

# Supplementary Materials for: Towards Memory- and Time-Efficient Backpropagation for Training Spiking Neural Networks

Anonymous ICCV submission

Paper ID 10305

## A. Derivation for Eqs. (9) and (10)

Recall that

$$\epsilon^l[t] \triangleq \frac{\partial \mathbf{u}^l[t+1]}{\partial \mathbf{u}^l[t]} + \frac{\partial \mathbf{u}^l[t+1]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]}, \quad (\text{S1})$$

which is the dependency between  $\mathbf{u}^l[t+1]$  and  $\mathbf{u}^l[t]$ . We derive Eq. (10) as below. We omit the derivation for Eq. (9) since it is a simple corollary of Eq. (10).

**Lemma 1** (Eq. (10)).

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]} &= \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t]} \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} \\ &+ \sum_{t'=t+1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t']} \frac{\partial \mathbf{u}^{l+1}[t']}{\partial \mathbf{s}^l[t']} \frac{\partial \mathbf{s}^l[t']}{\partial \mathbf{u}^l[t']} \prod_{t''=1}^{t'-t} \epsilon^l[t' - t'']. \end{aligned} \quad (\text{S2})$$

*Proof.* According to Fig. 2 and the chain rule, we have

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[T]} = \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[T]} \frac{\partial \mathbf{u}^{l+1}[T]}{\partial \mathbf{s}^l[T]} \frac{\partial \mathbf{s}^l[T]}{\partial \mathbf{u}^l[T]}, \quad (\text{S3})$$

and

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]} = \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t]} \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} + \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t+1]} \epsilon^l[t], \quad (\text{S4})$$

when  $t = T-1, \dots, 1$ . Eqs. (S3) and (S4) are standard steps in the Backpropagation through time (BPTT) algorithm.

Then we drive Eq. (S2) from Eq. (S4) by induction w.r.t.  $t$ . When  $t = T-1$ ,

$$\frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[T-1]} = \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[T-1]} \frac{\partial \mathbf{u}^{l+1}[T-1]}{\partial \mathbf{s}^l[T-1]} \frac{\partial \mathbf{s}^l[T-1]}{\partial \mathbf{u}^l[T-1]} + \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[T]} \epsilon^l[T-1], \quad (\text{S5})$$

which satisfies Eq. (S2). When  $t < T-1$ , we assume Eq. (S2) is satisfied for  $t+1$ , then show that Eq. (S2) is

satisfied for  $t$ :

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]} &= \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t]} \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} + \frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t+1]} \epsilon^l[t] \\ &= \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t]} \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} \\ &+ \left( \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t+1]} \frac{\partial \mathbf{u}^{l+1}[t+1]}{\partial \mathbf{s}^l[t+1]} \frac{\partial \mathbf{s}^l[t+1]}{\partial \mathbf{u}^l[t+1]} \right. \\ &\quad \left. + \sum_{t'=t+2}^T \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t']} \frac{\partial \mathbf{u}^{l+1}[t']}{\partial \mathbf{s}^l[t']} \frac{\partial \mathbf{s}^l[t']}{\partial \mathbf{u}^l[t']} \prod_{t''=1}^{t'-t-1} \epsilon^l[t' - t''] \right) \epsilon^l[t] \\ &= \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t]} \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} + \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t+1]} \frac{\partial \mathbf{u}^{l+1}[t+1]}{\partial \mathbf{s}^l[t+1]} \frac{\partial \mathbf{s}^l[t+1]}{\partial \mathbf{u}^l[t+1]} \epsilon^l[t] \\ &+ \sum_{t'=t+2}^T \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t']} \frac{\partial \mathbf{u}^{l+1}[t']}{\partial \mathbf{s}^l[t']} \frac{\partial \mathbf{s}^l[t']}{\partial \mathbf{u}^l[t']} \left( \prod_{t''=1}^{t'-t-1} \epsilon^l[t' - t''] \right) \epsilon^l[t] \\ &= \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t]} \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} \\ &+ \sum_{t'=t+1}^T \frac{\partial \mathcal{L}}{\partial \mathbf{u}^{l+1}[t']} \frac{\partial \mathbf{u}^{l+1}[t']}{\partial \mathbf{s}^l[t']} \frac{\partial \mathbf{s}^l[t']}{\partial \mathbf{u}^l[t']} \prod_{t''=1}^{t'-t} \epsilon^l[t' - t''], \end{aligned} \quad (\text{S6})$$

where the first equation is due to Eq. (S4), and the second equation is due to the assumption that Eq. (S2) is satisfied for  $t+1$ .  $\square$

## B. Time Complexity Analysis of SLTT and BPTT

### B.1. Pseudocode of the Backpropagation Though Time with Surrogate Gradinet Method

We present the pseudocode of one iteration of SNN training with the backpropagation though time (BPTT) with surrogate gradinet (SG) method in Algorithm S1. Note that the forward pass is defined by

$$\mathbf{u}^l[t] = (1 - \frac{1}{\tau})(\mathbf{u}^l[t-1] - V_{th} \mathbf{s}^l[t-1]) + \mathbf{W}^l \mathbf{s}^{l-1}[t], \quad (\text{S7})$$

where  $\mathbf{s}^l$  are the output spike trains of the  $l^{\text{th}}$  layer, which are calculated by:

$$\mathbf{s}^l[t] = H(\mathbf{u}^l[t] - V_{th}). \quad (\text{S8})$$

**Algorithm S1** One iteration of SNN training with the BPTT with SG method.

**Input:** Time steps  $T$ ; Network depth  $L$ ; Network parameters  $\{\mathbf{W}^l\}_{l=1}^L$ ; Training data  $(\mathbf{s}^0, \mathbf{y})$ ; Learning rate  $\eta$ .

- 1: **//Forward:**
- 2: **for**  $t = 1, 2, \dots, T$  **do**
- 3:   **for**  $l = 1, 2, \dots, L$  **do**
- 4:     Calculate  $\mathbf{u}^l[t]$  and  $\mathbf{s}^l[t]$  by Eqs. (S7) and (S8);
- 5:   **end for**
- 6: **end for**
- 7: Calculate the loss  $\mathcal{L}$  based on  $\mathbf{s}^L$  and  $\mathbf{y}$ .
- 8: **//Backward:**
- 9:  $\mathbf{e}_s^L[T] = \frac{\partial \mathcal{L}}{\partial \mathbf{s}^L[T]}$ ;
- 10: **for**  $l = L - 1, \dots, 1$  **do**
- 11:    $\mathbf{e}_u^l[T] = \mathbf{e}_s^{l+1}[T] \frac{\partial \mathbf{s}^{l+1}[T]}{\partial \mathbf{u}^l[T]}$ ,  $\mathbf{e}_s^l[T] = \mathbf{e}_u^l[T] \frac{\partial \mathbf{u}^l[T]}{\partial \mathbf{s}^l[T]}$ ;
- 12: **end for**
- 13: **for**  $t = T - 1, T - 2, \dots, 1$  **do**
- 14:   **for**  $l = L, L - 1, \dots, 1$  **do**
- 15:     **if**  $l = L$  **then**
- 16:        $\mathbf{e}_s^L[t] = \frac{\partial \mathcal{L}}{\partial \mathbf{s}^L[t]} + \mathbf{e}_u^L[t + 1] \frac{\partial \mathbf{u}^L[t + 1]}{\partial \mathbf{s}^L[t]}$ ;
- 17:     **else**
- 18:        $\mathbf{e}_s^l[t] = \mathbf{e}_u^{l+1}[t] \frac{\partial \mathbf{s}^{l+1}[t]}{\partial \mathbf{u}^l[t]} + \mathbf{e}_u^l[t + 1] \frac{\partial \mathbf{u}^l[t + 1]}{\partial \mathbf{s}^l[t]}$ ;
- 19:     **end if**
- 20:      $\mathbf{e}_u^l[t] = \mathbf{e}_s^l[t] \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]} + \mathbf{e}_u^l[t + 1] \frac{\partial \mathbf{u}^l[t + 1]}{\partial \mathbf{s}^l[t]}$ ;
- 21:      $\Delta \mathbf{W}^l += \mathbf{e}_u^l[t]^\top \mathbf{s}^{l-1}[t]^\top$ ;
- 22:   **end for**
- 23: **end for**
- 24:  $\mathbf{W}^l = \mathbf{W}^l - \eta \Delta \mathbf{W}^l$ ,  $l = 1, 2, \dots, L$ ;

**Output:** Trained network parameters  $\{\mathbf{W}^l\}_{l=1}^L$ .

## B.2. Time Complexity Analysis

The time complexity of each time step is dominated by the number of scalar multiplication operations. In this subsection, we analyze the required scalar multiplications of the Spatial Learning Through Time (SLTT) and BPTT with SG methods. We show the pseudocode of SLTT in Algorithm S2 again for better presentation.

Consider that each layer has  $d$  neurons. For simplicity, we only consider the scalar multiplications for one intermediate time step ( $t < T$ ) and one intermediate layer ( $l < L$ ). Regarding the BPTT with SG method, it requires scalar multiplications to update  $\mathbf{e}_s^l[t]$ ,  $\mathbf{e}_u^l[t]$ , and  $\Delta \mathbf{W}^l$  at lines 18, 20, and 21 respectively in Algorithm S1. To update  $\mathbf{e}_s^l[t]$ , two vector-Jacobian products are required. Since  $\frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]}$  is a diagonal matrix, the number of scalar multiplications for updating  $\mathbf{e}_s^l[t]$  is  $d^2 + d$ . To update  $\mathbf{e}_u^l[t]$  at line 20, the number is  $2d$ , since the Jacobians  $\frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]}$  and  $\frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]}$  are both diagonal. It requires  $d^2$  scalar multiplications to update  $\Delta \mathbf{W}^l$ . Regarding the SLTT method, it requires scalar multiplications to update  $\mathbf{e}_u^l[t]$  and  $\Delta \mathbf{W}^l$  at lines 12 and

13, respectively, in Algorithm S2. It requires  $d^2 + d$  scalar multiplications to update  $\mathbf{e}_u^l[t]$ , and requires  $d^2$  to update  $\Delta \mathbf{W}^l$ . Then compared with the BPTT with SG method, SLTT reduces the number of scalar multiplications by  $2d$  for one intermediate time step and one intermediate layer. For SLTT-K, it requires updating  $\mathbf{e}_u^l[t]$  and  $\Delta \mathbf{W}^l$  only at  $K$  randomly chosen time steps. Then the required number of scalar multiplication operations is proportional to  $K$ , which is proportional to  $T$  for SLTT and BPTT.

**Algorithm S2** One iteration of SNN training with the SLTT or SLTT-K methods.

**Input:** Time steps  $T$ ; Network depth  $L$ ; Network parameters  $\{\mathbf{W}^l\}_{l=1}^L$ ; Training data  $(\mathbf{s}^0, \mathbf{y})$ ; Learning rate  $\eta$ ; Required backpropagation times  $K$  (for SLTT-K).

**Initialize:**  $\Delta \mathbf{W}^l = 0$ ,  $l = 1, 2, \dots, L$ .

- 1: **if** using SLTT-K **then**
- 2:   Sample  $K$  numbers in  $[1, 2, \dots, T]$  w/o replacement to form *required\_bp\_steps*;
- 3: **else**
- 4:   *required\_bp\_steps* =  $[1, 2, \dots, T]$ ;
- 5: **end if**
- 6: **for**  $t = 1, 2, \dots, T$  **do**
- 7:   Calculate  $\mathbf{s}^L[t]$  by Eqs. (S7) and (S8);   **//Forward**
- 8:   Calculate the instantaneous loss  $\ell$ ;
- 9:   **if**  $t$  in *required\_bp\_steps* **then**   **//Backward**
- 10:      $\mathbf{e}_u^L[t] = \frac{1}{T} \frac{\partial \ell}{\partial \mathbf{s}^L[t]} \frac{\partial \mathbf{s}^L[t]}{\partial \mathbf{u}^L[t]}$ ;
- 11:     **for**  $l = L - 1, \dots, 1$  **do**
- 12:        $\mathbf{e}_u^l[t] = \mathbf{e}_u^{l+1}[t] \frac{\partial \mathbf{u}^{l+1}[t]}{\partial \mathbf{s}^l[t]} \frac{\partial \mathbf{s}^l[t]}{\partial \mathbf{u}^l[t]}$ ;
- 13:        $\Delta \mathbf{W}^l += \mathbf{e}_u^l[t]^\top \mathbf{s}^{l-1}[t]^\top$ ;
- 14:     **end for**
- 15:   **end if**
- 16: **end for**
- 17:  $\mathbf{W}^l = \mathbf{W}^l - \eta \Delta \mathbf{W}^l$ ,  $l = 1, 2, \dots, L$ ;

**Output:** Trained network parameters  $\{\mathbf{W}^l\}_{l=1}^L$ .

## C. Hard Reset Mechanism

In this work, we adopt the LIF model with soft reset as the neuron model, as shown in Eq. (2). Besides the soft reset mechanism, hard reset is also applicable for our method. Specifically, for the LIF model with hard reset

$$\begin{cases} u[t] = (1 - \frac{1}{\tau})v[t - 1] + \sum_i w_i s_i[t] + b, \\ s_{out}[t] = H(u[t] - V_{th}), \\ v[t] = u[t] \cdot (1 - s_{out}[t]), \end{cases} \quad (\text{S9})$$

we can also observe that the temporal components contribute a little to  $\frac{\partial \mathcal{L}}{\partial \mathbf{u}^l[t]}$  and the observation in Sec. 4.1 still holds. Consider the rectangle surrogate (Eq. (5)) with

Table S1: Comparison of accuracy between soft reset and hard reset on CIFAR-10 and DVS-CIFAR10.

Dataset	Reset Mechanism	Acc
CIFAR-10	Soft	94.44% $\pm$ 0.21%
	Hard	94.34%
DVS-CIFAR10	Soft	82.20 $\pm$ 0.95%
	Hard	81.40%

$\gamma = V_{th}$ , the diagonal elements of  $\epsilon^l[t]$ , which is the dependency between  $\mathbf{u}^l[t+1]$  and  $\mathbf{u}^l[t]$ , become

$$(\epsilon^l[t])_{jj} = \begin{cases} \lambda \left( 1 - (s^l[t])_j - \frac{(\mathbf{u}^l[t])_j}{V_{th}} \right), & \frac{1}{2}V_{th} < (\mathbf{u}^l[t])_j < \frac{3}{2}V_{th}, \\ \lambda(1 - (s^l[t])_j), & \text{otherwise.} \end{cases} \quad (\text{S10})$$

Since  $(s^l[t])_j \in \{0, 1\}$ , we have  $(\epsilon^l[t])_{jj} \in (-\frac{3}{2}\lambda, \frac{1}{2}\lambda) \cup \{\lambda\}$ , which is at least not large for commonly used small  $\lambda$ . As a result, the spatial components in Eqs. (9) and (10) dominate the gradients and then we can ignore the temporal components without much performance drop.

We conduct experiments with the hard reset mechanism on CIFAR-10 and DVS-CIFAR10, using the same training settings as for soft reset. And the comparison between hard reset and soft reset is shown in Tab. S1. We can see that our method can also achieve competitive results with hard reset.

## D. Dataset Description and Preprocessing

**CIFAR-10** The CIFAR-10 dataset [11] contains 60,000  $32 \times 32$  color images in 10 different classes, with 50,000 training samples and 10,000 testing samples. We normalize the image data to ensure that input images have zero mean and unit variance. We apply random cropping with 4 padding on each border of the image, random horizontal flipping, and cutout [6] for data augmentation. We apply direct encoding [18] to encode the image pixels into time series. Specifically, the pixel values are repeatedly applied to the input layer at each time step. CIFAR-10 is licensed under MIT.

**CIFAR-100** The CIFAR-100 dataset [11] is similar to CIFAR-10 except that it contains 100 classes of objects. There are 50,000 training samples and 10,000 testing samples, each of which is a  $32 \times 32$  color image. CIFAR-100 is licensed under MIT. We adopt the same data preprocessing and input encoding as CIFAR-10.

**ImageNet** The ImageNet-1K dataset [4] contains color images in 1000 classes of objects, with 1,281,167 training images and 50,000 validation images. This dataset is licensed under Custom (non-commercial). We normalize the image data to ensure that input images have zero mean and

unit variance. For training samples, we apply random re-sized cropping to get images with size  $224 \times 224$ , and then apply horizontal flipping. For validation samples, we re-size the images to  $256 \times 256$  and then center-cropped them to  $224 \times 224$ . We transform the pixels into time sequences by direct encoding [18], as done for CIFAR-10 and CIFAR-100.

**DVS-Gesture** The DVS-Gesture [1] dataset is recorded using a Dynamic Vision Sensor (DVS), consisting of spike trains with two channels corresponding to ON- and OFF-event spikes. The dataset contains 11 hand gestures from 29 subjects under 3 illumination conditions, with 1176 training samples and 288 testing samples. The license of DVS-Gesture is Creative Commons Attribution 4.0. For data preprocessing, we follow [8] to integrate the events into frames. The event-to-frame integration is handled with the Spiking-Jelly [7] framework.

**DVS-CIFAR10** The DVS-CIFAR10 dataset [12] is a neuromorphic dataset converted from CIFAR-10 using a DVS camera. It contains 10,000 event-based images with pixel dimensions expanded to  $128 \times 128$ . The dataset is licensed under CC BY 4.0. We split the whole dataset into 9000 training images and 1000 testing images. Regarding data preprocessing, we integrate the events into frames [8], and reduce the spatial resolution into  $48 \times 48$  by interpolation. For some experiments, we take random horizontal flip and random roll within 5 pixels as data augmentation, the same as [5].

## E. Network Architectures

### E.1. Scaled Weight Standardization

An important characteristic of the proposed SLTT method is the instantaneous gradient calculation at each time step, which enables time-steps-independent memory costs. Under our instantaneous update framework, an effective technique, batch normalization (BN) along the temporal dimension [23, 13, 5, 15], cannot be adopted to our method, since this technique requires gathering data from all time steps to calculate the mean and variance statistics. For some tasks, we still use BN components, but their calculated statistics are based on data from each time step. For other tasks, we replace BN with scaled weight standardization (sWS) [17, 2, 3], as introduced below.

The sWS component [2], which is modified from the original weight standardization [17], normalizes the weights according to:

$$\hat{\mathbf{W}}_{i,j} = \gamma \frac{\mathbf{W}_{i,j} - \mu_{\mathbf{W}_{i,\cdot}}}{\sigma_{\mathbf{W}_{i,\cdot}}}, \quad (\text{S11})$$

where the mean  $\mu_{\mathbf{W}_{i,\cdot}}$  and standard deviation  $\sigma_{\mathbf{W}_{i,\cdot}}$  are calculated across the fan-in extent indexed by  $i$ ,  $N$  is the dimension of the fan-in extent, and  $\gamma$  is a fixed hyperparameter. The hyperparameter  $\gamma$  is set to stabilize the signal propagation in the forward pass. Specifically, for one network layer  $\mathbf{z} = \hat{\mathbf{W}}g(\mathbf{x})$ , where  $g(\cdot)$  is the activation and the elements of  $\mathbf{x}$  are i.i.d. from  $\mathcal{N}(0, 1)$ , we determine  $\gamma$  to make  $\mathbb{E}(\mathbf{z}) = 0$ , and  $\text{Cov}(\mathbf{z}) = \mathbf{I}$ . For SNNs, the activation  $g(\cdot)$  at each time step can be treated as the Heaviside step function. Then we take  $\gamma \approx 2.74$  to stable the forward propagation, as calculated by [22]. Furthermore, sWS incorporate another learnable scaling factor for the weights [2, 22] to mimic the scaling factor of BN. sWS shares similar effects with BN, but introduces no dependence of data from different batches and time steps. Therefore, sWS is a convincing alternative for replacing BN in our framework.

## E.2. Normalization-Free ResNets

For deep ResNets [9], sWS cannot enjoy the similar signal-preserving property as BN very well due to the skip connection. Then in some experiments, we consider the normalization-free ResNets (NF-ResNets) [2, 3] that not only replace the BN components by sWS but also introduce other techniques to preserve the signal in the forward pass.

The NF-ResNets use the residual blocks of the form  $x_{l+1} = x_l + \alpha f_l(x_l/\beta_l)$ , where  $\alpha$  and  $\beta_l$  are hyperparameters used to stabilize signals, and the weights in  $f_l(\cdot)$  are imposed with sWS. The carefully determined  $\beta_l$  and the sWS component together ensure that  $f_l(x_l/\beta_l)$  has unit variance.  $\alpha$  controls the rate of variance growth between blocks, and is set to be 0.2 in our experiments. Please refer to [2] for more details on the network design.

## E.3. Description of Adopted Network Architectures

We adopt ResNet-18 [9] with the pre-activation residual blocks [10] to conduct experiments on CIFAR-10 and CIFAR-100. The channel sizes for the four residual blocks are 64, 128, 256, and 512, respectively. All the ReLU activations are substituted by the leaky integrate and fire (LIF) neurons. To make the network implementable for neuromorphic computing, we replace all the max pooling operations with average pooling. To enable instantaneous gradient calculation, we adopt the BN components that calculate the mean and variance statistics for each time step, not the total time horizon. Then for each iteration, a BN component is implemented for  $T$  times, where  $T$  is the number of total time steps.

We adopt VGG-11 [20] to conduct experiments on DVS-Gesture and DVS-CIFAR10. As done for Resnet-18, we substitute all the max poolings with average poolings and use the time-step-wise BN. We remove two fully connected layers [15, 5, 22] to reduce the computation. A dropout layer [21] is added behind each LIF neuron layer to bring

better generalization. The dropout rates for DVS-Gesture and DVS-CIFAR10 are set to be 0.4 and 0.3, respectively.

Table S2: Training hyperparameters about optimization. “WD” means weight decay, “LR” means initial learning rate, and “BS” means batch size.

Dataset	Epoch	LR	BS	WD
CIFAR-10	200	0.1	128	$5 \times 10^{-5}$
CIFAR-100	200	0.1	128	$5 \times 10^{-4}$
ImageNet (pre-train)	100	0.1	256	$1 \times 10^{-5}$
ImageNet (fine-tune)	30	0.001	256	0
DVS-Gesture	300	0.1	16	$5 \times 10^{-4}$
DVS-CIFAR10	300	0.05	128	$5 \times 10^{-4}$

We also adopt VGG-11 (WS) to conduct experiments on DVS-Gesture and achieve state-of-the-art performance. VGG-11 (WS) is a BN-free network that shares a similar architecture with VGG-11 introduced above. The difference between the two networks is that VGG-11 consists of the convolution-BN-LIF blocks while VGG-11 (WS) consists of the convolution-sWS-LIF blocks.

We adopt NF-ResNet-34, NF-ResNet-50, and NF-ResNet-101 to conduct experiments on ImageNet. Those networks are normalization-free ResNets introduced in Appendix E.2. In Figs. 1 and 3 of the main content, the used network for ImageNet is NF-ResNet-34.

## F. Training Settings

All the implementation is based on the PyTorch [16] and SpikingJelly [7] frameworks, and the experiments are carried out on one Tesla-V100 GPU or one Tesla-A100 GPU.

For CIFAR-10, CIFAR-100, and DVS-Gesture, we adopt the loss function proposed in [5]:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T (1 - \lambda) \ell_1(\mathbf{o}[t], y) + \lambda \ell_2(\mathbf{o}[t], V_{th}), \quad (\text{S12})$$

where  $T$  is the number of total time steps,  $\ell_1$  is the cross entropy function,  $\ell_2$  is the mean squared error (MSE) function,  $\mathbf{o}[t]$  is the network output at the  $t$ -th time step,  $y$  is the label,  $\lambda$  is a hyperparameter taken as 0.05, and  $V_{th}$  is the spike threshold which is set to be 1 in this work. For ImageNet, we use the same loss but simply set  $\lambda = 0$ . For DVS-CIFAR10, we also combine the cross entropy loss and the MSE loss, but the MSE loss does not act as a regularization term as in [5]:

$$\mathcal{L} = \frac{1}{T} \sum_{t=1}^T (1 - \lambda) \ell_1(\mathbf{o}[t], y) + \lambda \ell_2(\mathbf{o}[t], y), \quad (\text{S13})$$

where  $\alpha$  is also taken as 0.05. Our experiments show that such loss performs better than that defined in Eq. (S12) for DVS-CIFAR10.



For all the tasks, we use SGD [19] with momentum 0.9 to train the networks, and use cosine annealing [14] as the learning rate schedule. Other hyperparameters about optimization are listed in Tab. S2. For ImageNet, we first train the SNN with only 1 time step to get a pre-trained model, and then fine-tune the model for multiple time steps.

In Sec. 5.3 of the main content, we compare the proposed SLTT method and OTTT [22] following the same experimental settings as introduced in [22]. For both methods, we adopt the NF-ResNet-34 architecture for ImageNet and VGG-11 (WS) for other datasets. And the total number of time steps for CIFAR-10, CIFAR-100, ImageNet, DVS-Gesture, and DVS-CIFAR10 are 6, 6, 6, 20, and 10, respectively.

### References

- [1] Arnon Amir, Brian Taba, David Berg, Timothy Melano, Jeffrey McKinstry, Carmelo Di Nolfo, Tapan Nayak, Alexander Andreopoulos, Guillaume Garreau, Marcela Mendoza, et al. A low power, fully event-based gesture recognition system. In *CVPR*, 2017. 3
- [2] Andrew Brock, Soham De, and Samuel L Smith. Characterizing signal propagation to close the performance gap in unnormalized resnets. In *ICLR*, 2021. 3, 4
- [3] Andy Brock, Soham De, Samuel L Smith, and Karen Simonyan. High-performance large-scale image recognition without normalization. In *International Conference on Machine Learning*, pages 1059–1071. PMLR, 2021. 3, 4
- [4] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 3
- [5] Shikuang Deng, Yuhang Li, Shanghang Zhang, and Shi Gu. Temporal efficient training of spiking neural network via gradient re-weighting. In *ICLR*, 2022. 3, 4
- [6] Terrance DeVries and Graham W Taylor. Improved regularization of convolutional neural networks with cutout. *arXiv preprint arXiv:1708.04552*, 2017. 3
- [7] Wei Fang, Yanqi Chen, Jianhao Ding, Ding Chen, Zhaofer Yu, Huihui Zhou, Yonghong Tian, and other contributors. Spikingjelly. <https://github.com/fangwei123456/spikingjelly>, 2020. 3, 4
- [8] Wei Fang, Zhaofer Yu, Yanqi Chen, Timothée Masquelier, Tiejun Huang, and Yonghong Tian. Incorporating learnable membrane time constant to enhance learning of spiking neural networks. In *ICCV*, 2021. 3
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 4
- [10] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Identity mappings in deep residual networks. In *ECCV*, 2016. 4
- [11] Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. 2009. 3
- [12] Hongmin Li, Hanchao Liu, Xiangyang Ji, Guoqi Li, and Luping Shi. Cifar10-dvs: an event-stream dataset for object classification. *Frontiers in neuroscience*, 11:309, 2017. 3
- [13] Yuhang Li, Yufei Guo, Shanghang Zhang, Shikuang Deng, Yongqing Hai, and Shi Gu. Differentiable spike: Rethinking gradient-descent for training spiking neural networks. In *NeurIPS*, 2021. 3
- [14] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017. 5
- [15] Qingyan Meng, Mingqing Xiao, Shen Yan, Yisen Wang, Zhouchen Lin, and Zhi-Quan Luo. Training high-performance low-latency spiking neural networks by differentiation on spike representation. In *CVPR*, 2022. 3, 4
- [16] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. In *NeurIPS*, 2019. 4
- [17] Siyuan Qiao, Huiyu Wang, Chenxi Liu, Wei Shen, and Alan Yuille. Micro-batch training with batch-channel normalization and weight standardization. *arXiv preprint arXiv:1903.10520*, 2019. 3
- [18] Nitin Rathi and Kaushik Roy. Diet-snn: A low-latency spiking neural network with direct input encoding and leakage and threshold optimization. *TNNLS*, 2021. 3
- [19] David E Rumelhart, Geoffrey E Hinton, and Ronald J Williams. Learning representations by back-propagating errors. *Nature*, 323(6088):533–536, 1986. 5
- [20] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2014. 4
- [21] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014. 4
- [22] Mingqing Xiao, Qingyan Meng, Zongpeng Zhang, Di He, and Zhouchen Lin. Online training through time for spiking neural networks. In *NeurIPS*, 2022. 4, 5
- [23] Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li. Going deeper with directly-trained larger spiking neural networks. In *AAAI*, 2021. 3