultralytics/yolov5: v3.1 - Bug Fixes and Performance Improvements, `https://github.com/ultralytics/yolov5` (Oct. 2020). doi:10.5281/zenodo.4154370.
URL `https://doi.org/10.5281/zenodo.4154370`

[25] W. Fang, Z. Yu, Y. Chen, T. Huang, T. Masquelier, Y. Tian, Deep residual learning in spiking neural networks, Advances in Neural Information Processing Systems 34 (2021) 21056–21069.