

Towards High-performance Spiking Transformers from ANN to SNN Conversion

MM '24: Proceedings of the 32nd ACM International Conference on Multimedia

<https://github.com/h-z-h-cell/Transformer-to-SNN-ECMT>

Zihan Huang
Peking University
Beijing, China
hzh@stu.pku.edu.cn

Tong Bu
Peking University
Beijing, China
putong30@pku.edu.cn

Xinyu Shi
Peking University
Beijing, China
xyshi@pku.edu.cn

Jianhao Ding
Peking University
Beijing, China
djh01998@stu.pku.edu.cn

Zecheng Hao
Peking University
Beijing, China
1900012989@pku.edu.cn

Zhaofei Yu*
Peking University
Beijing, China
yuzf12@pku.edu.cn

Tiejun Huang
Peking University
Beijing, China
tjhuang@pku.edu.cn

Abstract

Spiking neural networks (SNNs) show great potential due to their energy efficiency, fast processing capabilities, and robustness. There are two main approaches to constructing SNNs. Direct training methods require much memory, while conversion methods offer a simpler and more efficient option. However, current conversion methods mainly focus on converting convolutional neural networks (CNNs) to SNNs. Converting Transformers to SNN is challenging because of the presence of non-linear modules. In this paper, we propose an **Expectation Compensation Module** to preserve the accuracy of the conversion. The core idea is to use information from the previous T time-steps to calculate the expected output at time-step T . We also propose a **Multi-Threshold Neuron** and the corresponding **Parallel Parameter normalization** to address the challenge of large time steps needed for high accuracy, aiming to reduce network latency and power consumption. Our experimental results demonstrate that our approach achieves state-of-the-art performance. For example, we achieve a top-1 accuracy of 88.60% with only a 1% loss in accuracy using 4 time steps while consuming only 35% of the original power of the Transformer. To our knowledge, this is the first successful Artificial Neural Network (ANN) to SNN conversion for Spiking Transformers that achieves high accuracy, low latency, and low power consumption on complex datasets. The source codes of the proposed method are available at <https://github.com/h-z-h-cell/Transformer-to-SNN-ECMT>.

摘要:

- 1.SNN有诸多优势
- 2.构建SNN的两种方法，各有利弊（主说Conversion弊端）

➤ **Transformer Conversion的挑战：** **非线性层/模块的存在**

- 3.本文提出的方法1：**Expectation Compensation Module**
用前 T 个时间步计算第 T 个时间的期望输出
- 4.本文提出的方法2：**Multi-Threshold Neuron**
- 5.本文提出的方法3：**Parallel Parameter Normalization**
效果：减少SNN推理延迟和能耗
- 6.实验结果.....
- 7.?

Introduction

Para.1

1. What is SNN?
2. Spike-based neuron dynamics: **Sparse**
3. SNN is hard to train.

Para.2

Training Strategies:

1. Direct training: surrogate gradient or STDP
2. ANN2SNN: similarity between ANN $\text{ReLU}(\cdot)$ and firing rate; **More timesteps, More Precision**

Para.3

Transformers are powerful, but no previous ANN2SNN works on such model.

- LayerNorm & GELU  require interaction between neurons **in the same layer**.
- Hard to apply piecewise quantization of individual neurons ?
- Non-linear characteristics

Ideally, layerwise
neurons are I.I.D maybe?

Introduction

Para.4

New Methods:

1. Expectation Compensation Module(ECM):



2. Multi-Threshold Neurons & Parallel Parameter normalization:

Irrelevant?

Para.5

Summary of Contributions:

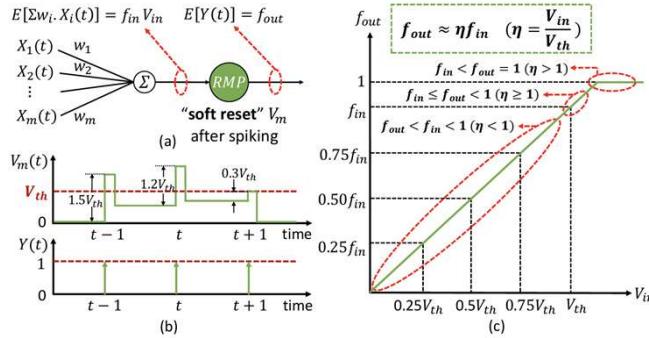
Related Works

Sec.1: ANN-SNN Conversion

2015 Spiking Deep Convolutional Neural Networks for Energy-Efficient Object Recognition.

2015 Fast-Classifying, High-Accuracy Spiking Deep Networks Through Weight and Threshold Balancing

2017 Conversion of Continuous-Valued Deep Networks to Efficient Event-Driven Networks for Image Classification



I. Optimizing thresholds

2019 Going Deeper in Spiking Neural Networks: VGG and Residual Architectures

2023 Low Latency and Sparse Computing Spiking Neural Networks With Self-Driven Adaptive Threshold Plasticity

1. Make output values in all layers positive:
 - a. Add an abs() function (absolute value) layer after preprocessing phase that includes color transformation and spatial normalization. The abs() function ensures that the input values to the first convolution layer are all non-negative. This step can be omitted if the preprocessing method itself gives positive output (e.g., the original RGB values are used);
 - b. Change the sigmoid activation function (after convolution) from tanh() to HalfRect(). HalfRect(x), also referred to as rectified linear units (ReLU), is defined as:

$$\text{HalfRect}(x) = \max(x, 0)$$

1. Make output values in all layers positive:
 - a. Add an abs() function (absolute value) layer after preprocessing phase that includes color transformation and spatial normalization. The abs() function ensures that the input values to the first convolution layer are all non-negative. This step can be omitted if the preprocessing method itself gives positive output (e.g., the original RGB values are used);
 - b. Change the sigmoid activation function (after convolution) from tanh() to HalfRect(). HalfRect(x), also referred to as rectified linear units (ReLU), is defined as:

$$\text{HalfRect}(x) = \max(x, 0)$$

Algorithm 1: Model-Based Normalization

```

1 for layer in layers:
2     max_pos_input = 0
3     # Find maximum input for this layer
4     for neuron in layer.neurons:
5         input_sum = 0
6         for input_wt in neuron.input_wts:
7             input_sum += max(0, input_wt)
8         max_pos_input = max(max_pos_input, input_sum)
9     # Rescale all weights
10    for neuron in layer.neurons:
11        for input_wt in neuron.input_wts:
12            input_wt = input_wt / max_pos_input

```

CSDN @Selena Lau

Algorithm 2: Data-Based Normalization

```

1 || previous_factor = 1
2 for layer in layers:
3     max_wt = 0
4     max_act = 0
5     for neuron in layer.neurons:
6         for input_wt in neuron.input_wts:
7             max_wt = max(max_wt, input_wt)
8             max_act = max(max_act, neuron.output_act)
9     scale_factor = max(max_wt, max_act)
10    applied_factor = scale_factor / previous_factor
11    # Rescale all weights
12    for neuron in layer.neurons:
13        for input_wt in neuron.input_wts:
14            input_wt = input_wt / applied_factor
15    previous_factor = scale_factor

```

CSDN @Selena Lau

$$r_i^1(t) = \begin{cases} a_i^1 r_{\max} \cdot \frac{V_{\text{thr}}}{V_{\text{thr}} + \epsilon_i^1} - \frac{V_i^1(t)}{t \cdot (V_{\text{thr}} + \epsilon_i^1)} & \text{reset to zero} \\ a_i^1 r_{\max} & - \frac{V_i^1(t)}{t \cdot V_{\text{thr}}} \end{cases} \quad (5a)$$

$$\text{reset by subtraction.} \quad (5b)$$

2020 RMP-SNN: Residual Membrane Potential Neuron for Enabling Deeper High-Accuracy and Low-Latency Spiking Neural Network

```

Algorithm 1: SPIKE-NORM
Input: Input Poisson Spike Train spikes, Number of Time-Steps #timesteps
Output: Weight-normalization / Threshold-balancing factors v_h_norm[i] for each neural layer (net.layer[i]) of the network net
1 initialization v_h_norm[i] = 0 \forall i = 1,...,net.layer;
2 // Set input of 1st layer equal to spike train
3 net.layer[1].input = spikes;
4 for t ← 1 to net.layer do
5     for i ← 1 to net.layer do
6         // Forward pass spike-train for neuron layer i characterized by membrane potential net.layer[i].v_norm and threshold net.layer[i].v_th
7         net.layer[i].forward(net.layer[i].input);
8         // Determine threshold-balancing factor according to maximum SNN activation,
9         max(net.layer[i].weight * net.layer[i].input[]),
10        where '*' represents the dot-product operation
11        v_h_norm[i] = max(v_h_norm[i], max(net.layer[i].weight * net.layer[i].input[]));
12    end
13    // Threshold-balance layer-
14    net.layer[i].v_h = v_h_norm[i];
15    // Record input spike-train for next layer
16    net.layer[i+1].input = net.layer[i].forward(net.layer[i].input);
17 end

```

```

Algorithm 1 SATP
Input: Inference Images  $\mathcal{X} = \{x_i\}_{i=1}^n$ .
Output: Predicted label  $C = \{c_i\}_{i=1}^n$  for each image.
Set the parameters  $m$ ,  $\lambda$ ,  $\sigma$  and  $\eta$ , initialize  $V_{th}$ ;
for time  $t$  in  $[1, 2, \dots, T]$  do
    for layer  $l$  in layers do
        for neuron  $p$  in layer.neurons do
            Calculate the input current  $I(t)$  by Eq. (4);
            Calculate  $\nabla v(t)$  by Eq. (2);
            if  $V(t) > V_{th}(t)$  then
                Neuron  $p$  fires;
                Update the membrane potential  $V(t)$  by Eq. (3);
                Calculate  $\Delta V_{th}(t)$  by Eq. (7);
                Update  $V_{th}(t) \leftarrow V_{th}(t) + \Delta V_{th}(t)$ .
            end
        end
    end
    Collect the spikes from the output layer;
    Obtain the classification result at time  $t$ .

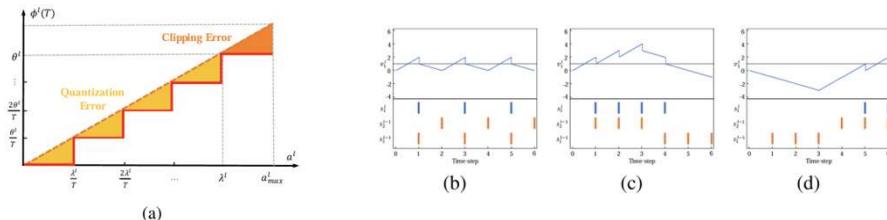
```

Related Works

Sec.1: ANN-SNN Conversion

II. Optimizing membrane potentials

2022 Optimized Potential Initialization for Low-latency Spiking Neural Networks



III. Optimizing pre-conversion ANN structure

2021 TCL: an ANN-to-SNN Conversion with Trainable Clipping Layers

Rate Norm Layer

$$\theta_l = p_l \cdot \max(W_{l-1}\hat{r}_{l-1} + b_{l-1}),$$

$$z_l = \text{clip}(W_{l-1}\hat{r}_{l-1} + b_{l-1}, 0, \theta_l),$$

$$\hat{r}_l = \frac{z_l}{\theta_l},$$

2022 Optimal ANN-SNN Conversion for High-accuracy and Ultra-low-latency Spiking Neural Networks

2023 Symmetric-threshold ReLU for Fast and Nearly Lossless ANN-SNN Conversion

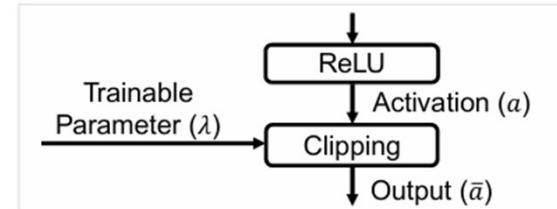
Optimal initialization of membrane potentials

Theorem 1. The expectation of square conversion error (Eq. (13)) reaches the minimum value when the initial value $v^l(0)$ is $V_{th}^l/2$, meanwhile the expectation of conversion error reaches 0, that is:

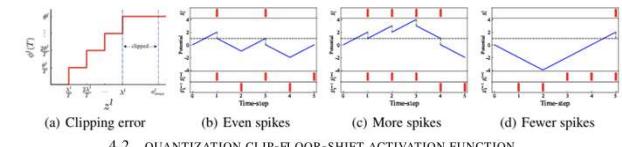
$$\arg \min_{v^l(0)} E_z \|f(z) - f'(z)\|_2^2 = \frac{V_{th}^l}{2}, \quad (14)$$

$$E_z (f(z) - f'(z)) \Big|_{v^l(0)=\frac{V_{th}^l}{2}} = 0. \quad (15)$$

2023 Reducing ANN-SNN Conversion Error through Residual Membrane Potential

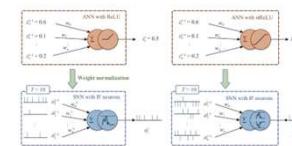


2021 Optimal ANN-SNN Conversion for Fast and Accurate Inference in Deep Spiking Neural Networks



We propose the quantization clip-floor-shift activation function to train ANNs.

$$a^l = \hat{h}(z^l) = \lambda^l \text{clip}\left(\frac{1}{L} \left[\frac{z^l L}{\lambda^l} + \varphi \right], 0, 1\right).$$



$$z^l = \text{stReLU}(\mathbf{W}^l z^{l-1} + b^l)$$

$$\text{stReLU}(x) = \begin{cases} V_{th}, & \text{if } x \geq V_{th} \\ x, & \text{if } -V_{th} < x < V_{th} \\ -V_{th}, & \text{if } x \leq -V_{th}. \end{cases}$$

Related Works

Sec.1: ANN-SNN Conversion

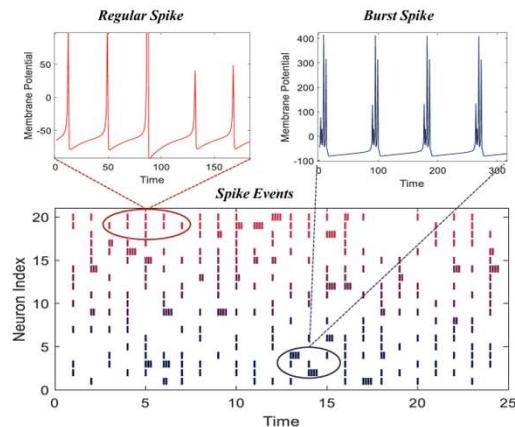
III. Optimizing pre-conversion ANN structure

2023 A Unified Optimization Framework of ANN-SNN Conversion: Towards Optimal Mapping from Activation Values to Firing Rates

2023 A New ANN-SNN Conversion Method with High Accuracy, Low Latency and Good Robustness

IV. Optimizing spiking neuronal models

2022 Efficient and Accurate Conversion of Spiking Neural Network with Burst Spikes



2022 Signed Neuron with Memory: Towards Simple, Accurate and High-Efficient ANN-SNNConversion

$$\tilde{v}_j^l(t) = v_j^l(t-1) + \sum_i w_{ij} s_i^{l-1}(t) + b_j, \quad (7)$$

$$\tilde{m}_j^l(t) = m_j^l(t-1), \quad (8)$$

$$s_j^l(t) = \begin{cases} \theta_j^l, & \tilde{v}_j^l(t) \geq \theta_j^l \\ -\theta_j^l, & \tilde{v}_j^l(t) \leq -\theta_j^l \quad \text{and} \quad \tilde{m}_j^l(t) > 0 \\ 0, & \text{otherwise} \end{cases}, \quad (9)$$

$$v_j^l(t) = \tilde{v}_j^l(t) - s_j^l(t), \quad (10)$$

$$m_j^l(t) = \tilde{m}_j^l(t) + s_j^l(t), \quad (11)$$

B.3. Special Cases of the SlipReLU Activation Function

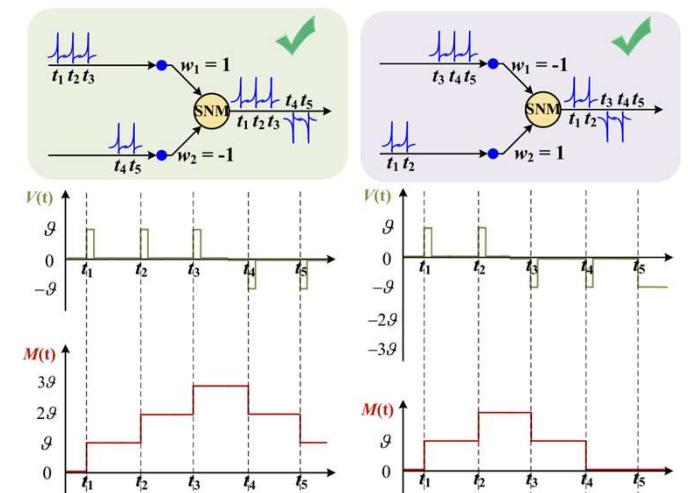
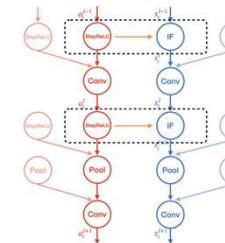
Here we list four different special cases of the proposed SlipReLU.

Threshold-ReLU When $c = 1$ and $\delta_1 = 0$, the SlipReLU becomes the threshold ReLU activation function which is studied in [Deng & Gu \(2021\)](#).

Shift-threshold-ReLU When $c = 1$ and $\delta_1 = -(1/(2N))$, the SlipReLU becomes the shift-threshold ReLU activation function which is studied in [Deng & Gu \(2021\)](#).

Quantization clip-floor (QCF) When $c = 0$ and $\delta = 0$, the SlipReLU becomes the quantization clip-floor (QCF) activation function which is studied in [Bu et al. \(2021\)](#).

Quantization clip-floor-shift (QCFs) When $c = 0$ and $\delta = 1/2$, the SlipReLU becomes the quantization clip-floor-shift (QCFs) activation function which is studied in [Bu et al. \(2021\)](#).



Related Works

Sec.1: ANN-SNN Conversion

IV. Optimizing spiking neuronal models

2022 Quantization Framework for Fast Spiking Neural Networks

V. Others:

2024 Spatio-Temporal Approximation: A Training-Free SNN Conversion for Transformers

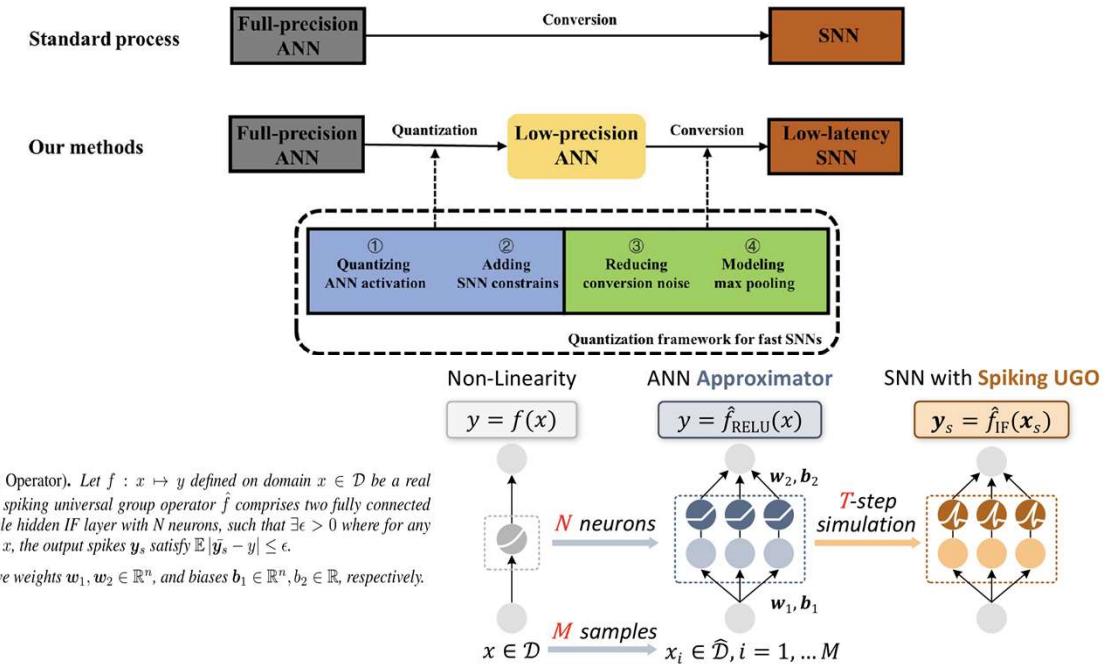


Figure 2: Spatial approximation process with UGO.

Sec.2: Direct Training for Transformer

Preliminaries

Sec.1: Theories of ANN-SNN Conversion

ANN Neuron: $\mathbf{a}^l = \text{ReLU}(\mathbf{W}^l \mathbf{a}^{l-1}) = \max(\mathbf{W}^l \mathbf{a}^{l-1}, 0),$

IF Neuron: $\mathbf{m}^l(t) = \mathbf{v}^l(t-1) + \mathbf{W}^l \mathbf{x}^{l-1}(t),$

$$\mathbf{s}^l(t) = H(\mathbf{m}^l(t) - \theta^l),$$

$$\mathbf{x}^l(t) = \theta^l \mathbf{s}^l(t),$$

$$\mathbf{v}^l(t) = \mathbf{m}^l(t) - \mathbf{x}^l(t). \quad \text{reset-by-subtraction}$$

Recurrent IF Neuron: $\mathbf{v}^l(t) - \mathbf{v}^l(t-1) = \mathbf{W}^l \mathbf{x}^{l-1}(t) - \mathbf{x}^l(t).$

Recurrent IF Neuron $\mathbf{v}^l(t) - \mathbf{v}^l(t-1) = \mathbf{W}^l \mathbf{x}^{l-1}(t) - \mathbf{x}^l(t).$

summed through T : $\mathbf{v}^l(t-1) - \mathbf{v}^l(t-2) = \mathbf{W}^l \mathbf{x}^{l-1}(t-1) - \mathbf{x}^l(t-2).$

.....

$$\frac{\mathbf{v}^l(T) - \mathbf{v}^l(0)}{T} = \frac{\mathbf{W}^l \sum_{i=1}^T \mathbf{x}^{l-1}(i)}{T} - \frac{\sum_{i=1}^T \mathbf{x}^l(i)}{T}.$$

Preliminaries

Sec.1: Theories of ANN-SNN Conversion

**Recurrent IF Neuron
summed through T :**

$$\frac{\mathbf{v}^l(T) - \mathbf{v}^l(0)}{T} = \frac{\mathbf{W}^l \frac{\sum_{i=1}^T \mathbf{x}^{l-1}(i)}{T} - \frac{\sum_{i=1}^T \mathbf{x}^l(i)}{T}}{T}.$$

$$\Phi^l(T) = \frac{\sum_{i=1}^T \mathbf{x}^l(i)}{T}$$

$$\Phi^l(T) = \mathbf{W}^l \Phi^{l-1}(T) - \frac{\mathbf{v}^l(T) - \mathbf{v}^l(0)}{T}.$$

Parameter normalization

Algorithm 1: Model-Based Normalization

$$W_{\text{SNN}}^l = W_{\text{ANN}}^l \frac{\lambda^{l-1}}{\lambda^l}.$$

Access all weights **for** input_wt **in** neuron.input_wts:
 Sum of weights(ignore neg. value) input_sum += **max**(0, input_wt)
 Update **max_pos_input** max_pos_input = **max**(max_pos_input, input_sum)
 # Rescale all weights
for neuron **in** layer.neurons:
for input_wt **in** neuron.input_wts:
 input_wt = input_wt / max_pos_input

Algorithm 2: Data-Based Normalization

```

1 previous_factor = 1
2 for layer in layers:
3     max_wt = 0
4     max_act = 0
5     for neuron in layer.neurons:
6         for input_wt in neuron.input_wts:
7             max_wt = max(max_wt, input_wt)
8             max_act = max(max_act, neuron.output_act)
9         scale_factor = max(max_wt, max_act)
10        applied_factor = scale_factor / previous_factor
11        # Rescale all weights
12        for neuron in layer.neurons:
13            for input_wt in neuron.input_wts:
14                input_wt = input_wt / applied_factor
15        previous_factor = scale_factor

```

Modifying Equation(9) and setting θ_j^l to 1 is equivalent to adjusting the firing threshold on the soft-reset neuron to λ^l

Methods

Sec.1 Error Analysis of Nonlinear Module in ANN-SNN Conversion

GeLU

Softmax

LN

Attention

matrix
product?

assume that the outputs of layer $l - 1$ in both ANNs and SNNs are identical, $a^{l-1} = \Phi^{l-1}(T) = \frac{\sum_{t=1}^T \mathbf{x}^{l-1}(t)}{T}$
Compare the outputs of layer l : a^l and Φ^l

An arbitrary non-linear module $F(\cdot)$, in ANN $a^l = F(a^{l-1})$,

The same arbitrary non-linear module $F(\cdot)$, in counterpart SNN $x^l(t) = F(x^{l-1}(t))$

Average output (0-T): $\Phi^l(T) = \frac{\sum_{t=1}^T x^l(t)}{T} = \frac{\sum_{t=1}^T F(x^{l-1}(t))}{T}$.

output in ANN: $a^l = F(a^{l-1}) = F\left(\frac{\sum_{t=1}^T x^{l-1}(t)}{T}\right)$

$$\frac{\sum_{t=1}^T F(x^{l-1}(t))}{T} \neq F\left(\frac{\sum_{t=1}^T x^{l-1}(t)}{T}\right)$$

Methods

Sec.2 Expectation Compensation Module

$$\frac{\sum_{t=1}^T F(x^{l-1}(t))}{T} \neq F\left(\frac{\sum_{t=1}^T x^{l-1}(t)}{T}\right)$$

output in ANN: $a^l = F(a^{l-1}) = F\left(\frac{\sum_{t=1}^T x^{l-1}(t)}{T}\right)$

$$a^l = \Phi^l(T) = \frac{\sum_{i=1}^T x^l(i)}{T}$$

General Expectation Compensation Module.

THEOREM 4.1. Consider a non-linear layer l with a function F . In SNNs, the output of this layer at time t is denoted as $O^l(t)$. Let $S^l(T)$ be the cumulative sum of layer l outputs up to time T , given by $S^l(T) = \sum_{t=1}^T O^l(t)$. The **expected output** of the SNNs at time T is given by:

$$O^l(T) = TF\left(\frac{S^{l-1}(T)}{T}\right) - (T-1)F\left(\frac{S^{l-1}(T-1)}{T-1}\right). \quad (15)$$

$$\begin{aligned} O^l(T) &= \sum_{t=1}^T O^l(t) - \sum_{t=1}^{T-1} O^l(t) \\ &= Ta_T^l - (T-1)a_{T-1}^l \quad \boxed{T\phi_T^l - (T-1)\phi_{T-1}^l} \\ &= TF(a_T^{l-1}) - (T-1)F(a_{T-1}^{l-1}) \\ &= TF\left(\frac{\sum_{t=1}^T O^{l-1}(t)}{T}\right) - (T-1)F\left(\frac{\sum_{t=1}^{T-1} O^{l-1}(t)}{T-1}\right) \\ &= TF\left(\frac{S^{l-1}(T)}{T}\right) - (T-1)F\left(\frac{S^{l-1}(T-1)}{T-1}\right) \end{aligned}$$

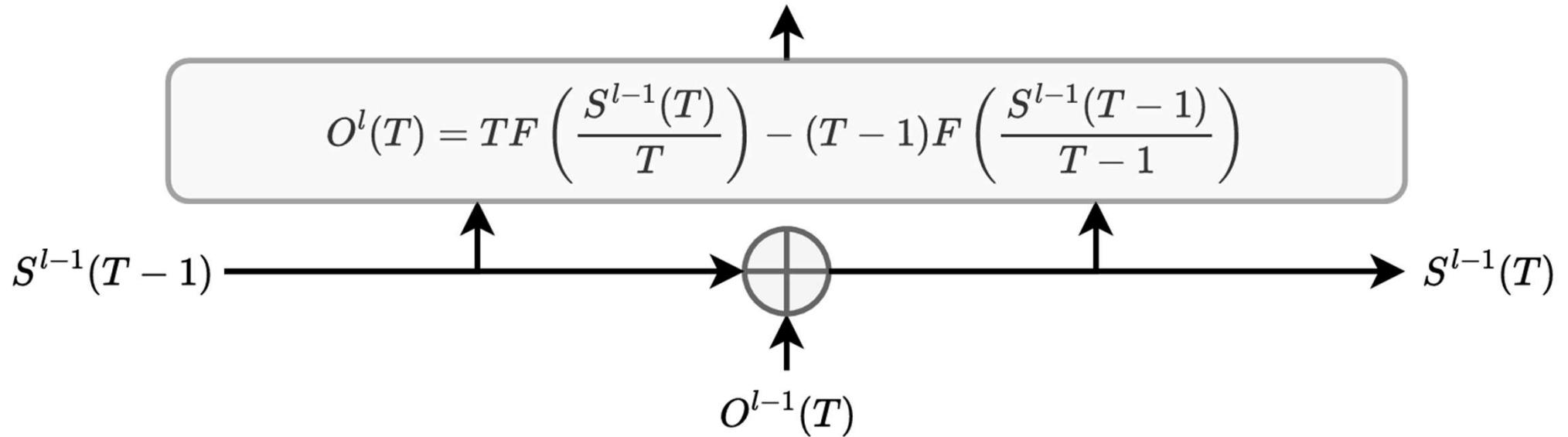
Methods

Sec.2 Expectation Compensation Module

$$O^l(T) = TF\left(\frac{S^{l-1}(T)}{T}\right) - (T-1)F\left(\frac{S^{l-1}(T-1)}{T-1}\right). \quad (15)$$

Expectation Compensation (EC)

$O^l(T)$



Methods

Sec.2 Expectation Compensation Module

Expectation Compensation Module for Matrix Product.

THEOREM 4.2. Consider a module for matrix product that receives two sets of spike inputs, denoted by $A_{v_a}(t)$ and $B_{v_b}(t)$. These inputs are generated by neurons A and B, respectively, and are characterized by multiple thresholds v_a and v_b , as described in Section 4.3.

We can integrate the input by $A(t) = \sum_{v_a} v_a A_{v_a}(t)$ and $B(t) = \sum_{v_b} v_b B_{v_b}(t)$. Here, $A(t)$ and $B(t)$ are the sum matrices weighted by multiple thresholds v_a and v_b , respectively.

Let $S_A(T) = \sum_{t=1}^T A(t)$ and $S_B(T) = \sum_{t=1}^T B(t)$ represent the cumulative sum of inputs up to time T. We define $S_K(T) = S_A(T)S_B(T)$. Then, the expected output at time T can be formulated as:

$$O(T) = \frac{1}{T} S_K(T) - \frac{1}{T-1} S_K(T-1), \quad (16)$$

where $S_K(T)$ can be calculated mainly using addition, as described by the following equation:

$$S_K(T) = S_K(T-1) + K(T) \quad (17)$$

$$\begin{aligned} K(T) &= \sum_{v_a, v_b} v_a v_b A_{v_a}(T) B_{v_b}(T) + \sum_{v_a} v_a A_{v_a}(T) S_B(T-1) \\ &\quad + \sum_{v_b} v_b S_A(T-1) B_{v_b}(T). \end{aligned} \quad (18)$$

$$\begin{aligned} A_T &= \frac{\sum_{t=1}^T A(t)}{T} \\ B_T &= \frac{\sum_{t=1}^T B(t)}{T} \\ O_T &= \frac{\sum_{t=1}^T O(t)}{T} \end{aligned}$$

approximate the value of ANNs using the mean value for the first T times in SNNs,

$$\begin{aligned} O(T) &= \sum_{t=1}^T O(t) - \sum_{i=t}^{T-1} O(t) \\ &= T O_T - (T-1) O_{T-1} \\ &= T A_T B_T - (T-1) A_{T-1} B_{T-1} \\ &= T \frac{\sum_{t=1}^T A(t)}{T} \frac{\sum_{t=1}^T B(t)}{T} \\ &\quad - (T-1) \frac{\sum_{t=1}^{T-1} A(t)}{T-1} \frac{\sum_{t=1}^{T-1} B(t)}{T-1} \\ &= \frac{1}{T} \sum_{t=1}^T A(t) \sum_{t=1}^T B(t) - \frac{1}{(T-1)} \sum_{t=1}^{T-1} A(t) \sum_{t=1}^{T-1} B(t) \\ &= \frac{1}{T} S_A(T) S_B(T) - \frac{1}{T-1} S_A(T-1) S_B(T-1) \\ &= \frac{1}{T} S_K(T) - \frac{1}{T-1} S_K(T-1) \end{aligned}$$

Methods

Sec.2 Expectation Compensation Module

Expectation Compensation Module for Matrix Product.

THEOREM 4.2. Consider a module for matrix product that receives two sets of spike inputs, denoted by $A_{v_a}(t)$ and $B_{v_b}(t)$. These inputs are generated by neurons A and B, respectively, and are characterized by multiple thresholds v_a and v_b , as described in Section 4.3.

We can integrate the input by $A(t) = \sum_{v_a} v_a A_{v_a}(t)$ and $B(t) = \sum_{v_b} v_b B_{v_b}(t)$. Here, $A(t)$ and $B(t)$ are the sum matrices weighted by multiple thresholds v_a and v_b , respectively.

Let $S_A(T) = \sum_{t=1}^T A(t)$ and $S_B(T) = \sum_{t=1}^T B(t)$ represent the cumulative sum of inputs up to time T. We define $S_K(T) = S_A(T)S_B(T)$. Then, the expected output at time T can be formulated as:

$$O(T) = \frac{1}{T} S_K(T) - \frac{1}{T-1} S_K(T-1), \quad (16)$$

where $S_K(T)$ can be calculated mainly using addition, as described by the following equation:

$$S_K(T) = S_K(T-1) + K(T) \quad (17)$$

$$\begin{aligned} K(T) &= \sum_{v_a, v_b} v_a v_b A_{v_a}(T) B_{v_b}(T) + \sum_{v_a} v_a A_{v_a}(T) S_B(T-1) \\ &\quad + \sum_{v_b} v_b S_A(T-1) B_{v_b}(T). \end{aligned} \quad (18)$$

$$\begin{aligned} S_K(T) &= S_A(T)S_B(T) \\ &= (\sum_{t=1}^T A(t))(\sum_{t=1}^T B(t)) \\ &= (S_A(T-1) + A(T))(S_B(T-1) + B(T)) \\ &= S_A(T-1)S_B(T-1) + A(T)B(T) \\ &\quad + A(T)S_B(T-1) + S_A(T-1)B(T) \\ &= S_K(T-1) + \sum_{v_a, v_b} v_a v_b A_{v_a}(T) B_{v_b}(T) \\ &\quad + \sum_{v_a} v_a A_{v_a}(T) S_B(T-1) + \sum_{v_b} v_b S_A(T-1) B_{v_b}(T) \\ &= S_K(T-1) + K(T). \end{aligned}$$

η_1 and η_2 are the firing rate of $A(T)$ and $B(T)$

$$ACs_{\text{SNN}}^{\max} = \eta_1 \eta_2 npm + \eta_1 npm + \eta_2 npm + 3nm$$

$$MACs_{\text{SNN}}^{\max} = \min(\eta_1, \eta_2) nm + \eta_1 nm + \eta_2 nm$$

$$ACs_{\text{SNN}}^{\max} \gg MACs_{\text{SNN}}^{\max}$$

Methods

Sec.2 Expectation Compensation Module

Multi-Threshold Neuron

If we only use the Expectation Compensation Module neuron communication will remain in a floating-point format.

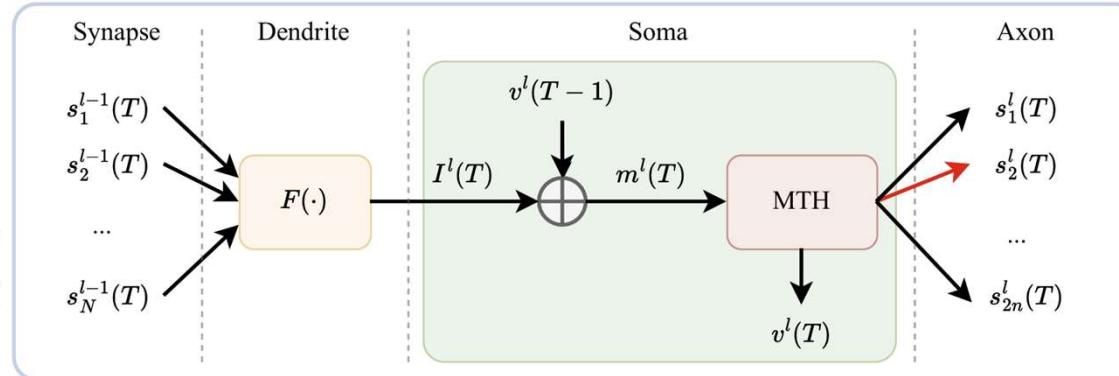
Introduce spiking neurons **before each linear layer and matrix product layer.**

Traditional SNN neurons takes several timesteps to generate at most one spike. →**Latency and power consumption**

Base threshold

$$\lambda_1^l = \theta_1^l, \lambda_2^l = 2\theta_1^l, \dots, \lambda_n^l = 2^{n-1}\theta_1^l,$$

$$\lambda_{n+1}^l = -\theta_2^l, \lambda_{n+2}^l = -2\theta_2^l, \dots, \lambda_{2n}^l = -2^{n-1}\theta_2^l,$$



$$\text{输入电流 } I_j^l(t) = F_j^l(s_{,1}^{l-1}(t), \dots, s_{,2n}^{l-1}(t)),$$

$$\text{膜电位更新 } m_j^l(t) = v_j^l(t-1) + I_j^l(t),$$

$$\text{脉冲输出 } s_{j,p}^l(t) = MTH_{\theta_1, \theta_2, n}(m_j^l(t))$$

$$\text{脉冲输出 } x_j^l(t) = \sum_p s_{j,p}^l(t) \lambda_p^l, \\ \text{的总和}$$

$$\text{膜电位重置 } v_j^l(t) = m_j^l(t) - x_j^l(t).$$

Methods

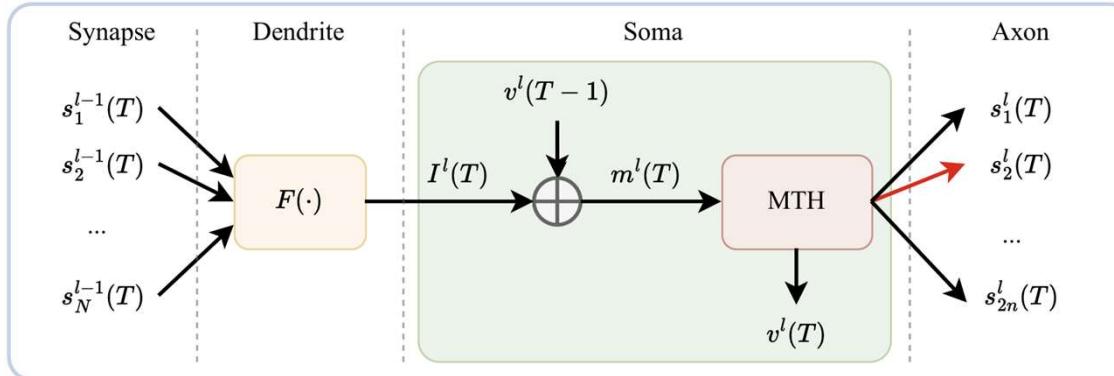
Sec.2 Expectation Compensation Module

Multi-Threshold Neuron

$$s_{j,p}^l(t) = MTH_{\theta_1, \theta_2, n}(m_j^l(t))$$

$MTH_{\theta_1, \theta_2, n}(x)$:

$$\left[\begin{array}{ll} \lambda_n^l - \frac{\lambda_1^l}{2} \leq x : & s_{j,n}^l(t) = 1, \\ \lambda_{n-1}^l - \frac{\lambda_1^l}{2} \leq x < \lambda_n^l - \frac{\lambda_1^l}{2} : & s_{j,n-1}^l(t) = 1, \\ \dots & \dots \\ \frac{\lambda_1^l}{2} \leq x < \lambda_2^l - \frac{\lambda_1^l}{2} : & s_{j,1}^l(t) = 1, \\ \text{未跨越阈值} & \text{all the } s_{j,p}^l(t) = 0, \\ \frac{\lambda_{n+1}^l}{2} \leq x < \frac{\lambda_1^l}{2} : & s_{j,n+1}^l(t) = 1, \\ \lambda_{n+2}^l - \frac{\lambda_{n+1}^l}{2} \leq x < \frac{\lambda_1^l}{2} : & \dots \\ \dots & \dots \\ \lambda_{2n}^l - \frac{\lambda_{n+1}^l}{2} \leq x < \lambda_{2n-1}^l - \frac{\lambda_{n+1}^l}{2} : & s_{j,2n-1}^l(t) = 1, \\ x < \lambda_{2n}^l - \frac{\lambda_{n+1}^l}{2} : & s_{j,2n}^l(t) = 1. \end{array} \right]$$



正阈值部分

The results of experiments presented in Section 5.4 indicate that although this neuron has multiple thresholds, most of the spikes it generated are concentrated in θ_1 and $-\theta_2$. The spikes generated by other thresholds are minimal, which reduces energy consumption and inference latency.

负阈值部分

Methods

Sec.2 Expectation Compensation Module

Parallel Parameter normalization for MT Neuron

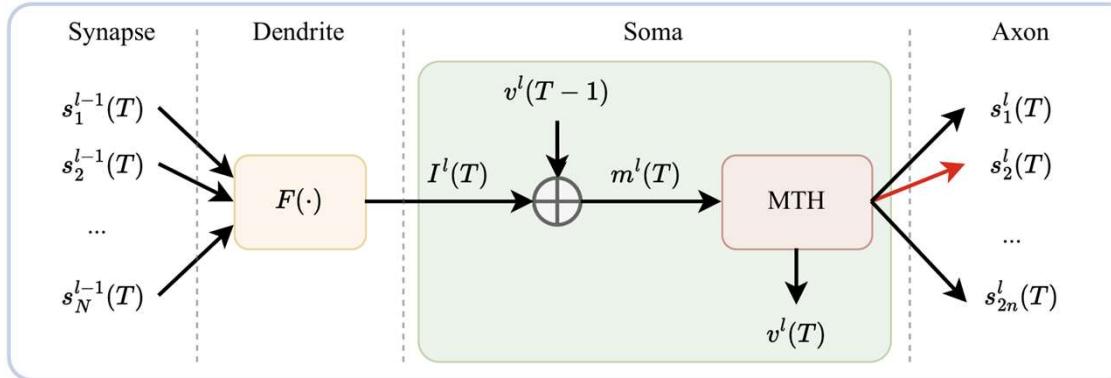
extends the ANN weight to $2n$ weights in the SNN corresponding to $2n$ thresholds of MT neurons

$$W_{\text{SNN},p}^l = W_{\text{ANN}}^l \frac{\lambda_p^{l-1}}{\lambda_1^l}$$

$$I_j^l(t) = \sum_{i,p} w_{ij\text{SNN},p}^l s_{i,p}^{l-1}(t)$$

$$\theta_{1,new} = 1, \theta_{2,new} = \eta$$

```
def main(args):
    args.distributed = False
    print("{}".format(args.replace(',', ',\n')))
    device = torch.device(args.device)
    cudnn.benchmark = True
    dataset_val, args.nb_classes = build_dataset(is_train=False, args=args)
    num_tasks = utils.get_world_size()
    global_rank = utils.get_rank()
    sampler_val = torch.utils.data.DistributedSampler(dataset_val, num_replicas=num_tasks, rank=global_rank, shuffle=False)
    data_loader_val = torch.utils.data.DataLoader(dataset_val, sampler=sampler_val, batch_size=int(args.batch_size), num_workers=args.num_workers, pin_memory=True, drop_last=False)
    model = create_model(args.model, pretrained=False, img_size=args.input_size, num_classes=args.nb_classes)
```



Algorithm 1 The conversion method using Expectation Compensation Module and Multi-Threshold Neuron(ECMT)

Input: Pre-trained Transformer ANN model $f_{\text{ANN}}(W)$; Dataset D ; Time-step T to test dataset; Threshold percent p .
Output: SNN model $f_{\text{SNN}}(W, \theta_1, \theta_2, v)$

- 1: **step1:** Obtain the base thresholds θ_1 and θ_2
- 2: **for** length of Dataset D **do**
- 3: Sample minibatch data from D
- 4: Run the data on f_{ANN} and static the activation values before linear and matrix product module at $p\%$ and $(1-p\%)$, setting them as θ_1 and $-\theta_2$ respectively.
- 5: **end for**
- 6: **step2:** Converted to SNN model
- 7: **for** module m in $f_{\text{ANN}}.\text{Module}$ **do**
- 8: **if** m is Linear Module **then**
- 9: Add a Multi-Threshold Neuron before m
- 10: **else if** m is Matrix Product **then**
- 11: replace m by two Multi-Threshold Neurons followed by a Matrix Product EC Module
- 12: **else if** m is Other Nonlinear Module **then**
- 13: replace m by an EC Module
- 14: **end if**
- 15: **end for**
- 16: Set the base thresholds of MT neurons to corresponding $\theta_1, -\theta_2$ and set the initial membrane potential v to 0.
- 17: $f_{\text{SNN}} = \text{Parallel Parameter normalization}(f_{\text{ANN}})$
- 18: **return** f_{SNN}

Experiments

