Efficient Spiking Neural Networks with Radix Encoding

Zhehui Wang¹, **Xiaozhe Gu**², **Rick Goh**¹, **Joey Tianyi Zhou**¹ and **Tao Luo**¹ Institute of High Performance Computing, A*STAR, Singapore ²Chinese University of Hong Kong, Shenzhen, China {zhehui.w, leto.luo}@gmail.com,

Abstract

Spiking neural networks (SNNs) have advantages in latency and energy efficiency over traditional artificial neural networks (ANNs) due to its eventdriven computation mechanism and replacement of energy-consuming weight multiplications with additions. However, in order to reach accuracy of its ANN counterpart, it usually requires long spike trains to ensure the accuracy. Traditionally, a spike train needs around one thousand time steps to approach similar accuracy as its ANN counterpart. This offsets the computation efficiency brought by SNNs because longer spike trains mean a larger number of operations and longer latency. In this paper, we propose a radix encoded SNN with ultrashort spike trains. In the new model, the spike train takes less than ten time steps. Experiments show that our method demonstrates 25X speedup and 1.1% increment on accuracy, compared with the state-of-the-art work on VGG-16 network architecture and CIFAR-10 dataset.

1 Introduction

Spiking neural networks (SNNs) increasingly attract people's attention in recent years because of their similarity with the biological neural system. Theoretically, SNN shows higher energy efficiency and lower latency than traditional artificial neural networks (ANNs) due to its special computation mechanism. The inference of SNN consists of only simple addition, which utilizes much less hardware resource than the complicated, energy-consuming multiplication operations used by ANN. Today, the state-of-the-art works have pushed the accuracy of SNN to a competitive level with ANN. For example, Rueckauer's work shows over 90% accuracy under the CIFAR-10 dataset [Rueckauer et al., 2017]. In the past decade, SNN technology is fast developing. Researchers not only advanced the model accuracy but also fabricated many types of neuromorphic computing chips to further improve the computation efficiency of SNN, such as SpiN-Naker [Khan et al., 2008], TrueNorth [Merolla et al., 2014], and Loihi [Davies et al., 2018]

However, in order to ensure high accuracy performance, traditional rate encoded SNN models require long spiking trains. A spike train can be considered as a sequence of time steps. The length of the spike train is defined as the number of time steps. At each time step, the neuron

may choose to fire a spike or not. Traditional SNN is rate encoded [Diehl *et al.*, 2015] [Sengupta *et al.*, 2019], which means that the signal strength in SNN is represented by the number of spikes over the number of time steps (i.e., the length of the spike train). Therefore, the precision of the signal strength is proportional to the length of the spike train. For example, ten time steps correspond to 10% error, and one hundred time steps correspond to 1% error. To pursue high model accuracy, the signal strength expressed in SNN should be as precise as possible. Hence, a long spike train becomes necessary for rate encoded SNN. Experiments show that to achieve competitive accuracy with the ANN, we need around one thousand time steps to encode the signal strength [Sengupta *et al.*, 2019].

Unfortunately, long spike trains offset the advantages of SNN over ANN on computation efficiency. At each time step, neurons need to accumulate all the spikes from the input source and decide whether it should fire or not. More time steps mean more addition operations during the simulation. Today, many works tried to shorten the spike train [Wu et al., 2019][Zhang et al., 2019]. However, it is difficult to do so without obvious accuracy loss. To solve this issue, we proposed radix encoding method, where spikes at different time steps have different weights. In the new design, some spikes are more important than others when representing the signal strength. This is different from the traditional rate encoding method, where all the spikes are identical and have the same weight. Therefore, with only a few time steps, we can represent high-precision signal strength.

We apply our radix encoded method into the ANN-to-SNN process [Pérez-Carrasco *et al.*, 2013] and train several SNN models using popular network architectures and datasets. We then compare our model with the state-of-the-art works in terms of accuracy, the number of operations, and latency. Large improvements are observed in experiments. Overall, the contributions of this paper can be concluded in the following three points.

- We propose a radix encoded SNN with ultra-short spike trains. Compared with the state-of-the-art works, our new design shows at least 25X speedup and 1.1% increment on accuracy.
- We develop a technique to simplify the training of SNN by transforming the non-differentiable spikes signals into trainable data. We mathematically prove the advantage of our method for SNN training.
- We speed up the training of SNN using mature platforms and accelerators. Experiments show that we can train

the SNN model within one hour and achieve even higher accuracy than the state-of-the-art works.

The rest of the paper is organized as follows. Section 2 introduces the details of our radix encoded SNN. Section 3 discusses the experiments. Section 4 gives the related work. Section 5 concludes.

2 SNN with Radix Encoding

Different from traditional ANN, where the input and output of the neurons are real numbers. In SNN, the input/output is a group of spikes, called spike trains. The spike trains are normally equal in length. Mathematically, we can express a spike train in Equation (1). Each element in the spike train stands for the existence of the spike. By definition, s(t) equals 1 if the spike exists at time step t, and equals 0 if not. T is the length of the spike train.

$$s(t) \in \{0, 1\}$$
 $t = 0, 1, 2, ..., T - 1$ (1)

2.1 Leaky Integrate-and-Fire (LIF) Model

Biologically, a neuron in SNN receives signals from upstream neurons via many synapses. Each synapse can receive one spike train sequentially and independently. We denote the input sequence from n-th synapse as $i_n(t)$, and express it in Equation (2). Here $s_n^{in}(t)$ stands for the input spike train from n-th synapse. If the time step t is greater than the spike train length T, $i_n(t)$ always equals 0.

$$i_n(t) = \begin{cases} s_n^{in}(t) & \text{if } 0 \le t < T \\ 0 & \text{if } t \ge T \end{cases}$$
 (2)

The operation mechanism of neuron in the radix encoded SNN is shown in Algorithm 1. Each neuron has an internal parameter called the membrane potential. At the beginning, this value is initialized to be b (line 1). Next, the value of the membrane potential is updated every time step. The neuron behavior at time step t is based on the value of membrane potential at the previous time step t-1. In the leaky integrateand-fire (LIF) model [Hunsberger and Eliasmith, 2015], we can divide a single time step into three stages. The first stage is called integration. The neuron receives spikes from all the synapses and accumulates these spikes with different weights into the neuron's membrane potential v (line 3-4). The second stage is called firing. If the membrane potential meets a certain condition, the neuron fires a spike via the axon (line 5-6). If the spike is fired, the membrane potential will go through the refractory period (line 7). Here \hat{v} denotes the least significant digit of v in binary format, V_{th} is the threshold, and V_r is the drop value of the membrane potential. The third stage is called leakage. The membrane value would decrease by κ (line 10). Here κ is the decay factor.

We denote the membrane potential of m-th neuron after the integration stage at time step t by u(t). It can be expressed in Equation (3). Here $v_m(t-1)$ is the final membrane potential of m-th neuron after three stages at time step t-1. w_{nm} is the weight associated with n-th synapse and m-th neuron. N is the total number of synapses. Noted that w_{nm} is an integer.

$$u_m(t) = v_m(t-1) + \sum_{n=0}^{N-1} i_n(t) \cdot w_{nm}$$
 (3)

Algorithm 1: Leaky integrate-and-fire (LIF) neuron model for radix encoded SNN

Input: binary sequence $i_n[t]$, weight w_n , bias bOutput: binary sequence o[t]1 v=b; // Initialization
2 for $t\leftarrow 0$ to T'-1 do
3 for $n\leftarrow 0$ to N-1 do
4 $v\leftarrow v+i_n[t]*w_n$; // Integration
5 if $\hat{v}>V_{th}$ then
6 o[t]=1; // Firing
7 $v=v-V_r$; // Refractory Period
8 else
9 $v\leftarrow v*\kappa$; // Leakage

After integration, the neuron will go through the firing stage and the leakage stage. We denote the membrane potential of m-th neuron after those two stages at time step t by $v_m(t)$. It can be expressed in Equation (4). Here $o_m(t)$ is the output sequence, which can be either 0 or 1.

$$v_m(t) = \kappa \cdot (u_m(t) - V_r \cdot o_m(t)) \tag{4}$$

At the very beginning, the membrane potential is initialized by a bias value b_m , expressed in Equation (5). Note that b_m is an integer. Here the time step -1 is used to index the states of the neuron before the first time step 0.

$$v_m(-1) = b_m \tag{5}$$

2.2 Radix Encoding Method

We propose the radix encoding method for SNN, where the hyper-parameters $V_{th}=0$, $V_r=1$ and $\kappa=0.5$. Based on Equation (3-5), we can derive $v_m(t)$ as Equation (6).

$$v_m(t) = \sum_{n=0}^{N-1} \sum_{\tau=0}^{t} 2^{-\tau-1} i_n(t-\tau) \cdot w_{nm} + 2^{-t-1} b_m$$

$$-\sum_{\tau=0}^{t} 2^{-\tau-1} o_m(t-\tau)$$
(6)

Lemma 1. There exists an integer T' > T, such that when $t \ge T' - 1$, then $v_m(t) = c$, where c = 0 or c = -1.

Proof. Since T' > T, we have t > T - 1, the relationship between the membrane potential at time step t + 1 and time step t can be expressed in Equation (7).

$$v_m(t+1) = 2^{-1}(v_m(t) - o_m(t)) \tag{7}$$

We denote $\Delta v_m(t)$ as the difference between $v_m(t+1)$ and $v_m(t)$, and express it in Equation (8)

$$\Delta v_m(t) = v_m(t+1) - v_m(t) = -2^{-1}(v_m(t) + o_m(t))$$
 (8)

Since $i_n(t)$ is a binary sequence, w_{nm} and b_m are integers, $v_m(T-1)$ must be an integer. Its value can be either greater than or equals to 0, or it can be less than or equals to -1. We discuss these two possibilities one by one. If $v_m(T-1)$ is

greater than or equals to 0, we can infer from Equation (7) that $v_m(t)$ is always greater than or equals to 0, for any t greater than T-1. On the other hand, we can infer from Equation (8) that $\Delta v_m(t)$ is less than or equals to -1, if $v_m(t)$ is greater than 0. We list these two conclusions in Equation (9) and (10). Combine them together, it can be proved that series $v_m(t)$ will eventually converge to 0.

$$v_m(t) \ge 0 \qquad \qquad \text{if } v_m(T-1) \ge 0 \tag{9}$$

$$\Delta v_m(t) \le -1 \qquad \text{if } v_m(t) > 0 \tag{10}$$

If $v_m(T-1)$ is less than or equals to -1, based on the similar logic, we can have two conclusions, listed in Equation (11) and (12). Combine them together, it can be proved that series $v_m(t)$ will eventually converge to -1.

$$v_m(t) \le -1$$
 if $v_m(T-1) \le -1$ (11)

$$v_m(t) \le -1$$
 if $v_m(T-1) \le -1$ (11)
 $\Delta v_m(t) \ge 1$ if $v_m(t) < -1$ (12)

Noted that T' is greater than T, for alignment reasons, the output spikes from the m-th neuron $s_m^{out}(t)$ can be expressed in Equation (13). Here the time step t ranges from T' - T to T'-1. This is to guarantee that the input spike $s_n^{in}(t)$ and the output spike $s_m^{out}(t)$ have the same length.

$$s_m^{out}(t) = o(t + T' - T)$$
 (13)

From Lemma 1, we show that the sequence of $v_m(t)$ would finally converge to either 0 or -1, when t = T' - 1. We assume that the input sequence $i_n(t)$, weight matrix w_{nm} and bias b_m meet Condition (14).

$$\sum_{n=0}^{N-1} \sum_{\tau=0}^{T'-1} 2^{-\tau-1} i_n (T'-\tau-1) \cdot w_{nm} + 2^{-T'} b_m \ge 0$$
 (14)

From Equation (6) we can infer that $v_m(T'-1) > -1$ under Condition (14). This is because of Inequality (15)

$$\sum_{\tau=0}^{T'-1} 2^{-\tau-1} o_m(T'-\tau-1) < 1 \tag{15}$$

Since $v_m(T'-1)$ can be either 0 or -1, we have $v_m(T'-1) =$ 0, under Condition (14). Put it back to Equation (6), the relationship between input sequence $i_n(t)$ and output sequence $o_m(t)$ can be expressed in Equation (16).

$$\sum_{\tau=0}^{T'-1} 2^{-\tau-1} o_m(T'-\tau-1) = \sum_{n=0}^{N-1} \sum_{\tau=0}^{T'-1} 2^{-\tau-1} i_n(T'-\tau-1) \cdot w_{nm} + 2^{-T'} b_m$$
(16)

By scaling Equation (16) on both sides with $2^{T'}$, and substituting $T' - \tau - 1$ with t, we have Equation (17).

$$\sum_{t=0}^{T'-1} 2^t o_m(t) = \sum_{n=0}^{N-1} \sum_{t=0}^{T'-1} 2^t i_n(t) \cdot w_{nm} + b_m$$
 (17)

Table 1: Data and Model Transformation Functions

Input x, X	Data Transform & Inverse Transform
$s_n^{in}(t), s_m^{out}(t)$	$WDT(x) = \sum_{t=0}^{T-1} 2^t x(t)$
S_n^{in}, S_m^{out}	$\mathrm{WDT}^{-1}(X) = \lfloor X/2^t \rfloor \bmod 2$
Input x, X	Parameter Transform & Inverse Transform
w_{nm}, b_m	$WPT(x) = 2^{-\Delta T}x$
W_{nm}, B_m	$\mathrm{WPT}^{-1}(X) = 2^{\Delta T} X$

By putting Equation (2) and (13) back to Equation (17), substituting T'-T by ΔT , and scaling Equation (17) on both sides with $2^{-\Delta T}$, we have Equation (18).

$$\sum_{t=0}^{T-1} 2^t s_m^{out}(t) + \sum_{t=0}^{\Delta T-1} 2^{t-\Delta T} o_m(t) = \sum_{n=0}^{N-1} \sum_{t=0}^{T-1} 2^t s_n^{in}(t) \cdot 2^{-\Delta T} w_{nm} + 2^{-\Delta T} b_m$$
(18)

2.3 **Data and Parameter Transformation**

To make the Equation (18) easy to read, we define some new parameters equalling those complicated terms, shown in Equations (19). Noted that $s_n^{in}(t)$ and $s_m^{out}(t)$ are binary sequences. Each item can be only 0 or 1. The capitalized S_n^{in} or S_n^{in} is a single real number. The weight w_{nm} and bias b_m can only be integers, and the capitalized weight W_{nm} and bias B_m can be any real number.

$$S_n^{in} = \sum_{t=0}^{T-1} 2^t s_n^{in}(t) \qquad S_m^{out} = \sum_{t=0}^{T-1} 2^t s_m^{out}(t)$$

$$W_{nm} = 2^{-\Delta T} w_{nm} \qquad B_m = 2^{-\Delta T} b_m$$
(19)

By putting new denotations listed in Equations (19) back to Equation (18), we obtain Equation (20).

$$S_m^{out} + r_m = \sum_{n=0}^{N-1} S_n^{in} \cdot W_{nm} + B_m$$
 (20)

Here the residue r_m is related to the first few elements in the sequence $o_m(t)$, which can be expressed in Equation (21). It is easy to prove that r_m is less than 1.

$$r_m = \sum_{t=0}^{\Delta T - 1} 2^{t - \Delta T} o_m(t) < 1$$
 (21)

Since S_m^{out} is on the scale of 2^T , which is much larger than r_m , we do some approximations, and rewrite Equation (20) as approximate Equation (22).

$$S_m^{out} \approx \sum_{n=0}^{N-1} S_n^{in} \cdot W_{nm} + B_m \tag{22}$$

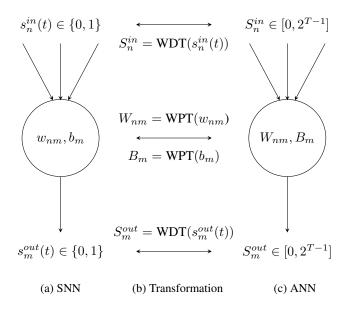


Figure 1: Data transformation and parameter transformation between the SNN neuron and the ANN neuron

From Equation (16), our derivations are all based on Condition (14). By putting new denotations listed in Equations (19) back to Condition (14), we obtain Condition (23).

$$\sum_{n=0}^{N-1} S_n^{in} \cdot W_{nm} + B_m \ge 0 \tag{23}$$

In ANN, the relationship between the input and the output can also be expressed as Equation (22), if we assume S_n^{in} and S_m^{out} to be input and output data, W_{nm} and B_m to be weight and bias. This similarity tells us that the relationship between the input and the output in the SNN domain can be expressed in the ANN domain by transforming the data and the parameters. Table 1 shows the transformation methods. WDT and WDT $^{-1}$ are the transformation and inverse transformation function for input and output data. WPT and WPT $^{-1}$ are the transformation and inverse transformation function for weights and bias parameters. Noted that the input of WDT and the output of WDT $^{-1}$ are time series with index t. Fig. 1 clearly shows the relationship of the neurons in the SNN domain and the neurons in the ANN domain. We can describe their relationship in the following two aspects.

Input/Output: The input/output data in SNN are time series whose elements are either 0 or 1, while the input/output data in ANN are real number ranging from 0 to 2^{T-1} . We can transform them by WDT and WDT $^{-1}$. In traditional ANN with batch normalization layer, the elements in the input and output feature maps are normalized, which range from 0 to 1. Since T is a fixed value, we can scale the data back by 2^T .

Weight/Bias: The weight/bias in SNN are integers while the weight/bias in ANN are real numbers. We can transform them by WPT and WPT⁻¹. In traditional ANN with the batch normalization layer, the normalization process is directly applied to the output feature map. Alternatively, we can apply the normalization on trained weights and biases, to achieve similar effects [Rueckauer *et al.*, 2016].

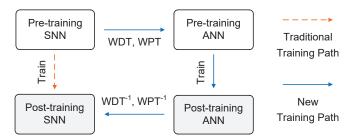


Figure 2: The traditional training path and the proposed training path for SNN models using transformation functions

2.4 Efficient Training Path for SNN

Training SNN directly is more difficult than training CNN because of the lack of training platforms and hardware accelerators. To solve this problem, we propose a new training path for SNN training. Fig. 2 shows the difference between these two training paths. In this figure, the dashed line denotes the traditional training path for SNN. Following this path, the pre-training SNN model is trained to the posttraining SNN models directly. Alternatively, there is a new training path for SNN, shown as the solid line. Following this new path, we first transform the data and parameters from the pre-training SNN model to the pre-training ANN model using functions WDT and WPT. Then we train the pre-training ANN model to the post-training ANN model using mature platforms and hardware accelerators. Afterward, we transform the post-training ANN model back to the post-training SNN model using functions WDT $^{-1}$ and WPT $^{-1}$. Compared with the traditional training path, this new path is more efficient in time and energy consumption.

3 Experiment

We implement our radix encoding method in CUDAaccelerated PyTorch. For comparison reasons, we test several popular ANN architectures, including ResNet-18, VGG-16, and MobileNet. For the ANN model, we train the network with full-precision weights and activations for several epochs until the accuracy converges. For the SNN model, we implement a customized quantization function and normalization function to simulate the proposed training path shown in Fig 2. The weights are quantized to eight bits while the activations are quantized by T bits. Here T is the number of time steps. We train the SNN model from a well-trained checkpoint and fine-tune the model by several epochs using the STE [Yin et al., 2019] method until the accuracy converges. The neuron networks are trained by the SGD optimizer with a dynamic learning rate, which decreases by half for every thirty epochs. We do the experiments on a single graphic card Nvidia's 2080Ti, and each experiment can finish within one hour.

We compare the number of operations between ANN and SNN during inference. For SNN, the number of operations is calculated based on Equation 24, where N and M are the number of input spike trains and output spike trains in a neuron array. T is the number of time steps. To obtain the total number of operations, we accumulated the number of oper-

Table 2: Comparison between SNN and ANN models on the CIFAR-10/100 dataset

			CIFAR-10				CIFAR-100				
Model	Type	Time Steps	Accuracy	Δ Accuracy	#Operations	Latency	Accuracy	Δ Accuracy	#Operations	Latency	
ResNet-18	ANN	N/A	95.2%	N/A	36.7×10^9	1.77	77.4%	N/A	36.7×10^9	1.77	
	SNN	8	95.2%	0.0%	4.44×10^9	0.21	77.6%	0.2%	4.45×10^9	0.21	
	SNN	4	92.1%	-3.1%	2.22×10^{9}	0.11	74.2%	-3.2%	$2.22{\times}10^9$	0.11	
VGG-16	ANN	N/A	93.8%	N/A	20.7×10^9	1.00	74.1%	N/A	20.7×10^9	1.00	
	SNN	8	93.8%	0.0%	2.51×10^{9}	0.12	73.9%	-0.2%	2.51×10^{9}	0.12	
	SNN	4	92.2%	-1.6%	1.25×10^{9}	0.06	70.1%	-4.0%	1.25×10^{9}	0.06	
MobileNet	ANN	N/A	92.0%	N/A	3.07×10^9	0.15	65.8%	N/A	3.07×10^9	0.15	
	SNN	8	91.6%	-0.4%	0.37×10^{9}	0.02	64.7%	-1.1%	0.37×10^{9}	0.02	
	SNN	4	90.5%	-1.5%	0.19×10^{9}	0.01	62.7%	-3.1%	0.19×10^{9}	0.01	

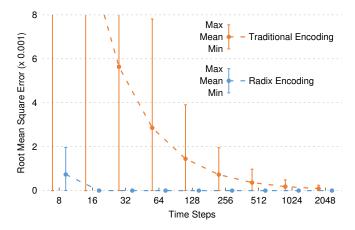


Figure 3: Root mean square error between the traditional encoding method and the radix encoding method, data collected from VGG-16 on CIFAR-10 dataset, and normalized to the range [0, 1]

ations from all the neuron arrays. Since ANN uses a different computation mechanism from SNN, we cannot compare their operation counts directly. To make the comparison, we find the number of effective operations in ANN, assuming that each effective operation in ANN uses the same amount of hardware resource as to each operation in SNN.

$$\Sigma N \cdot M \cdot T$$
 (24)

We also compare the inference latencies between ANN and SNN, which includes the time to execute all the operations plus the extras time for data processing. We assume both SNN and ANN use the same amount of hardware resources with an equal number of processing elements (PEs). For convenience, we normalized all the latencies in this paper. Specifically, the latency of the VGG-16 ANN network is normalized to be one.

3.1 Root Mean Square Error

The derivation of the transformation method WDT and WPT is based on the condition that the residual (expressed in Equation 21) is relatively small, and can be ignored. To check the real value of this residual, we evaluate the VGG-16 network

on the CIFAR-10 dataset, and show the root means square of the residual under different time steps in Fig. 3. In statistics, the root means square of the residual is know as the root mean square error (RMSE) [Chai and Draxler, 2014]. For comparison reasons, we normalize the RMSE by 2^T , where T is the number of time steps. We study the RMSE in two cases. In the first case, we use our radix encoding method. In the second case, we use the traditional rate encoding method, where each activation is represented by the number of spikes over the spike train length. On one hand, we want the RMSE value to be as small as possible to achieve the high accuracy of the model. On the other hand, we prefer fewer time steps because it results in fewer operations and smaller latency.

In Fig. 3, we show the RMSE as well as the maximum and minimum residuals. The dashed lines show the trend of the RMSE with the increment of time steps. From the figure, we can see that our radix encoding method shows much less RMSE than that of the traditional encoding method. For example, if the number of time steps is eight, the RMSE of the radix encoding method is 96% less than that of the traditional encoding method (this data is not shown in the figure because it is out of the figure boundary). From another perspective, to achieve the same RMSE, the required time steps of the radix encoding method is much less than that of the traditional encoding method. For example, the radix encoding method with 16 time steps has a similar RMSE as the traditional encoding method with 256 time steps.

3.2 Validation of radix Encoding

To validate our radix encoding method, we test three networks on the CIFAR-10 and the CIFAR-100 datasets. In Table. 2, we compare the inference performance of ANN and SNN in terms of accuracy, the number of operations, and latency. From the table, we can see that under similar accuracies, SNN shows fewer operations than ANN. For example, on average of the three networks and two datasets, the SNN with eight time steps show 88% less latency than the ANN using the same architecture, with almost the same accuracy. This is due to the advantage of SNN over ANN on computation efficiency. In SNN, the complicated weight multiplications are replaced by simple additions, so that its utilization of hardware resources can be substantially reduced.

Table 3: Comparison between this work and the state-of-the-art on the CIFAR-10/100 dataset

Model	Architecture	Dataset	Time Steps	Accuracy	Δ Accuracy	#Operations	Latency	Speedup
GD-SNN [Sengupta et al., 2019]	VGG-16	CIFAR-10	2500	91.6%	-2.2%	783×10 ⁹	37.8	0.04X
Hybrid [Rathi et al., 2020]	VGG-16	CIFAR-10	100	91.1%	-2.7%	31.3×10^{9}	1.51	1X
Radix Encoding (this work)	VGG-16	CIFAR-10	8	93.8%	0.0%	2.51×10^9	0.12	12.5X
Radix Encoding (this work)	VGG-16	CIFAR-10	4	92.2%	-1.6%	1.25×10^9	0.06	25X
SBP-SNN [Lee et al., 2020]	VGG-9	CIFAR-10	100	90.5%	-2.7%	19.1×10^9	0.92	1X
Radix Encoding (this work)	VGG-9	CIFAR-10	4	91.7%	-1.5%	0.76×10^9	0.04	25X
S-ResNet [Hu et al., 2018]	ResNet-44	CIFAR-100	350	68.6%	-1.6%	34.1×10^9	1.65	1X
Radix Encoding (this work)	ResNet-44	CIFAR-100	8	69.7%	-0.5%	0.78×10^9	0.04	43.8X

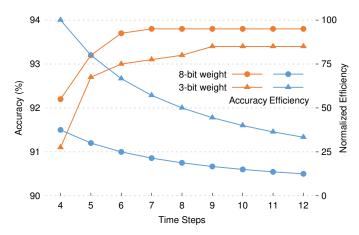


Figure 4: Accuracy and efficiency of radix encoded SNN versus the number of time steps. Efficiency is inversely proportional to latency. Data collected from VGG-16 with CIFAR-10 dataset

We also compare the accuracy and efficiency of radix encoded SNN versus times steps in Fig. 4. Data are collected from the VGG-16 network on the CIFAR-10 datasets. We assume that the weights are quantized by 8 bits and 3 bits. Here the efficiency is defined as the number of processed images during inference in a fixed time. From the figure, we can see that the SNN model with 8-bit weights has higher accuracy than that with 3-bit weights, but its efficiency is 63% lower. Another observation from the figure is that the increment of time steps does not guarantee the increment of accuracy. If the number of time steps is greater than a certain point, the model accuracy reaches its maximum value. This certain point can be regarded as the optimal time step for this radix encoded SNN. In this example, the optimal time steps are eight and nine, for the 8-bit weight model and 3-bit weight model, respectively.

3.3 Comparison with the State-of-the-Art

We compare our radix encoded SNN with the state-of-theart work in Table 3. The comparison is conducted under the same network architecture and the same dataset. From the table, we can see that our work can achieve not only higher accuracy but also significantly higher efficiency. For example, compared with the hybrid method [Rathi *et al.*, 2020], the radix encoded SNN with eight time steps shows at least 12.5X speedup and 2.7% increment in accuracy. If the number of time steps is reduced to four, we can observe 25X speedup with around 1.1% increment on accuracy. This is mainly due to the reason that the radix encoding method used in our work shows a much lower residual than the traditional method. With only a few time steps, the radix encoded SNN can achieve very high precision. As we discussed, fewer time steps means less number of operations and lower latency. Hence, the efficiency of the SNN model can be substantially improved by the radix encoding method.

4 Related Work

The training of spiking neural networks can be classified into unsupervised and supervised training. One popular unsupervised SNN is spike-timing-dependent plasticity (STDP) [Hebb, 1949]. For supervised learning, although it demonstrates competitive accuracy than unsupervised learning, the training process is difficult because of the nondifferentiability of spike trains. There are mainly two solutions to alleviate this problem. One solution is to approximate the spiking function as a differentiable function [Huh and Sejnowski, 2018] or to define a surrogate gradient function for the backpropagation of spiking function [Neftci et al., 2019]. This type of solution has a high computation cost because of the complicated approximating functions implemented during training. Today, many of these works are limited to shallow networks on simple dataset [Shrestha and Orchard, 2018] [Kaiser et al., 2020], and [Bellec et al., 2018].

A more widely used solution today is to first train ANN using the standard training method and then converted the well-trained ANN to event-based SNN [Rueckauer et al., 2017]. This type of solution does not suffer from the non-differentiability of spike trains, and the high computation cost of the approximation functions. There are several works in the literature showing relatively high accuracies. Sengupta et al. proposed a weight normalization method for the SNN-to-ANN conversion [Sengupta et al., 2019]. Rathi et al. proposed a hybrid conversion method for SNN [Rathi et al., 2020]. Lee et al. proposed a spike-based backpropagation method for SNN training [Lee et al., 2020]. Hu et al. proposed a method to convert the ANN with a shortcut to the SNN model [Hu et al., 2018]. These work use long spike trains. Different from them, our work uses the radix

encoding method, resulting in ultra-short spike trains.

5 Conclusion

In this paper, we proposed an efficient radix encoded SNN with ultra-short spike trains. To simplify the training of SNN, we develop a transform function to convert non-differentiable spikes into trainable data. We also develop a mechanism to speed up the training process by using mature platforms and hardware accelerators. Experiments show that compared with the state-of-the-art, our work achieves 25X speedup, with 1.1% increment on accuracy.

References

- [Bellec et al., 2018] Guillaume Bellec, Darjan Salaj, Anand Subramoney, Robert Legenstein, and Wolfgang Maass. Long Short-term Memory and Learning-to-learn in Networks of Spiking Neurons. Advances in Neural Information Processing Systems, 31:787–797, 2018.
- [Chai and Draxler, 2014] Tianfeng Chai and Roland R Draxler. Root Mean Square Error (RMSE) or Mean Absolute Error (MAE)? *GMDD*, 7(1):1525–1534, 2014.
- [Davies *et al.*, 2018] Mike Davies, Narayan Srinivasa, Tsung-Han Lin, Gautham Chinya, Yongqiang Cao, Sri Harsha Choday, Georgios Dimou, Prasad Joshi, Nabil Imam, Shweta Jain, et al. Loihi: A Neuromorphic Manycore Processor with On-chip Learning. *IEEE Micro*, 38(1):82–99, 2018.
- [Diehl et al., 2015] Peter U Diehl, Daniel Neil, Jonathan Binas, Matthew Cook, Shih-Chii Liu, and Michael Pfeiffer. Fast-classifying, High-accuracy Spiking Deep Networks through Weight and Threshold Balancing. In 2015 International Joint Conference on Neural Networks (IJCNN), pages 1–8. ieee, 2015.
- [Hebb, 1949] Donald Olding Hebb. *The Organization of Behavior: a Neuropsychological Theory*. J. Wiley; Chapman & Hall, 1949.
- [Hu et al., 2018] Yangfan Hu, Huajin Tang, Yueming Wang, and Gang Pan. Spiking Deep Residual Network. arXiv preprint arXiv:1805.01352, 2018.
- [Huh and Sejnowski, 2018] Dongsung Huh and Terrence J Sejnowski. Gradient Descent for Spiking Neural Networks. In *Advances in Neural Information Processing Systems*, pages 1433–1443, 2018.
- [Hunsberger and Eliasmith, 2015] Eric Hunsberger and Chris Eliasmith. Spiking Deep Networks with LIF Neurons. *arXiv preprint arXiv:1510.08829*, 2015.
- [Kaiser et al., 2020] Jacques Kaiser, Hesham Mostafa, and Emre Neftci. Synaptic Plasticity Dynamics for Deep Continuous Local Learning (DECOLLE). Frontiers in Neuroscience, 14:424, 2020.
- [Khan et al., 2008] Muhammad Mukaram Khan, David R Lester, Luis A Plana, A Rast, Xin Jin, Eustace Painkras, and Stephen B Furber. SpiNNaker: Mapping Neural Networks onto a Massively-parallel Chip Multiprocessor. In

- 2008 IEEE International Joint Conference on Neural Networks (IEEE World Congress on Computational Intelligence), pages 2849–2856. Ieee, 2008.
- [Lee et al., 2020] Chankyu Lee, Syed Shakib Sarwar, Priyadarshini Panda, Gopalakrishnan Srinivasan, and Kaushik Roy. Enabling Spike-based Backpropagation for Training Deep Neural Network Architectures. Frontiers in Neuroscience, 14, 2020.
- [Merolla et al., 2014] Paul A Merolla, John V Arthur, Rodrigo Alvarez-Icaza, Andrew S Cassidy, Jun Sawada, Filipp Akopyan, Bryan L Jackson, Nabil Imam, Chen Guo, Yutaka Nakamura, et al. A Million Spiking-neuron Integrated Circuit with a Scalable Communication Network and Interface. Science, 345(6197):668–673, 2014.
- [Neftci et al., 2019] Emre O Neftci, Hesham Mostafa, and Friedemann Zenke. Surrogate Gradient Learning in Spiking Neural Networks. IEEE Signal Processing Magazine, 36:61–63, 2019.
- [Pérez-Carrasco et al., 2013] José Antonio Pérez-Carrasco, Bo Zhao, Carmen Serrano, Begona Acha, Teresa Serrano-Gotarredona, Shouchun Chen, and Bernabé Linares-Barranco. Mapping from Frame-driven to Frame-free Event-driven Vision Systems by Low-rate Rate Coding and Coincidence Processing—Application to Feedforward ConvNets. IEEE transactions on pattern analysis and machine intelligence, 35(11):2706–2719, 2013.
- [Rathi et al., 2020] Nitin Rathi, Gopalakrishnan Srinivasan, Priyadarshini Panda, and Kaushik Roy. Enabling Deep Spiking Neural Networks with Hybrid Conversion and Spike Timing Dependent Backpropagation. arXiv preprint arXiv:2005.01807, 2020.
- [Rueckauer et al., 2016] Bodo Rueckauer, Iulia-Alexandra Lungu, Yuhuang Hu, and Michael Pfeiffer. Theory and Tools for the Conversion of Analog to Spiking Convolutional Neural Networks. arXiv preprint arXiv:1612.04052, 2016.
- [Rueckauer et al., 2017] Bodo Rueckauer, Iulia-Alexandra Lungu, Yuhuang Hu, Michael Pfeiffer, and Shih-Chii Liu. Conversion of Continuous-valued Deep Networks to Efficient Event-driven Networks for Image Classification. Frontiers in neuroscience, 11:682, 2017.
- [Sengupta *et al.*, 2019] Abhronil Sengupta, Yuting Ye, Robert Wang, Chiao Liu, and Kaushik Roy. Going Deeper in Spiking Neural Networks: VGG and Residual Architectures. *Frontiers in neuroscience*, 13:95, 2019.
- [Shrestha and Orchard, 2018] Sumit B Shrestha and Garrick Orchard. Slayer: Spike Layer Error Reassignment in Time. *Advances in neural information processing systems*, 31:1412–1421, 2018.
- [Wu et al., 2019] Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, Yuan Xie, and Luping Shi. Direct Training for Spiking Neural Networks: Faster, Larger, Better. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 1311–1318, 2019.

[Yin et al., 2019] Penghang Yin, Jiancheng Lyu, Shuai Zhang, Stanley Osher, Yingyong Qi, and Jack Xin. Understanding Straight-through Estimator in Training Activation Quantized Neural Nets. arXiv preprint arXiv:1903.05662, 2019.

[Zhang et al., 2019] Lei Zhang, Shengyuan Zhou, Tian Zhi,

Zidong Du, and Yunji Chen. TDSNN: From Deep Neural Networks to Deep Spike Neural Networks with Temporal-coding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 1319–1326, 2019.