

---

# Stepwise Weighted Spike Coding for Deep Spiking Neural Networks

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Spiking Neural Networks (SNNs) seek to mimic the spiking behavior of biological  
2 neurons and are expected to play a key role in the advancement of neural computing  
3 and artificial intelligence. The efficiency of SNNs is often determined by the  
4 neural coding schemes. Existing coding schemes either cause huge delays and  
5 energy consumption or necessitate intricate neuron models and training techniques.  
6 To address these issues, we propose a novel Stepwise Weighted Spike (SWS)  
7 coding scheme to enhance the encoding of information in spikes. This approach  
8 compresses the spikes by weighting the significance of the spike in each step of  
9 neural computation, achieving high performance and low energy consumption. A  
10 Ternary Self-Amplifying (TSA) neuron model with a silent period is proposed for  
11 supporting SWS-based computing, aimed at minimizing the residual error resulting  
12 from stepwise weighting in neural computation. Our experimental results show  
13 that the SWS coding scheme outperforms the existing neural coding schemes in  
14 very deep SNNs, and significantly reduces operations and latency.

## 15 1 Introduction

16 Spiking Neural Networks (SNNs) are known as the third generation of neural network models  
17 inspired by the biological structures and functions in the brain [32]. Unlike traditional Artificial  
18 Neural Networks (ANNs) that use continuous activation functions, SNNs incorporate discrete spiking  
19 events, enabling them to capture temporal dynamics and process information in a manner that closely  
20 mimics the brain's functioning [31]. This event-driven paradigm aligns with the brain's energy-  
21 efficient computation and has the potential for more efficient and lower-power computing systems.  
22 [33].

23 Various coding schemes have been proposed to describe neural activities, including rate coding and  
24 temporal coding [9]. Rate coding counts the number of spikes fired within a broad time window  
25 [23, 3, 18, 6], which effectively mitigates the impact of short-term interference on the signal. It was  
26 widely accepted in the early days and typically outperformed temporal coding [11, 34, 4, 29, 20].  
27 However, the rate coding scheme disregards the information in the temporal domain of the input  
28 spike sequence and requires many pulses to represent the input signal value, making it an inefficient  
29 coding method that negates the low-power benefits of SNN. Due to the functional similarity to the  
30 biological neural network, spiking neural networks can embrace the sparsity found in biology and  
31 are highly compatible with temporal coding [31, 33, 27, 28, 21, 15]. Temporal coding relies on  
32 the specific timing or patterns of input spikes, allowing for greater information capacity in a single  
33 pulse. However, it requires a large number of time steps to provide fine-grained timing, which  
34 increases inference latency. Its sensitivity to variations in spike timing also makes it more vulnerable  
35 to temporal jitter or delays [25, 24]. Additionally, decoding temporal-coded information usually  
36 requires more complex neuron models [30, 36] and training methodologies [17, 26].

37 In the study of the temporal information dynamics of spikes, Kim et al. [16] discovered a phenomenon  
38 of temporal information concentration in SNNs. It is found that after training, information becomes  
39 highly concentrated in the first few timesteps. Based on this observation, we hypothesize that, from  
40 the perspective of the postsynaptic neuron, the first arriving spikes contain more information and  
41 require stronger responses. Consequently, we propose a mechanism whereby the neuron augments  
42 its own membrane potential with a specific coefficient prior to processing the subsequent input.  
43 This enhancement serves to increase the importance of preceding pulses on neurons, which is why  
44 the spikes are designated as Stepwise Weighted Spikes (SWS). Nevertheless, the amplification of  
45 the membrane potential makes it difficult for neurons to reduce its value through traditional "soft  
46 reset" (i.e. subtracted by an amount equal to the firing threshold), which can result in residual errors  
47 after neuron firing. To address this issue, we make the membrane potential reduced by a magnitude  
48 exceeding the threshold after firing. As a result, the membrane potential has both positive and negative  
49 residual values, which will generate both positive and negative spikes. This neuron is designated as a  
50 Ternary Self-Amplifying (TSA) neuron. To further reduce the error caused by the weighting process,  
51 a silent period is incorporated into the TSA neuron, allowing it to receive more input information  
52 before firing. We perform the classification tasks with SWS-based SNN on MNIST, CIFAR10, and  
53 ImageNet. The results show that the SWS coding scheme can achieve better performance with much  
54 fewer coding and computing steps. Even in very deep SNN, SWS coding scheme still performs well  
55 and achieves similar accuracy to the ANN with the same structure. Our major contributions to this  
56 paper can be summarized as follows:

- 57 • We propose the SWS coding scheme, which enables easy implementation of SNNs with  
58 low energy consumption and high accuracy. The stepwise weighting process enhances  
59 the information-carrying capacity of the preceding pulses, greatly reducing the number  
60 of coding spikes. Negative pulses are introduced in SWS coding to ensure an accurate  
61 information transmission.
- 62 • A novel TSA neuron model is proposed. TSA neuron progressively weights the input by  
63 augmenting its residual membrane potential before receiving the subsequent spike. The  
64 introduction of negative residual membrane potential and negative thresholds enhances the  
65 accuracy of the model's output.
- 66 • A silent period is added to TSA neuron to markedly improve accuracy at minimal latency  
67 cost. By adjusting the silent period step and coding step, SWS-based SNNs can exhibit  
68 performance advantages in different aspects, improving the flexibility of applications.

## 69 **2 Related work**

70 SNNs use spike sequences to convey information, making the encoding of real data into pulses a  
71 crucial step. Currently, the mainstream schemes of neural coding are rate coding and temporal coding  
72 [9, 33, 32]. Rate coding represents different activities with the number of spikes emitted within  
73 a specific time window. Due to its simplicity, rate coding is commonly used in deep learning of  
74 SNNs. However, it distributes information uniformly across a large number of spikes, resulting in an  
75 inefficient transmission process that increases network latency and energy consumption. Numerous  
76 researchers have proposed solutions to optimize inference latency in rate coding. Han et al.[11]  
77 proposed a "soft reset" spiking neuron model that retains a residual membrane potential after firing  
78 to better mimic the ReLU functionality. They demonstrated near lossless ANN-SNN conversion by  
79 using 2-8 times fewer inference time steps. Still, a delay of thousands of steps is required in large  
80 datasets or deep networks. In [14], Hu et al. reduced the encode time steps by converting a quantized  
81 low-precision ANN to a rate-coded SNN. They also proposed a layer-wise fine-tuning mechanism  
82 to minimize the inference latency. However, their neuron model and the subsequent fine-tuning  
83 algorithm are relatively complex. Furthermore, in deeper neural networks such as ResNet56, a 1.5%  
84 drop in accuracy can be observed. The above rate encoding solutions are limited because they do not  
85 consider the significance of each spike.

86 In [15], Kim et al. proposed phase coding, which assigns different weights to spikes based on their  
87 time phase. However, the transmission amount of information is bounded by the global phase, which  
88 causes inefficiency in hidden layers, resulting in a latency of up to three thousand steps for a 32-layer  
89 network. Burst coding [21] attempts to overcome this issue by introducing burst spikes, which  
90 utilize Inter-Spike Interval (ISI). Burst spikes are capable of conveying more information quickly and  
91 accurately by inducing Post-Synaptic Potential (PSP) dramatically. Nevertheless, it is still deficient

Table 1: Common symbols and their meanings in this paper.

Symbol	Meaning
$S_i^l(t)$	The spike train fired by the $i^{th}$ neuron in the $l^{th}$ layer
$u_i^l(t)$	The membrane potential of the $i^{th}$ neuron in the $l^{th}$ layer
$z_i^l(t)$	The integrated inputs to the $i^{th}$ neuron in the $l^{th}$ layer
$V_{th}^l$	The firing threshold of the neurons in the $l^{th}$ layer
$\theta^l$	The amplitude of the spikes fired by the neurons in the $l^{th}$ layer

92 in terms of latency and efficiency. Rueckauer and Liu [27] proposed an efficient temporal encoding  
 93 scheme where the analog activation values of the ANN neurons are represented by the inverse Time-  
 94 To-First-Spike (TTFS) in the SNN neurons. Their new spiking network model generates 7-10 times  
 95 fewer pulses by utilizing temporal information carried by a single spike. However, as pointed out  
 96 in [10], TTFS coding scheme incurs expensive memory access and computational overhead, which  
 97 diminishes the benefit of reduced pulse count. Furthermore, TTFS necessitates a large number of time  
 98 steps to differentiate between various time points, which also increases network latency. Han and  
 99 Roy [10] proposed the Temporal-Switch-Coding (TSC) scheme, in which each input image pixel is  
 100 represented by two spikes, and its intensity is proportional to the timing between the two pulses. Their  
 101 results showed a reduction in energy expenditure. However, TSC coding requires a large number of  
 102 time steps to provide distinguishable time intervals, rendering it an ineffective approach to addressing  
 103 the issue of the long latency.

104 Overall, rate coding employs a large number of pulses to encode information, which results in a  
 105 considerable energy overhead and inference delays. On the other hand, temporal coding allows for  
 106 greater information capacity in a single spike, but this does not reduce the computing latency as a  
 107 precise time point or period can be identified only with a sufficient number of time steps. Therefore,  
 108 new neural coding schemes should be developed.

### 109 3 Stepwise weighted spike coding scheme

#### 110 3.1 Stepwise weighting

111 The spike train  $S_i^l(t)$  of the  $i^{th}$  neuron in the  $l^{th}$  layer can be expressed as follows:

$$S_i^l(t) = \sum_{t_i^{l,(f)} \in F_i^l} \theta^l \delta(t - t_i^{l,(f)}) \quad (1)$$

112 where  $\delta(t)$  is the Dirac delta function,  $\theta^l$  is the spike amplitude of the  $l^{th}$  layer, which is usually set  
 113 to the same value as the firing threshold.  $f$  is the index of the spike in the sequence, and  $F_i^l$  denotes a  
 114 set of spike times which satisfies the firing condition:

$$t_i^{l,(f)} : u_i^l(t_i^{l,(f)}) \geq V_{th}^l \quad (2)$$

115 where  $u_i^l(t)$  denotes the membrane potential and  $V_{th}^l$  denotes the firing threshold of the neurons in  
 116 the  $l^{th}$  layer.

117 Our basic idea is to amplify the membrane potential before the receipt of the subsequent input, which  
 118 amplifies and prolongs the impact of the preceding input spikes on membrane potential, emulating the  
 119 phenomenon of information concentration identified in [16]. For clarity, the meanings of important  
 120 symbols are provided in table 1. The action of a neuron in SWS-SNN can be described as follows:

$$u_j^l(t) = \beta u_j^l(t-1) + z_j^l(t) - S_j^l(t) \quad (3)$$

121 where  $\beta$  is the amplification factor which should be greater than one,  $z_j^l(t)$  denotes the PSP (i.e.  
 122 integrated inputs):

$$z_j^l(t) = \sum_i \omega_{ij}^l S_i^{l-1}(t) + b_j^l \quad (4)$$

123 where  $\omega_{ij}$  is the synaptic weight and  $b_j^l$  is the bias. Begin with the initial value  $u_j^l(0) = 0$  and  
 124 iteratively apply eq. (3) for each subsequent value until  $u_j^l(n)$  and substitute eq. (1) and eq. (4) into it,

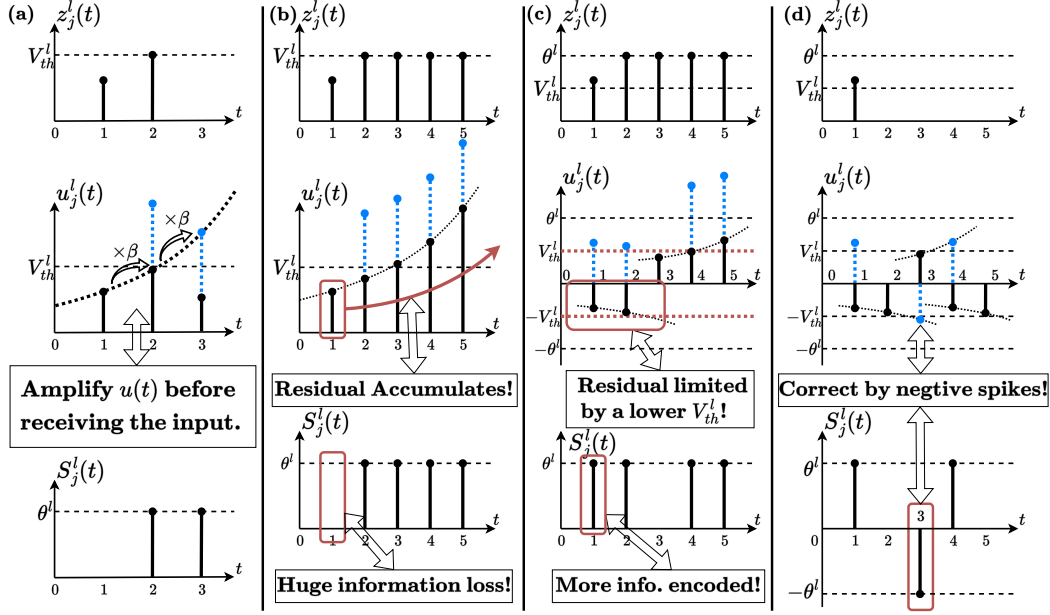


Figure 1: (a) Illustration of the stepwise weighting process. The meanings of the symbol  $z_j^l(t)$ ,  $u_j^l(t)$  and  $S_j^l(t)$  can be found in table 1. The blue dotted line represents the membrane potential prior to the spike firing, and the black exponential function-like dotted line is employed to illustrate the trend of membrane potential amplification. (b) A  $V_{th}^l$  equal to  $\theta^l$  results in residual errors, leaving a lot of information unencoded. (c)  $V_{th}^l$  is set to  $\frac{1}{2}\theta^l$ , which increases the possibility to fire spikes early to better limit the residual. (d) Use negative spikes to correct the excessively emitted information.

125 eq. (3) can be written as:

$$u_j^l(n) = \beta^n u_j^l(0) + \sum_{\tau=1}^n \beta^{n-\tau} z_j^l(\tau) = \sum_{t_i^{l-1,(f)}} \sum_i \sum_{\tau=1}^n \beta^{n-\tau} \omega_{ij}^l \theta^{l-1} \delta(\tau - t_i^{l-1,(f)}) + \beta^{n-\tau} b_j^l \quad (5)$$

126 Note that  $S_j^l(t)$  is set to zero for simplicity. From eq. (5), it can be seen that the stepwise augment of  
 127 the membrane potential results in the spike input at time  $t_i^{l-1,(f)}$  encoding the value  $\theta^{l-1} \beta^{n-t_i^{l-1,(f)}}$ .  
 128 This process is thus referred to as stepwise weighting, and  $\beta^{n-t_i^{l-1,(f)}}$  serves as the weight. The  
 129 earlier the input pulse, the greater its ability to carry information. This solves the problem of excessive  
 130 encoding steps in previous schemes, allowing faster information transmission.

### 131 3.2 Residual error

132 Stepwise weighting effectively assigns more weight to earlier arriving pulses, but it also makes spike  
 133 generation more tricky. To ensure that input information is efficiently encoded and transmitted to  
 134 the next layer, the residual membrane potential should be minimized after neural computation is  
 135 completed. The stepwise weighting, however, amplifies the residual potential from the previous  
 136 time step. If  $z_j^l(t)$  remains high in subsequent steps, reducing the membrane potential becomes  
 137 challenging, as shown in fig. 1(b). This vicious cycle ultimately leads to a persistently high membrane  
 138 potential, indicating that a substantial amount of information remains unencoded.

139 We refer to this phenomenon as residual error. One contributing factor is that the threshold is set  
 140 too high, resulting in a pulse being emitted only when the membrane potential exceeds the value  $\theta^l$ .  
 141 While this prevents excessive information transmission, it results in missed opportunities to bring  
 142 down  $u_j^l(t)$  by firing a spike.

143 To address this issue, we propose setting the firing threshold  $V_{th}^l$  to  $\frac{1}{2}\theta^l$ . This adjustment facilitates  
 144 pulse generation and reduces the residual membrane potential. After the neuron firing, the membrane  
 145 potential is subtracted by  $\theta^l$ , which leads to the emergence of a negative residual that will be stepwise

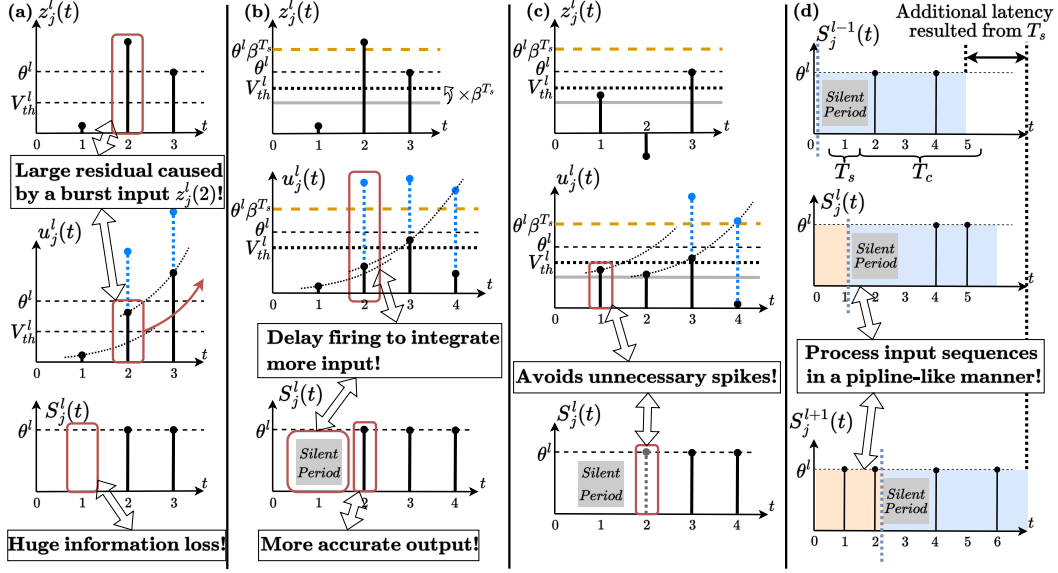


Figure 2: (a) Uncertainty in the input distribution leads to residual errors. (b) The silent period allows more information to be known when firing pulses.  $T_s$  is set to 1 here.  $V_{th}^l$  is amplified by  $\beta^{T_s}$ , and the original threshold is represented by a gray solid line. The orange dashed line represents the amount of membrane potential reduction after firing. (c) The silent period also avoids some unnecessary spikes and increases sparsity. Without the silent period, since  $u_j^l(1)$  exceeds the original threshold, a pulse will be generated at  $t = 1$ , which will later be corrected by another negative spike. (d) The impact of the silent period on network latency. The output spike sequences corresponding to different inputs are drawn in blocks of different colors. The pulses drawn in the spike sequence are for illustrative purposes only.

146 weighted over time. The coefficient  $1/2$  is selected as it is capable of controlling both positive and  
 147 negative residuals within a narrow and balanced range. A negative threshold  $-V_{th}^l$  is introduced into  
 148 the neuron model, which initiates a negative spike when the membrane potential falls below this  
 149 threshold. This mechanism allows the excessively emitted information to be corrected by the negative  
 150 spike, as shown in fig. 1(d). Given the above characteristics, we designate this neuron model as a  
 151 TSA neuron.

### 152 3.3 Silent period

153 Another contributing factor to residual error is the imbalanced distribution of  $z_j^l(t)$ . A burst input of  
 154  $z_j^l(t)$  at time point  $\tau$  results in a sharp rise in membrane potential, making it difficult for subsequent  
 155 spikes to reduce it, as shown in fig. 2(a).

156 This can be addressed by incorporating a silent period  $T_s$  into the TSA neuron model. The neurons  
 157 only integrates input and performs stepwise weighting, but are not allowed to fire in the first  $T_s$   
 158 steps. This enables the acquisition of more known information before spike generation, resulting  
 159 in increased accuracy, as illustrated in fig. 2(b). Since the preceding input information has been  
 160 amplified by  $\beta^{T_s}$  after the silent period,  $V_{th}^l$  also needs to be adjusted accordingly, which is set to  
 161  $\frac{\beta^{T_s}}{2}\theta^l$ . Similarly, after firing, the membrane potential should be subtracted by  $\theta^l\beta^{T_s}$ . Note that the  
 162 fired spike amplitude remains unchanged, that is,  $\theta^l$ .

163 The impact of the silent period on network latency is shown in fig. 2(d). The output results for  
 164 different input sequences are distinguished by blocks of different colors. It can be observed that  
 165 as network depth increases, the silent period accumulates, leading to a higher output latency. The  
 166 inference latency of SWS-SNN can be calculated as follows:

$$T_{inf} = T_c + T_s \cdot L_{TSA} \quad (6)$$

167 where  $T_{inf}$  is the inference delay,  $T_c$  is the coding time steps,  $T_s$  is the length of the silent period and  
 168  $L_{TSA}$  is the number of TSA neuron layers. The neuron model in other coding schemes yields a zero  
 169  $T_s$ , leading to an output delay equal to the coding time step, which is consistent with the definition in  
 170 the previous scheme. From fig. 2(d), it can be seen that different input sequences are processed in a  
 171 pipeline-like manner, and the value of  $T_c + T_s$  determines the throughput rate of SWS-SNN.

### 172 3.4 Input encoding

173 According to eq. (5), the value that can be losslessly encoded under the SWS coding scheme can be  
 174 expressed as follows:

$$A_j = \sum_{\tau=1}^{T_c} a_j^\tau \cdot \theta^0 \beta^{T_c - \tau} \quad (7)$$

175 where  $A_j$  denotes the encoded value.  $a_j^\tau \in \{-1, 0, 1\}$  indicates the type of the output spike at time  $\tau$ :  
 176 1 for a positive pulse,  $-1$  for a negative pulse and 0 for no pulse.  $T_c$  denoted the time steps used for  
 177 encoding. The weight  $\beta^{T_c - \tau}$  results from the stepwise weighting process described in section 3.1.  $\theta^0$   
 178 denotes the spike amplitude of the input encoding layer, which can be assigned an appropriate value  
 179 based on the range to be encoded.

180 According to eq. (7), given a fixed  $T_c$  and  $\theta^0$ , the distribution of  $A_j$  is determined by  $\beta$ . Setting  $\beta$  to  
 181 2 is reasonable, as it ensures  $A_j$  is evenly distributed within the codable range. Compared to rate  
 182 coding, which necessitates  $2^{T_c}$  coding steps to encode the same range with same precision, SWS  
 183 coding significantly enhances coding efficiency. Note that with the introduction of negative pulses,  
 184 setting  $\beta$  to 3 can also achieve a uniform distribution of  $A_j$  and offers even more values for accurate  
 185 encoding compared to  $\beta = 2$ .<sup>1</sup> When  $\beta$  is less than 2, the distribution of  $A_j$  becomes denser at  
 186 smaller values, which may be suitable for encoding data that follows a similar distribution.

187 For static image classification tasks, the pixel value  $p_j$  can be encoded by applying a constant input  
 188  $z_j^0(t)$  to the TSA neuron. Considering the stepwise weighting process, we can write:

$$p_j = \sum_{\tau=1}^{T_c} |z_j^0| \beta^{T_c - \tau} \quad (8)$$

189 where  $|z_j^0|$  denotes the amplitude of the constant input  $z_j^0(t)$ . Solve for  $|z_j^0|$  and we have:

$$z_j^0(t) = \sum_{\sigma=1}^{T_c} \frac{p_j}{\sum_{\tau=1}^{T_c} \beta^{T_c - \tau}} \cdot \delta(t - \sigma) \quad (9)$$

190 Given that  $z_j^0(t)$  is a constant at each step,  $T_s$  can be set to 0 for the encoding layer. However, the  
 191 neuron must await  $T_s$  time steps after the completion of an encoding. This allows neurons in the  
 192 subsequent layer to complete the previous neural computing before receiving the next encoded input.

## 193 4 Experiments

194 In this section, we convert quantized ANNs to SWS-based SNNs<sup>2</sup> and conduct experiments on  
 195 MNIST, CIFAR10, and ImageNet. Firstly, an overview of SWS-SNN’s performance across various  
 196 datasets is provided. Subsequently, the network’s inference latency and energy consumption is  
 197 compared with other spike coding schemes. Finally, an ablation study is conducted to investigate the  
 198 impact of lowered thresholds and silent periods on reducing residuals and enhancing accuracy.

199 ANNs used for conversion are all quantized to 8 bits.  $\beta$  is set to 2 in the experiments to ensure that  
 200 codable values are evenly distributed. Compared to  $\beta = 3$ , a smaller amplification factor reduces the  
 201 impact of residual errors, resulting in more accurate output.

<sup>1</sup>Setting  $\beta$  to 2 introduces some coding redundancy. E.g.,  $a_j^1 = 1, a_j^2 = -1$  and  $a_j^1 = 0, a_j^2 = 1$  encodes the same amount of information.

<sup>2</sup>Details of the conversion process can be found in appendix A.1 and appendix A.2

Table 2: Performance on CIFAR10 and ImageNet.

	Category	Methods	Architecture	Time Step	$T_s$	SNN Acc	$\Delta\text{Acc}^\dagger$
CIFAR10	Directly Learning	STBP-tdBN[35]	ResNet-19	6	-	93.16%	-
		TET[5]	ResNet-19	6	-	94.50%	-
	ANN-SNN	TTRBR[20]	ResNet-18	64	-	95.04%	-0.13%
		DSR[19]	PreAct-ResNet-18	20	-	95.24%	-
		Calibration[18]	VGG-16	256	-	95.79%	+0.05%
		OPI[1]	VGG-16	256	-	94.49%	-0.08%
		Opt Conversion[4]	ResNet-20	128	-	93.56%	+1.25%
	ANN-SNN	<b>SWS (ours)</b>	ResNet-18	8	1	95.67%	+0.22%
			VGG-16	8	2	95.86%	-0.04%
	ImageNet	Directly Learning	TET[5]	SEW-ResNet-34	4	-	68.00%
STBP-tdBN[35]			SEW-ResNet-34	4	-	67.04%	-
SEW Resnet[8]			SEW-ResNet-152	4	-	69.26%	-
ANN-SNN		Hybrid training[26]	ResNet-34	250	-	61.48%	-8.72%
		Spiking ResNet[13]	ResNet-50	350	-	72.75%	-2.70%
		QCFS[2]	VGG-16	64	-	72.85%	-1.44%
		Fast-SNN[14]	VGG-16	7	-	72.95%	-0.41%
		COS[12]	ResNet-34	8	-	74.17%	-0.05%
		RMP-SNN[11]	ResNet-34	4096	-	69.89%	-0.75%
		TTRBR[20]	ResNet-50	512	-	75.04%	-0.98%
ANN-SNN		<b>SWS (ours)</b>	VGG-16	8	2	75.27%	-0.11%
			ResNet-34	8	2	76.10%	-0.08%
			Inception-v3	8	2	76.70%	-0.70%
			ResNet-50	8	2	80.34%	-0.35%
			ResNeXt101_32x8d	8	1	81.32%	-1.17%
			ResNeXt101_32x8d	8	2	82.06%	-0.42%

 $\dagger \Delta\text{Acc} = \text{Acc}_{\text{SNN}} - \text{Acc}_{\text{ANN}}$ 

## 202 4.1 Overall performance

203 For simple classification tasks such as CIFAR10, our proposed SWS coding scheme has a faster  
204 inference speed than other ANN-SNN models while achieving similar classification accuracy, or has  
205 higher classification accuracy than direct learning at similar inference speeds. For example, ResNet18  
206 with SWS improves throughput seven times over [20] while simultaneously improving accuracy.  
207 Although the network in [5] has a slightly higher throughput, its accuracy is 1.17% lower than our  
208 scheme. To fully test the potential of our proposed coding scheme, we conducted experiments on  
209 ImageNet using networks with various structures. The experimental results demonstrate that SWS  
210 coding has distinct advantages on extremely deep SNNs. Our SWS-based ResNet50 and ResNeXt101  
211 achieved over 80% accuracy on ImageNet with only eight coding steps. The model in [12] achieves  
212 an almost lossless conversion with eight time steps. However, their method has to adjust the resting  
213 potential of neurons layer by layer, and the calibration effect for deeper networks is unclear. In [14],  
214 the original ANN needs to be quantized to 3 bits, resulting in a larger conversion loss. Directly trained  
215 SNNs typically achieve higher throughput, but their accuracy still requires improvement. In addition,  
216 the SWS coding scheme is easy to implement. No further fine-tuning is required after the conversion.

## 217 4.2 Accuracy vs. latency

218 The comparison of latency results between SWS-SNN and other ANN-converted SNNs[1, 11, 10,  
219 4, 18, 2, 7] is illustrated in fig. 3. The latency of the network is calculated with eq. (6). In the  
220 counterpart models, the variation of delay is mainly caused by the changes in  $T_c$ . In contrast,  
221  $T_s$  determines latency in deep SWS-SNNs. Therefore, SWS-SNN has an upper limit on latency:  
222  $T_{inf}^{max} = T_c(1 + L_{TSA})$ , which causes our curve to terminate earlier in fig. 3.

223 To ensure a fair comparison, we represent the ANN accuracy of each counterpart with dotted lines of  
224 the same color. The experimental results indicate that SWS-SNN can achieve optimal performance  
225 with minimal latency. Specifically, SWS-based VGG-16 can converge to the ANN performance

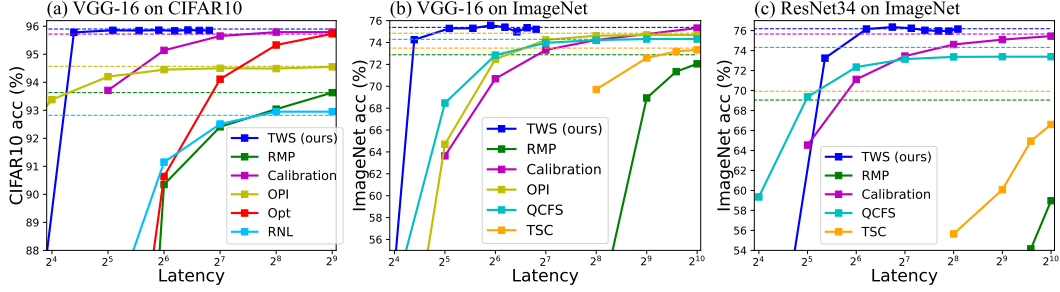


Figure 3: Latency versus accuracy. The ANN accuracy of each compared SNN is marked by dotted lines of the same colour. (a) VGG-16 on CIFAR10. (b) VGG-16 on ImageNet. (c) ResNet34 on ImageNet.

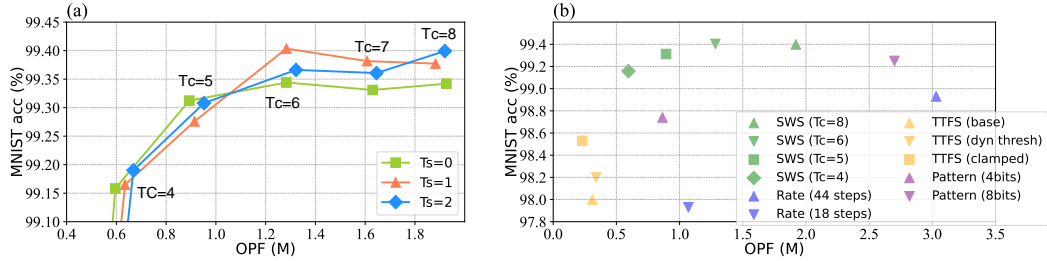


Figure 4: (a) Accuracy versus OPF with different combinations of  $T_c$  and  $T_s$ . (b) Comparison of accuracy and energy consumption of SWS-SNN with other SNNs.

226 in the shortest time on CIFAR10 and reduce the inference latency on ImageNet by more than one  
 227 order. Even though the silent period accumulates when the network gets deeper, the results in fig. 3(c)  
 228 demonstrate that our scheme still achieves the fastest inference speed with the highest accuracy in a  
 229 34-layer network. Note that  $T_s$  is set to the same value for each TSA layer for simplicity, resulting in  
 230 discontinuous  $T_{inf}$  values. This causes a sharp drop in accuracy at smaller delays.

### 231 4.3 Operation counting

232 To compare the energy consumption of SWS-SNN with SNNs under other encoding schemes, we  
 233 adopt the method as in [29, 27, 28] to count operations:

$$OPF = (T_c + T_s)N_{TSA} + \sum_{l=1}^{L_{TSA}} \sum_{\tau=T_s l+1}^{T_s l+T_c} f_{out}^l n^l(\tau) \quad (10)$$

234 where  $OPF$  (Operations Per Frame) denotes the number of operations for the classification of one  
 235 frame,  $T_c$  and  $T_s$  denotes the coding steps and the length of the silent period, respectively.  $L_{TSA}$   
 236 denotes the number of TSA layers,  $f_{out}^l$  denotes the fan-out of neurons in layer  $l$ ,  $n^l(t)$  denotes the  
 237 number of spikes fired in layer  $l$  at time  $\tau$  and  $N_{TSA}$  denotes the number of TSA neurons. The  
 238 first term on the right-hand side of the equation arises from the TSA's requirement to amplify the  
 239 membrane potential. Note that due to the accumulation of  $T_s$  over the network depth, the time period  
 240 for counting  $n^l(t)$  varies with  $l$ .

241 Experiments were conducted on MNIST using LeNet-5. We varied the silent periods and adjusted  
 242 the coding steps to study their effects on OPF. The results are presented in fig. 4(a). As indicated in  
 243 eq. (10), reducing  $T_c$  lowers energy overhead. This presents a trade-off between energy consumption  
 244 and inference accuracy, as fewer coding steps also reduce the number of values that can be accurately  
 245 encoded. A larger  $T_s$  requires TSA neurons to perform more operations to amplify the membrane  
 246 potential. On the other hand, it reduces the number of unnecessary pulse emissions. Overall, silent  
 247 period has a negligible impact on OPF.



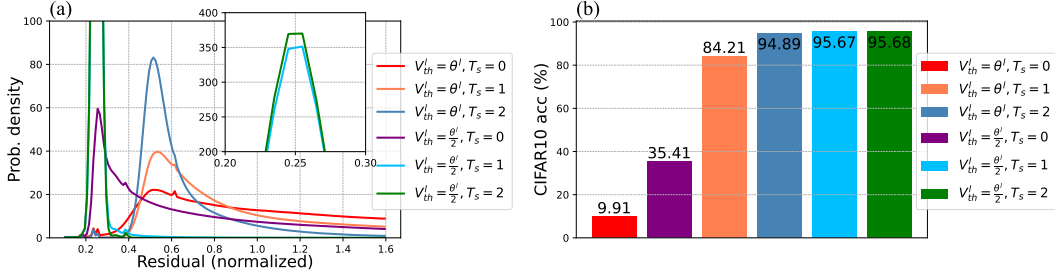


Figure 5: (a) The probability density of the residuals with/without a lowered  $V_{th}^l$  and a silent period. (b) Inference accuracy of SWS-ResNet18 on CIFAR10 with/without a lowered  $V_{th}^l$  and a silent period.

248 In fig. 4(b), the energy consumption of SWS-based SNN is compared with that of other SNNs. The  
 249 experimental results demonstrate that our coding scheme can achieve a favorable balance between  
 250 accuracy and energy consumption. The SWS coding scheme is superior to rate coding and temporal  
 251 pattern coding in that it requires fewer operations and achieves higher accuracy. In TTFS encoding,  
 252 each neuron fires at most one spike at a time, theoretically demanding the least OPF. With  $T_c = 4$ ,  
 253 SWS-SNN can achieve significantly higher accuracy with minimal increase in OPF. Note that if  
 254 the ANN is quantized to a lower number of bits (e.g., 4 bits), the error caused by the reduced  $T_c$   
 255 can actually be compensated by the quantization algorithm, which can potentially result in a higher  
 256 performance.

#### 257 4.4 Ablation study

258 In section 3.2 and section 3.3, we proposed reducing the firing threshold and introducing a silent  
 259 period to mitigate residual error. To assess the impact of these two adjustments, we conducted  
 260 experiments on CIFAR10 using ResNet18. After the neural computation, the residuals (absolute  
 261 values) of the TSA neurons were analyzed. We first scaled the residuals by  $1/\beta^{T_s}$  to counteract the  
 262 effect of membrane potential amplification caused by the silent period, and then normalized them in  
 263 units of  $\theta^l$ . The probability density of the residuals is shown in fig. 5(a).

264 The results demonstrate that lowering  $V_{th}^l$  shifts the residual distribution from around  $0.5\theta^l$  to  
 265 approximately  $0.25\theta^l$ , corresponding to the quantization errors (i.e. rounding errors) under their  
 266 respective thresholds. The addition of silent periods further concentrates the distribution and reduces  
 267 large deviations. As can be seen from the green curve in fig. 5(a), setting  $T_s$  to 2 and  $V_{th}^l$  to  $\theta^l/2$   
 268 makes the residuals almost all distributed around the quantization error. Compared to the red curve  
 269 (without a lowered  $V_{th}^l$  or a silent period), the residuals are greatly reduced, which fully proves the  
 270 effectiveness of lowering the threshold and adding a silent period. The inference results on CIFAR10  
 271 is shown in fig. 5(b). When setting  $V_{th}^l$  to  $\theta^l$  and  $T_s$  to zero, the network's output is almost random.  
 272 Lowering the threshold and adding a silent period improve the accuracy to 35.41% and 84.21%,  
 273 respectively. Ultimately, the combination of both adjustments enabled SWS-ResNet18 to achieve an  
 274 accuracy of 95.68% on CIFAR10.

## 275 5 Conclusion

276 In this work, we have proposed a novel SWS spike coding scheme. The stepwise weighting process  
 277 enhances the information-carrying capacity of the preceding pulses, greatly reducing the number of  
 278 time steps for encoding. Combined with a silent period, our proposed TSA neuron model solves the  
 279 problem of residual errors and achieves fast and accurate information transmission. Our experimental  
 280 results have demonstrated that SWS coding is highly effective in extremely deep SNNs and achieves  
 281 state-of-the-art accuracy. The SWS coding scheme is also highly flexible and can adapt to various  
 282 needs.

## 283 References

- 284 [1] Bu, T., Ding, J., Yu, Z., Huang, T.: Optimized potential initialization for low-latency spiking  
285 neural networks (2022)
- 286 [2] Bu, T., Fang, W., Ding, J., Dai, P., Yu, Z., Huang, T.: Optimal ann-snn conversion for high-  
287 accuracy and ultra-low-latency spiking neural networks (2023)
- 288 [3] Cao, Y., Chen, Y., Khosla, D.: Spiking deep convolutional neural networks for  
289 energy-efficient object recognition. *International Journal of Computer Vision* **113**(1), 54–  
290 66 (May 2015). <https://doi.org/10.1007/s11263-014-0788-3>, <https://doi.org/10.1007/s11263-014-0788-3>
- 292 [4] Deng, S., Gu, S.: Optimal conversion of conventional artificial neural networks to spiking neural  
293 networks (2021)
- 294 [5] Deng, S., Li, Y., Zhang, S., Gu, S.: Temporal efficient training of spiking neural network via  
295 gradient re-weighting (2022)
- 296 [6] Diehl, P.U., Neil, D., Binas, J., Cook, M., Liu, S.C., Pfeiffer, M.: Fast-classifying,  
297 high-accuracy spiking deep networks through weight and threshold balancing. In:  
298 2015 International Joint Conference on Neural Networks (IJCNN). pp. 1–8 (2015).  
299 <https://doi.org/10.1109/IJCNN.2015.7280696>
- 300 [7] Ding, J., Yu, Z., Tian, Y., Huang, T.: Optimal ann-snn conversion for fast and accurate inference  
301 in deep spiking neural networks (2021)
- 302 [8] Fang, W., Yu, Z., Chen, Y., Huang, T., Masquelier, T., Tian, Y.: Deep residual learning in  
303 spiking neural networks (2022)
- 304 [9] Guo, W., Fouda, M.E., Eltawil, A.M., Salama, K.N.: Neural coding in spiking neural net-  
305 works: A comparative study for robust neuromorphic systems. *Frontiers in Neuroscience*  
306 **15** (2021). <https://doi.org/10.3389/fnins.2021.638474>, [https://www.frontiersin.org/  
307 journals/neuroscience/articles/10.3389/fnins.2021.638474](https://www.frontiersin.org/journals/neuroscience/articles/10.3389/fnins.2021.638474)
- 308 [10] Han, B., Roy, K.: Deep spiking neural network: Energy efficiency through time based coding.  
309 In: Vedaldi, A., Bischof, H., Brox, T., Frahm, J.M. (eds.) *Computer Vision – ECCV 2020*. pp.  
310 388–404. Springer International Publishing, Cham (2020)
- 311 [11] Han, B., Srinivasan, G., Roy, K.: Rmp-snn: Residual membrane potential neuron for enabling  
312 deeper high-accuracy and low-latency spiking neural network (2020)
- 313 [12] Hao, Z., Ding, J., Bu, T., Huang, T., Yu, Z.: Bridging the gap between anns and snns by  
314 calibrating offset spikes (2023)
- 315 [13] Hu, Y., Tang, H., Pan, G.: Spiking deep residual networks. *IEEE Transactions on  
316 Neural Networks and Learning Systems* **34**(8), 5200–5205 (2023).  
317 <https://doi.org/10.1109/TNNLS.2021.3119238>
- 318 [14] Hu, Y., Zheng, Q., Jiang, X., Pan, G.: Fast-snn: Fast spiking neural network by converting  
319 quantized ann. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **45**(12), 14546–  
320 14562 (2023). <https://doi.org/10.1109/TPAMI.2023.3275769>
- 321 [15] Kim, J., Kim, H., Huh, S., Lee, J., Choi, K.: Deep neural networks with weighted spikes. *Neuro-  
322 computing* **311**, 373–386 (2018). <https://doi.org/https://doi.org/10.1016/j.neucom.2018.05.087>,  
323 <https://www.sciencedirect.com/science/article/pii/S0925231218306726>
- 324 [16] Kim, Y., Li, Y., Park, H., Venkatesha, Y., Hambitzer, A., Panda, P.: Exploring temporal  
325 information dynamics in spiking neural networks (2022)
- 326 [17] Lee, C., Sarwar, S.S., Panda, P., Srinivasan, G., Roy, K.: Enabling spike-based back-  
327 propagation for training deep neural network architectures. *Frontiers in Neuroscience*  
328 **14** (Feb 2020). <https://doi.org/10.3389/fnins.2020.00119>, [http://dx.doi.org/10.3389/  
329 fnins.2020.00119](http://dx.doi.org/10.3389/fnins.2020.00119)

- 330 [18] Li, Y., Deng, S., Dong, X., Gong, R., Gu, S.: A free lunch from ann: Towards efficient, accurate  
331 spiking neural networks calibration (2021)
- 332 [19] Meng, Q., Xiao, M., Yan, S., Wang, Y., Lin, Z., Luo, Z.Q.: Training high-performance low-  
333 latency spiking neural networks by differentiation on spike representation (2023)
- 334 [20] Meng, Q., Yan, S., Xiao, M., Wang, Y., Lin, Z., Luo, Z.Q.: Training much deeper spiking  
335 neural networks with a small number of time-steps. *Neural Networks* **153**, 254–268 (2022).  
336 <https://doi.org/https://doi.org/10.1016/j.neunet.2022.06.001>, <https://www.sciencedirect.com/science/article/pii/S0893608022002064>
- 338 [21] Park, S., Kim, S., Choe, H., Yoon, S.: Fast and efficient information transmission with burst  
339 spikes in deep spiking neural networks (2019)
- 340 [22] Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin, Z.,  
341 Gimelshein, N., Antiga, L., Desmaison, A., Köpf, A., Yang, E., DeVito, Z., Raison, M., Tejani,  
342 A., Chilamkurthy, S., Steiner, B., Fang, L., Bai, J., Chintala, S.: Pytorch: An imperative style,  
343 high-performance deep learning library (2019)
- 344 [23] Pérez-Carrasco, J.A., Zhao, B., Serrano, C., Acha, B., Serrano-Gotarredona, T., Chen, S.,  
345 Linares-Barranco, B.: Mapping from frame-driven to frame-free event-driven vision systems  
346 by low-rate rate coding and coincidence processing—application to feedforward convnets.  
347 *IEEE Transactions on Pattern Analysis and Machine Intelligence* **35**(11), 2706–2719 (2013).  
348 <https://doi.org/10.1109/TPAMI.2013.71>
- 349 [24] Querlioz, D., Bichler, O., Dollfus, P., Gamrat, C.: Immunity to device variations in a spiking  
350 neural network with memristive nanodevices. *IEEE Transactions on Nanotechnology* **12**(3),  
351 288–295 (2013). <https://doi.org/10.1109/TNANO.2013.2250995>
- 352 [25] Querlioz, D., Bichler, O., Gamrat, C.: Simulation of a memristor-based spiking neural network  
353 immune to device variations. In: *The 2011 International Joint Conference on Neural Networks*.  
354 pp. 1775–1781 (2011). <https://doi.org/10.1109/IJCNN.2011.6033439>
- 355 [26] Rathi, N., Srinivasan, G., Panda, P., Roy, K.: Enabling deep spiking neural networks with hybrid  
356 conversion and spike timing dependent backpropagation (2020)
- 357 [27] Rueckauer, B., Liu, S.C.: Conversion of analog to spiking neural networks using sparse temporal  
358 coding. In: *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*. pp. 1–5  
359 (2018). <https://doi.org/10.1109/ISCAS.2018.8351295>
- 360 [28] Rueckauer, B., Liu, S.C.: Temporal pattern coding in deep spiking neural networks.  
361 In: *2021 International Joint Conference on Neural Networks (IJCNN)*. pp. 1–8 (2021).  
362 <https://doi.org/10.1109/IJCNN52387.2021.9533837>
- 363 [29] Rueckauer, B., Lungu, I.A., Hu, Y., Pfeiffer, M., Liu, S.C.: Conversion of continuous-valued  
364 deep networks to efficient event-driven networks for image classification. *Frontiers in Neuro-*  
365 *science* **11** (2017). <https://doi.org/10.3389/fnins.2017.00682>
- 366 [30] Stöckl, C., Maass, W.: Optimized spiking neurons can classify images with high accu-  
367 racy through temporal coding with two spikes. *Nature Machine Intelligence* **3**(3), 230–  
368 238 (Mar 2021). <https://doi.org/10.1038/s42256-021-00311-4>, <https://doi.org/10.1038/s42256-021-00311-4>
- 370 [31] Taherkhani, A., Belatreche, A., Li, Y., Cosma, G., Maguire, L.P., McGinnity, T.: A re-  
371 view of learning in biologically plausible spiking neural networks. *Neural Networks* **122**,  
372 253–272 (2020). <https://doi.org/https://doi.org/10.1016/j.neunet.2019.09.036>, <https://www.sciencedirect.com/science/article/pii/S0893608019303181>
- 374 [32] Wang, X., Lin, X., Dang, X.: Supervised learning in spiking neural networks:  
375 A review of algorithms and evaluations. *Neural Networks* **125**, 258–280 (2020).  
376 <https://doi.org/https://doi.org/10.1016/j.neunet.2020.02.011>, <https://www.sciencedirect.com/science/article/pii/S0893608020300563>
- 377

- 378 [33] Yamazaki, K., Vo-Ho, V.K., Bulsara, D., Le, N.: Spiking neural networks and their applications:  
379 A review. *Brain Sci* **12**(7) (Jun 2022)
- 380 [34] Yan, Z., Zhou, J., Wong, W.: Near lossless transfer learning for spiking neural networks. In:  
381 AAAI Conference on Artificial Intelligence (2021), [https://api.semanticscholar.org/](https://api.semanticscholar.org/CorpusID:235349069)  
382 CorpusID:235349069
- 383 [35] Zheng, H., Wu, Y., Deng, L., Hu, Y., Li, G.: Going deeper with directly-trained larger spiking  
384 neural networks (2020)
- 385 [36] Zhou, S., LI, X., Chen, Y., Chandrasekaran, S.T., Sanyal, A.: Temporal-coded deep spiking  
386 neural network with easy training and robust performance (2021)

## 387 A Appendix

### 388 A.1 Convert quantized ANNs to SWS-SNNs

389 A pretrained ANN was first obtained from torchvision, which is part of the PyTorch[22] project, and  
390 then quantized into  $n$  bits following the Quantization-Aware Training (QAT) Workflow provided by  
391 PyTorch (8 bits in the actual experiment, with  $n$  bits used here for generality). The quantized ANN  
392 can be characterized by the parameters listed in table 3, and the basic idea of the conversion process is  
393 illustrated in fig. 6(a). The activations of the quantized ANN can be mapped to an integer  $Q$  between  
394  $[0, 2^n - 1]$  using a scaling factor  $C$  and a zero point  $Z$ . With the same weight and bias between  $Q_i^l$   
395 and  $Q_o^l$ , the TSA layer can generate  $S^l$ , which encodes  $Q_o^l$ , provided that  $S^{l-1}$  encodes  $Q_i^l$  and no  
396 residual error occurs. In the actual SNN, the pulse amplitude  $\theta^l$  is normalized to 1. Therefore, the  
397 bias need to be further scaled to derive the final weight  $W^l$  and bias  $b^l$  for the SWS-SNN.

Table 3: The notations and meanings of parameters in the quantized network.

Notation	Meaning
$\hat{X}_i^l$	The quantized input of the $l^{th}$ layer
$\hat{X}_o^l$	The quantized output of the $l^{th}$ layer
$C_i^l$	The scaling factor of the quantized input of the $l^{th}$ layer
$Z_i^l$	The zero point of the quantized input of the $l^{th}$ layer
$C_o^l$	The scaling factor of the quantized output of the $l^{th}$ layer
$Z_o^l$	The zero point of the quantized output of the $l^{th}$ layer
$\hat{W}^l$	The quantized weight of layer $l$
$C_w^l$	The scaling factor of the quantized weight of layer $l$
$Z_w^l$	The zero point of the quantized weight of layer $l$
$\hat{b}^l$	The bias of layer $l$

398 The derivation is as follows. After QAT, we have:

$$\hat{W}^l \hat{X}_i^l + \hat{b}^l = \hat{X}_o^l, \quad (11)$$

$$Q_i^l = \frac{\hat{X}_i^l}{C_i^l} + Z_i^l, \quad (12)$$

$$Q_o^l = \frac{\hat{X}_o^l}{C_o^l} + Z_o^l, \quad (13)$$

399 where  $Q_i^l, Q_o^l$  represent the integers to which the quantized input and output are mapped, respectively.  
400 Substitute eq. (12) and eq. (13) into eq. (11), and we can write:

$$\hat{W}^l (Q_i^l - Z_i^l) C_i^l + \hat{b}^l = (Q_o^l - Z_o^l) C_o^l, \quad (14)$$

401 which gives:

$$\begin{aligned} Q_o^l &= \hat{W}^l \frac{C_i^l}{C_o^l} Q_i^l + \frac{\hat{b}^l}{C_o^l} + Z_o^l - \frac{\hat{W}^l Z_i^l C_i^l}{C_o^l} \\ &= \tilde{W}^l Q_i^l + \tilde{b}^l, \end{aligned} \quad (15)$$

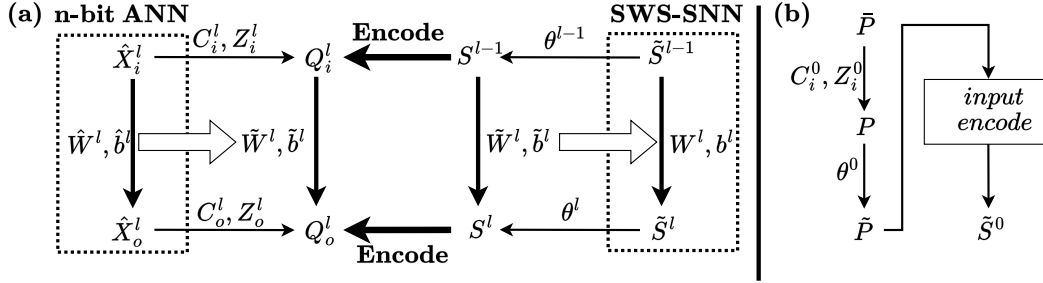


Figure 6: (a) Convert quantized ANNs to SWS-SNNs.  $Q_i^l$  and  $Q_o^l$  represent the integers to which  $\hat{X}_i^l$  and  $\hat{X}_o^l$  are mapped, respectively.  $\tilde{W}^l$  and  $\tilde{b}^l$  denotes the weight and bias to get  $Q_o^l$  from  $Q_i^l$ .  $W^l$  and  $b^l$  denotes the weight and bias in SWS-SNN. The process of transferring weights and biases from the quantized ANN to SWS-SNN is indicated by white arrows. The core of the conversion is that the distribution of the integer  $Q_o^l$  is known and can be easily encoded by  $S^l$ . (b) Process the input pixels to encode by pulses with an amplitude of 1.  $\bar{P}$  denotes the original pixel value,  $P$  denotes the mapped value and  $\tilde{P}$  denotes the value after scaled by  $1/\theta^0$ .

402 where

$$\tilde{W}^l = \hat{W}^l \frac{C_i^l}{C_o^l}, \quad (16)$$

$$\tilde{b}^l = \frac{\hat{b}^l}{C_o^l} + Z_o^l - \frac{\hat{W}^l Z_i^l C_i^l}{C_o^l}. \quad (17)$$

403 As seen in eq. (15), with the weight and bias set to  $\tilde{W}^l$  and  $\tilde{b}^l$  respectively, the layer outputs  $Q_o^l$  when  
 404 receiving  $Q_i^l$ . The pulse amplitude  $\theta^l$  can be set to any value as long as the codable range calculated  
 405 by eq. (7) covers  $[0, 2^n - 1]$ . Then we have:

$$W^l = \tilde{W}^l \frac{\theta^{l-1}}{\theta^l} = \hat{W}^l \frac{C_i^l \theta^{l-1}}{C_o^l \theta^l} \quad (18)$$

406 Considering the membrane potential amplification,  $b^l$  can be calculated as follows:

$$b^l = \frac{1}{\sum_{\tau=1}^{T_c} \beta^{T_c-\tau}} \tilde{b}^l = \frac{1}{\sum_{\tau=1}^{T_c} \beta^{T_c-\tau}} \left( \frac{\hat{b}^l}{C_o^l} + Z_o^l - \frac{\hat{W}^l Z_i^l C_i^l}{C_o^l} \right) \quad (19)$$

407 Once the  $T_c$ ,  $\beta$  and  $\theta^l$  ( $\theta^{l-1}$  is given by the previous layer) have been determined, all values on the  
 408 right side of eq. (18) and eq. (19) are known. Consequently,  $W_l$  and  $b^l$  in the SWS-SNN can be  
 409 readily calculated from the weight and bias of the quantized ANN.

410 After configuring the weights and biases as described above, the input pixel must be encoded into a  
 411 pulse sequence with an amplitude of 1 as well. This process is illustrated in fig. 6(b). First, map the  
 412 pixel value to  $[0, 2^n - 1]$  using  $C_i^0$  and  $Z_i^0$  obtained from QAT. Assuming this range can be encoded  
 413 by SWSs with an amplitude of  $\theta^0$ , scaling the pixel value by  $1/\theta^0$  allows the use of a sequence with  
 414  $\theta^0 = 1$  for encoding. Finally, encode the scaled pixels following section 3.4, and the required input  
 415 spike sequence is obtained.

## 416 A.2 Details for QAT

417 QAT is the quantization method that typically results in the highest accuracy. We basically follows  
 418 the workflow provided by PyTorch. The default QAT quantization configuration is chosen to specify  
 419 the kind of fake-quantization inserted after weights and activations. We choose Stochastic Gradient  
 420 Descent (SGD) optimizer in QAT, with the value of momentum set to 0.9 and the learning rate set to  
 421  $1 \times 10^{-4}$  since the weights only need to be fine-tuned. QAT is done for 12 epochs and 20 batches in  
 422 each epoch. We freeze the batch norm mean and variance estimates after three epochs and freeze the  
 423 quantizer parameters (scaling factor and zero point) after another two epochs.

## 424 **NeurIPS Paper Checklist**

### 425 **1. Claims**

426 Question: Do the main claims made in the abstract and introduction accurately reflect the  
427 paper's contributions and scope?

428 Answer: [\[Yes\]](#)

429 Justification: Stepwise weighting enhances the encoding of information in spikes, as is  
430 proved in eq. (5) in section 3.1. Our proposed SWS coding scheme achieves high perfor-  
431 mance and low energy consumption, which is supported by our experimental results in  
432 section 4. The TSA neuron model effectively minimizes the residual error, which can be  
433 proved from the ablation study in section 4.4.

434 Guidelines:

- 435 • The answer NA means that the abstract and introduction do not include the claims  
436 made in the paper.
- 437 • The abstract and/or introduction should clearly state the claims made, including the  
438 contributions made in the paper and important assumptions and limitations. A No or  
439 NA answer to this question will not be perceived well by the reviewers.
- 440 • The claims made should match theoretical and experimental results, and reflect how  
441 much the results can be expected to generalize to other settings.
- 442 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
443 are not attained by the paper.

### 444 **2. Limitations**

445 Question: Does the paper discuss the limitations of the work performed by the authors?

446 Answer: [\[Yes\]](#)

447 Justification: The inclusion of silent periods can lead to increased latency, as noted in  
448 section 3.3, which is a limitation we've found so far. However, our experimental results  
449 demonstrate that our delay performance still surpasses that of other SNNs.

450 Guidelines:

- 451 • The answer NA means that the paper has no limitation while the answer No means that  
452 the paper has limitations, but those are not discussed in the paper.
- 453 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 454 • The paper should point out any strong assumptions and how robust the results are to  
455 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
456 model well-specification, asymptotic approximations only holding locally). The authors  
457 should reflect on how these assumptions might be violated in practice and what the  
458 implications would be.
- 459 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
460 only tested on a few datasets or with a few runs. In general, empirical results often  
461 depend on implicit assumptions, which should be articulated.
- 462 • The authors should reflect on the factors that influence the performance of the approach.  
463 For example, a facial recognition algorithm may perform poorly when image resolution  
464 is low or images are taken in low lighting. Or a speech-to-text system might not be  
465 used reliably to provide closed captions for online lectures because it fails to handle  
466 technical jargon.
- 467 • The authors should discuss the computational efficiency of the proposed algorithms  
468 and how they scale with dataset size.
- 469 • If applicable, the authors should discuss possible limitations of their approach to  
470 address problems of privacy and fairness.
- 471 • While the authors might fear that complete honesty about limitations might be used by  
472 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
473 limitations that aren't acknowledged in the paper. The authors should use their best  
474 judgment and recognize that individual actions in favor of transparency play an impor-  
475 tant role in developing norms that preserve the integrity of the community. Reviewers  
476 will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [Yes]

Justification: The membrane potential amplification enhances the information-carrying capacity of the preceding pulses and is proved in eq. (5).

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We set specific random number seeds when conducting experiments to ensure that all the results of section 4 are reproducible.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility.

530 In the case of closed-source models, it may be that access to the model is limited in  
531 some way (e.g., to registered users), but it should be possible for other researchers  
532 to have some path to reproducing or verifying the results.

## 533 5. Open access to data and code

534 Question: Does the paper provide open access to the data and code, with sufficient instruc-  
535 tions to faithfully reproduce the main experimental results, as described in supplemental  
536 material?

537 Answer: [No]

538 Justification: Code will be released when the paper is accepted.

539 Guidelines:

- 540 • The answer NA means that paper does not include experiments requiring code.
- 541 • Please see the NeurIPS code and data submission guidelines ([https://nips.cc/  
542 public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 543 • While we encourage the release of code and data, we understand that this might not be  
544 possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not  
545 including code, unless this is central to the contribution (e.g., for a new open-source  
546 benchmark).
- 547 • The instructions should contain the exact command and environment needed to run to  
548 reproduce the results. See the NeurIPS code and data submission guidelines ([https://  
549 nips.cc/public/guides/CodeSubmissionPolicy](https://nips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 550 • The authors should provide instructions on data access and preparation, including how  
551 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 552 • The authors should provide scripts to reproduce all experimental results for the new  
553 proposed method and baselines. If only a subset of experiments are reproducible, they  
554 should state which ones are omitted from the script and why.
- 555 • At submission time, to preserve anonymity, the authors should release anonymized  
556 versions (if applicable).
- 557 • Providing as much information as possible in supplemental material (appended to the  
558 paper) is recommended, but including URLs to data and code is permitted.

## 559 6. Experimental Setting/Details

560 Question: Does the paper specify all the training and test details (e.g., data splits, hyper-  
561 parameters, how they were chosen, type of optimizer, etc.) necessary to understand the  
562 results?

563 Answer: [Yes]

564 Justification: The details for acquiring different delays and the OPF calculation method  
565 are provided in section 4.2 and section 4.3, respectively. The parameters used during QAT  
566 training is outlined in appendix A.2.

567 Guidelines:

- 568 • The answer NA means that the paper does not include experiments.
- 569 • The experimental setting should be presented in the core of the paper to a level of detail  
570 that is necessary to appreciate the results and make sense of them.
- 571 • The full details can be provided either with the code, in appendix, or as supplemental  
572 material.

## 573 7. Experiment Statistical Significance

574 Question: Does the paper report error bars suitably and correctly defined or other appropriate  
575 information about the statistical significance of the experiments?

576 Answer: [No]

577 Justification: We believe it is not necessary to include error bars in the results because each  
578 experimental result itself is already the average of a large number of tests (E.g., the test  
579 accuracy for an epoch is averaged over  $Num\_of\_batches \times Batch\_size$  input images,  
580 and is therefore very close to each other in every test epoch).

581 Guidelines:



- 582 • The answer NA means that the paper does not include experiments.
- 583 • The authors should answer "Yes" if the results are accompanied by error bars, confi-
- 584 dence intervals, or statistical significance tests, at least for the experiments that support
- 585 the main claims of the paper.
- 586 • The factors of variability that the error bars are capturing should be clearly stated (for
- 587 example, train/test split, initialization, random drawing of some parameter, or overall
- 588 run with given experimental conditions).
- 589 • The method for calculating the error bars should be explained (closed form formula,
- 590 call to a library function, bootstrap, etc.)
- 591 • The assumptions made should be given (e.g., Normally distributed errors).
- 592 • It should be clear whether the error bar is the standard deviation or the standard error
- 593 of the mean.
- 594 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 595 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 596 of Normality of errors is not verified.
- 597 • For asymmetric distributions, the authors should be careful not to show in tables or
- 598 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 599 error rates).
- 600 • If error bars are reported in tables or plots, The authors should explain in the text how
- 601 they were calculated and reference the corresponding figures or tables in the text.

## 602 8. Experiments Compute Resources

603 Question: For each experiment, does the paper provide sufficient information on the com-  
604 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
605 the experiments?

606 Answer: [No]

607 Justification: We found it difficult to quantify the computing resources used in every  
608 experiments.

609 Guidelines:

- 610 • The answer NA means that the paper does not include experiments.
- 611 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 612 or cloud provider, including relevant memory and storage.
- 613 • The paper should provide the amount of compute required for each of the individual
- 614 experimental runs as well as estimate the total compute.
- 615 • The paper should disclose whether the full research project required more compute
- 616 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 617 didn't make it into the paper).

## 618 9. Code Of Ethics

619 Question: Does the research conducted in the paper conform, in every respect, with the  
620 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

621 Answer: [Yes]

622 Justification: We have read the NeurIPS Code of Ethics and the research conducted in this  
623 paper conforms with it.

624 Guidelines:

- 625 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 626 • If the authors answer No, they should explain the special circumstances that require a
- 627 deviation from the Code of Ethics.
- 628 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 629 eration due to laws or regulations in their jurisdiction).

## 630 10. Broader Impacts

631 Question: Does the paper discuss both potential positive societal impacts and negative  
632 societal impacts of the work performed?

633 Answer: [NA]

634 Justification: There is no societal impact of the work performed.

635 Guidelines:

- 636 • The answer NA means that there is no societal impact of the work performed.
- 637 • If the authors answer NA or No, they should explain why their work has no societal
- 638 impact or why the paper does not address societal impact.
- 639 • Examples of negative societal impacts include potential malicious or unintended uses
- 640 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 641 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 642 groups), privacy considerations, and security considerations.
- 643 • The conference expects that many papers will be foundational research and not tied
- 644 to particular applications, let alone deployments. However, if there is a direct path to
- 645 any negative applications, the authors should point it out. For example, it is legitimate
- 646 to point out that an improvement in the quality of generative models could be used to
- 647 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 648 that a generic algorithm for optimizing neural networks could enable people to train
- 649 models that generate Deepfakes faster.
- 650 • The authors should consider possible harms that could arise when the technology is
- 651 being used as intended and functioning correctly, harms that could arise when the
- 652 technology is being used as intended but gives incorrect results, and harms following
- 653 from (intentional or unintentional) misuse of the technology.
- 654 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 655 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 656 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 657 feedback over time, improving the efficiency and accessibility of ML).

## 658 11. Safeguards

659 Question: Does the paper describe safeguards that have been put in place for responsible  
660 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
661 image generators, or scraped datasets)?

662 Answer: [NA]

663 Justification: The paper poses no such risks.

664 Guidelines:

- 665 • The answer NA means that the paper poses no such risks.
- 666 • Released models that have a high risk for misuse or dual-use should be released with
- 667 necessary safeguards to allow for controlled use of the model, for example by requiring
- 668 that users adhere to usage guidelines or restrictions to access the model or implementing
- 669 safety filters.
- 670 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 671 should describe how they avoided releasing unsafe images.
- 672 • We recognize that providing effective safeguards is challenging, and many papers do
- 673 not require this, but we encourage authors to take this into account and make a best
- 674 faith effort.

## 675 12. Licenses for existing assets

676 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
677 the paper, properly credited and are the license and terms of use explicitly mentioned and  
678 properly respected?

679 Answer: [Yes]

680 Justification: The pretrained ANN model and the QAT workflow is provided by PyTorch  
681 and we cited the original paper in appendix A.1 as [22].

682 Guidelines:

- 683 • The answer NA means that the paper does not use existing assets.
- 684 • The authors should cite the original paper that produced the code package or dataset.
- 685 • The authors should state which version of the asset is used and, if possible, include a  
686 URL.

- 687 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 688 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 689 service of that source should be provided.
- 690 • If assets are released, the license, copyright information, and terms of use in the
- 691 package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets)
- 692 has curated licenses for some datasets. Their licensing guide can help determine the
- 693 license of a dataset.
- 694 • For existing datasets that are re-packaged, both the original license and the license of
- 695 the derived asset (if it has changed) should be provided.
- 696 • If this information is not available online, the authors are encouraged to reach out to
- 697 the asset's creators.

### 698 13. **New Assets**

699 Question: Are new assets introduced in the paper well documented and is the documentation  
700 provided alongside the assets?

701 Answer: [NA]

702 Justification: The paper does not release new assets.

703 Guidelines:

- 704 • The answer NA means that the paper does not release new assets.
- 705 • Researchers should communicate the details of the dataset/code/model as part of their
- 706 submissions via structured templates. This includes details about training, license,
- 707 limitations, etc.
- 708 • The paper should discuss whether and how consent was obtained from people whose
- 709 asset is used.
- 710 • At submission time, remember to anonymize your assets (if applicable). You can either
- 711 create an anonymized URL or include an anonymized zip file.

### 712 14. **Crowdsourcing and Research with Human Subjects**

713 Question: For crowdsourcing experiments and research with human subjects, does the paper  
714 include the full text of instructions given to participants and screenshots, if applicable, as  
715 well as details about compensation (if any)?

716 Answer: [NA]

717 Justification: The paper does not involve crowdsourcing nor research with human subjects.

718 Guidelines:

- 719 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 720 human subjects.
- 721 • Including this information in the supplemental material is fine, but if the main contribu-
- 722 tion of the paper involves human subjects, then as much detail as possible should be
- 723 included in the main paper.
- 724 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 725 or other labor should be paid at least the minimum wage in the country of the data
- 726 collector.

### 727 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human** 728 **Subjects**

729 Question: Does the paper describe potential risks incurred by study participants, whether  
730 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)  
731 approvals (or an equivalent approval/review based on the requirements of your country or  
732 institution) were obtained?

733 Answer: [NA]

734 Justification: The paper does not involve crowdsourcing nor research with human subjects.

735 Guidelines:

- 736 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 737 human subjects.

738  
739  
740  
741  
742  
743  
744  
745

- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.