# A APPENDIX

## A.1 EXPERIMENTS ON VGG STRUCTURES

We evaluated our method on the CIFAR100 dataset using VGG architectures, besides ResNets. We treated each layer as a stage in the VGG series. During experimentation, we observed that VGG16 with three fully connected (fc) layers could not be trained effectively using the standard direct training approach (its accuracy remained limited at 1%). To tackle this issue, we merged the last three fc layers of VGG16 into one and named the resulting architecture VGG14. We set $T_{max} = 5$, $s = 3$, and computed the mean and standard deviation of three runs. We used the SGD optimizer to train the model, with a learning rate of 0.1, a weight decay of 0.0005, and a batch size of 256. The results, presented in Table 9, indicate the effectiveness of our approach on VGG structures.

**Table 9:** Accuracy of VGG on CIFAR100

| Methods | Model | T=2 | T=3 | T=4 | T=5 |
|---|---|---|---|---|---|
| MTT | VGG14 | 73.53±0.09 | 74.52±0.07 | 75.27±0.14 | 75.72±0.10 |
| InfLoR-SNN | VGG16 | - | - | - | 71.56±0.10 |

## A.2 TRAINING DETAILS

**CIFAR** The CIFAR10/CIFAR100 dataset comprises 50K training images and 10K test images with a 32×32 pixel resolution. For CIFAR100, we train a ResNet-19 TFSNN using the MTT pipeline for 300 epochs with a batch size of 256 and a $T_{max}$ of 6. Following the practice in GLIF (Yao et al., 2022), the last 2 fully connected layers of ResNet-19 are replaced with a single fully connected layer. We employ the SGD optimizer with a weight decay of 0.0005 and a learning rate of 0.1 cosine decayed to 0. To make a fair comparison with the state-of-the-art (SOTA) work (Li et al., 2021b; Guo et al., 2022; Yao et al., 2022), AutoAugment (Cubuk et al., 2018) and Cutout (DeVries & Taylor, 2017) are applied to both CIFAR10 and CIFAR100 datasets. However, these augmentation techniques are only used for comparative experiments and temporal flexibility experiments, and not for other experiments.

**ImageNet** ImageNet (Deng et al., 2009) contains more than 1280k training images and 50k test images. We use the standard data processing flow to crop each image to a size of 224×224. We deploy the ResNet-34 structure, however, with the removal of the first max-pooling layer and changing the stride of the first basic block from 1 to 2 (Zheng et al., 2021; Yao et al., 2022). We train the model for 160 epochs with a batch size of 512 and a $T_{max}$ of 6. We utilize the AdamW optimizer with a weight decay of 0.02 and a learning rate of 0.004 cosine decayed to 0.

**DVS-Dataset** CIFAR10-DVS and N-Caltech101 are neuromorphic datasets widely used in SNN experimentation. We divide the dataset into a 9:1 ratio and merge all events to form ten frames, which, similar to previous work (Li et al., 2021b; Deng et al., 2022), we resize to 48×48. For both these datasets, we adopt a random horizontal flip and rotate the frames up to 5 pixels as augmentation techniques. We employ the additional temporal inversion policy (Shen et al., 2023) uniquely for N-Caltech101. For these datasets, we use a $T_{max}$ of 10, a batch size of 50, and train the TFSNN ResNet-18 model for 300 epochs. The optimizer we choose is SGD, with a weight decay of 0.0005, and a learning rate of 0.1, which we cosine decay to 0. While training the DVS dataset, we take only the first $t$ frames of the ten frames where $t$ denotes the time step of the input stage, to feed into the network.

## A.3 DETAILS OF COMBINATORIAL OPTIMIZATION

As previously mentioned, the accuracy formula of TFSNN is given by the expression $\sum_{i=1}^{I} K_i \sqrt{\log_2 t_i} + c$, where $K_i$ represents the contribution of each block, $I$ is the number of blocks, and $c$ is a bias. We randomly select 18 distinct temporal configurations ($\boldsymbol{t}$) and evaluate their accuracies on the test set, resulting in 18 pairs of temporal configurations and their corresponding accuracies. Using the least squares method, we compute the values of $K_i$ and $c$ from the collected data. Next, we estimate the average firing rate ($R_i$) of each block in a unified SNN of T=6. Then, the energy consumption of a specified temporal configuration $t$ can be approximated as $\sum_{i=1}^{I} t_i \cdot R_i$. For example, the estimated energy consumption of a unified SNN with T=4 is calculated as $EC_4 = \sum_{i=1}^{I} 4R_i$. Based on this, we can obtain a group of TFSNNs with lower

energy consumption ($EC$) for a given T=$T_g$, and we aim to identify the TFSNN with the maximum estimated accuracy from this set. This is formulated as the following optimization problem:

$$\text{maximize} \quad \text{ACC}_{\text{estimated}} = \sum_{i=1}^{I} K_i \sqrt{\log_2 t_i} + c$$

$$\text{s.t.} \quad \sum_{i=1}^{I} t_i \cdot R_i \leq EC_{T_g}$$

$$T_{min} \leq t_1, t_2, \ldots, t_l \leq T_{max},$$

where $t_i$ is the $i$-th component of temporal configurations $\boldsymbol{t}$ and $EC_{T_g}$ is the given uper bound of the TFSNN energy.

In order to solve the above problem, we adopt the depth-first search (DFS) algorithm to search in the solution space. To obtain a more accurate accuracy for each temporal configuration $\boldsymbol{t}$, we perform three times of BN calibrations and take the average of the accuracies.

### A.4  DESIGN TTM GROUPING POLICY

We previously stated our policy's objective is to partition $l$ frames into $k$ groups ($l \geq k$) as evenly as possible. In this section, we mathematically interpret the design. We will start by describing the grouping process in a different manner. The $l$ frames are viewed as $l$ adjacent intervals of length 1 over the rational number domain with the $i$-th frame starting at $i-1$ and ending at $i$. We define $c_i$ as the boundary between group $i$ and group $i-1$. Here, $i$ ranges from 1 to $k$, and $c_1$ is 0. Ideally, $c_i = (i-1) \cdot l / k$ is set to group frames most evenly. Nevertheless, this strategy produces non-integer $c_i$, which results in atomic frames' division when $l$ is not a multiple of $k$. To solve this problem, we retreat and set $c_i$ to the nearest integer and get

$$c_i = \lfloor \frac{(i-1) \cdot l}{k} - \varepsilon \rceil, \tag{11}$$

where $\varepsilon$ is a small constant used to determine $c_i$ when the distances to the closest two integers are equal. As $b_i$ must be the frame directly following the boundary $c_i$ ($b_i = c_i + 1$), we obtain Eq. 8.

### A.5  DERIVATION OF THE BACKPROPAGATION FORMULA FOR LIF

In this section, we derive Eq. 7 from the forwarding formula. We derive $\partial u(t)/\partial v(t)$ from Eq. 4 first:

$$\frac{\partial u(t)}{\partial v(t)} = 1 - s(t) - v(t) \cdot \frac{\partial s(t)}{\partial v(t)}. \tag{12}$$

Then, we consider the derivation of $\partial u(t)/\partial v(t-1)$. According to Fig. 2, $\partial u(t)/\partial v(t-1)$ can be calculated as follows

$$\frac{\partial u(t)}{\partial v(t-1)} = \frac{\partial u(t)}{\partial v(t)} \frac{\partial v(t)}{\partial u(t-1)} \frac{\partial u(t-1)}{\partial v(t-1)} = \tau \frac{\partial u(t)}{\partial v(t)} \frac{\partial u(t-1)}{\partial v(t-1)}. \tag{13}$$

By combining multiple Eq. 13, we get

$$\frac{\partial u(t)}{\partial v(t-n)} = \tau^n \prod_{i=t-n}^{t} \frac{\partial u(i)}{\partial v(i)}. \tag{14}$$

Finally, we get the complete expression for $\partial s(t)/\partial I(t-n)$ as follows

$$\frac{\partial s(t)}{\partial I(t-n)} = \frac{\partial s(t)}{\partial v(t)} \frac{\partial v(t)}{\partial u(t-1)} \frac{\partial u(t-1)}{\partial v(t-n)} \frac{\partial v(t-n)}{\partial I(t-n)}$$

$$= \frac{\partial s(t)}{\partial v(t)} \cdot \tau^n \prod_{i=t-n}^{t-1} [(1 - s(i)) - v(i) \cdot \frac{\partial s(i)}{\partial v(i)}] \tag{15}$$

## A.6 Details of Handmade Model Structures

To demonstrate MTT's capacity for training deep networks, we created two deep networks, B-ResNet-50 and B-ResNet-72. We will present their specific structures in this section. Table 10 displays the architectures of these two networks. For convenience, we still use basic blocks for each block, and simply reconfigure the number of basic blocks in each part. Since the model is large, we use $T_{max} = 4$.

**Table 10:** Structures of our handmade deep networks on CIFAR.

| Stage | Output Size | B-ResNet-50 | | B-ResNet-72 | |
|---|---|---|---|---|---|
| conv1 | 32×32 | 3x3, 64 | | | |
| conv2_x | 32×32 | 3x3, 64 <br> 3x3, 64 | * 3 | 3x3, 64 <br> 3x3, 64 | * 3 |
| conv3_x | 16×16 | 3x3, 128 <br> 3x3, 128 | * 6 | 3x3, 128 <br> 3x3, 128 | * 8 |
| conv4_x | 8×8 | 3x3, 256 <br> 3x3, 256 | * 9 | 3x3, 256 <br> 3x3, 256 | * 16 |
| conv5_x | 4×4 | 3x3, 512 <br> 3x3, 512 | * 6 | 3x3, 512 <br> 3x3, 512 | * 8 |
| FC | 1×1 | average pool, fc, softmax | | | |

## A.7 T = T$_{MAX}$ BN Statistics vs. Recalculated BN Statistics

As previously mentioned, we discovered that the BN statistics of the T=T$_{max}$ network can be applied to other TFSNN with uniform T across blocks. In this section, we provide experimental verification for this observation. Our experiment involves testing the accuracy of one of our trained ResNet-19 models with two distinct BN layer information approaches. The first approach utilizes the statistics of the T=T$_{max}$ network, while the second approach recalculates the BN statistics individually for each time step T. For the latter approach, we calibrate the BN layer three times and report the average accuracy. Our experimental results demonstrated in Table 12 show that the mean and variance calculated at T=T$_{max}$ is applicable directly to other T values. Therefore, we can utilize the BN statistics of T=T$_{max}$ for other T directly, which saves the time required to calibrate the BN layers for other T values.

**Table 11:** Accuracy given sampling number $s$ and training epochs $e$.

| Sampling Num | $e \times s$ | | | |
|---|---|---|---|---|
| | 300 | 600 | 900 | 1200 |
| $s = 1$ | **75.61** | 76.13 | 76.23 | 75.78 |
| $s = 2$ | 75.53 | 76.37 | 76.31 | 76.36 |
| $s = 3$ | 75.25 | **76.45** | **76.47** | **76.44** |
| $s = 4$ | 74.81 | 75.74 | 76.41 | 76.42 |

**Table 12:** Test accuracy of a single model with two kinds of BN statistics.

| Method | TimeStep | | | |
|---|---|---|---|---|
| | 2 | 3 | 4 | 5 |
| T=T$_{max}$ stat | 80.21 | 81.06 | **81.51** | 81.82 |
| Recalculated stat | 80.21 | **81.23** | 81.44 | **81.85** |

## A.8 Impact of Different Sampling Number and Training Epochs

In most previous experiments, we employed sampling number $s = 3$. In this section, we experiment with varying values of $s$, assess their effects at different epochs, and explain why we chose $s = 3$. We train ResNet-18s with $T_{max} = 6$ on CIFAR100 with varied $s$ and list their test accuracy at T=6. The results are displayed in Table 11. When $e \times s$ is constrained, the model requires more epochs to converge, necessitating a lower $s$. However, if trained across a sufficient number of epochs, sampling of $s$ structurally diverse networks can smooth the optimization of network parameters and improve performance. Specifically, we find $s = 3$ performs well and adopt $s = 3$ for most of the experiments.
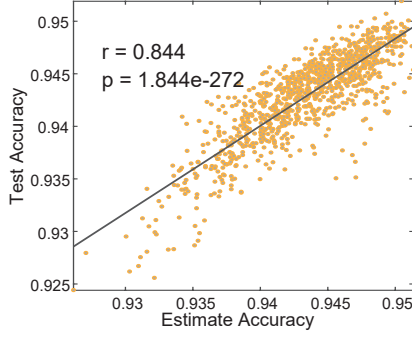
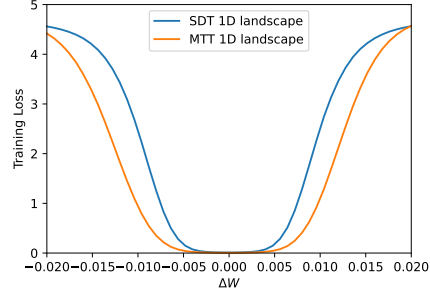**Figure 8:** Correlation curve for estimate accuracy and test accuracy on CIFAR10



**Figure 9:** The 1D landscapes of ResNet18 trained by SDT and MTT on CIFAR100.

### A.9 TFSNN ACCURACY ESTIMATION ON CIFAR10

In addition to CIFAR100, we also conducted experiments with ResNet-18 on CIFAR10. We randomly sample 18 temporal configurations as usual and solve the equation in appendix A.3. The weight parameters are $\{0.0049, 0.0028, 0.0017, 0.0037, 0.0037, 0.0004, 0.0005, 0.0017\}$, and the constant value $c$ is 92.13. Notes that block 1,4,5 have a higher contribution, and the 1,5 blocks are also highly weighted on CIFAR100, which may imply that the weights are partly related to the network structure. We then resample 1000 TFSNN and plot their estimate accuracy and test accuracy on CIFAR10 in Fig. 8. We also use the aforementioned combinatorial optimization strategy to search the optimal TFSNN under the energy consumption of T=3 and find $\{4, 2, 2, 3, 5, 2, 2, 4\}$, which achieves an accuracy of 94.70%, 0.21% higher than its T=3 counterpart.

### A.10 LOSS LANDSCAPES OF MTT AND SDT

To visually confirm the flatter minimum achieved by the model trained with MTT, we trained ResNet18 using SDT and MTT on CIFAR100 and plotted their loss landscapes in Fig. 9. We observed that MTT led the model to a flatter minimum which indicates improved generalizability.

### A.11 VERIFYING GENERALIZABILITY THROUGH GRADIENT METRICS

Apart from noise injection, another famous metric that indicates the generalizability is the length of the gradient on weights $||\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}}||$ and the inputs $||\frac{\partial \mathcal{L}}{\partial \boldsymbol{x}_i}||$. For $||\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}}||$, we evaluate the length of the gradient of loss over the entire training set for the convolution layers. For $||\frac{\partial \mathcal{L}}{\partial \boldsymbol{x}_i}||$, we calculate the mean value of the length of each input gradient. The model trained by MTT exhibits a shorter gradient of both weights and inputs (see Table 13), which implies the model's strong robustness and generalizability.

**Table 13:** The gradient statistics of the model trained by SDT and MTT.

| Methods | $||\frac{\partial \mathcal{L}}{\partial \boldsymbol{W}}||$ | $||\frac{\partial \mathcal{L}}{\partial \boldsymbol{x}_i}||$ |
|---|---|---|
| MTT | 11.59 | 1.81 |
| SDT | 38.08 | 7.78 |

**Table 14:** Results on seqMNIST.

| Methods | Acc |
|---|---|
| Our RNN | 56.22 |
| SNN SDT | 55.75 |
| SNN MTT | **64.56** |

**Table 15:** Results on Spiking Heidelberg Digits

| Methods | Acc |
|---|---|
| l=3 Repro by code of Hammouamri et al. (2023) | 75.26 |
| l=3 our SDT | 74.43 |
| l=3 our MTT | **79.68** |

**Table 16:** Results on Spiking Speech Commands

| Methods | Acc |
|---|---|
| Our SDT l=3 | 57.75 |
| Our MTT l=3 | **60.15** |

### A.12 EXPERIMENTS ON AUDIO AND SEQUENTIAL DATASETS

Our research reveals that TFSNN can function effectively as a time encoder when the temporal configuration vectors used for training are monotonically non-increasing.

To illustrate this adaptability, we present the performance of TFSNN on three distinct temporal tasks: seqMNIST, Spiking Heidelberg Digits and Spiking Speech Commands.

For the sequential task seqMNIST, we utilized a simple fc LIF SNN with 2 hidden layers of width 64 and set the time constant $\tau$ to 0.99. We also trained an RNN with 2 hidden layers of width 64 for comparison. The results are as shown in Table 14.

For the Spiking Heidelberg Digits, we adopt the plain 3-layer feed-forward SNN architecture proposed by Cramer et al. (2020), a fully connected SNN with an input width of 70 and 128 LIF neurons in each of the 3 hidden layers. The timestep of the first layer is fixed to the input timestep, while the timesteps of subsequent layers are restricted to monotonically non-increasing. We set $\tau = 0.9753$, which is equivalent to the parameter $\lambda$ in the work of Cramer et al. (2020), namely $1 - 1/\tau$ in most other articles, and train the model for 150 epochs. To ensure the validity of our results, we also reproduce the result using the code provided by Hammouamri et al. (2023). The results are as shown in Table 15.

Spiking Speech Commands (SSC) (Cramer et al., 2020) is a spiking dataset converted from Google Speech Commands v0.2 and is tailored for SNN. For SSC, we continue using the same architecture and the same parameters as we used in SHD, except that here we only train the model for 60 epochs. The results are as shown in Table 16.