# Efficient Neuromorphic Data Processing: A Unified Framework for Pipeline Optimization and Lightweight Preprocessing*

1st Given Name Surname
*dept. name of organization (of Aff.)*
*name of organization (of Aff.)*
City, Country
email address or ORCID

2nd Given Name Surname
*dept. name of organization (of Aff.)*
*name of organization (of Aff.)*
City, Country
email address or ORCID

3rd Given Name Surname
*dept. name of organization (of Aff.)*
*name of organization (of Aff.)*
City, Country
email address or ORCID

4th Given Name Surname
*dept. name of organization (of Aff.)*
*name of organization (of Aff.)*
City, Country
email address or ORCID

5th Given Name Surname
*dept. name of organization (of Aff.)*
*name of organization (of Aff.)*
City, Country
email address or ORCID

6th Given Name Surname
*dept. name of organization (of Aff.)*
*name of organization (of Aff.)*
City, Country
email address or ORCID

*Abstract*—To address the computational efficiency bottleneck in real-time processing of neuromorphic data (e.g., DVS128 Gesture), this paper proposes an end-to-end acceleration solution that integrates data pipeline optimization and lightweight preprocessing. At the system level, the use of `num_workers` for parallel loading and `prefetch_factor` for prefetching significantly enhances the throughput of the data pipeline. At the algorithmic level, a dynamic time window event frame aggregation method based on spatiotemporal sparsity is introduced, along with a low-computation pulse noise filtering module leveraging local spatiotemporal correlations. Additionally, structured model pruning is employed to reduce redundant connections and lower computational overhead. Experimental results demonstrate that on an NVIDIA 4070 GPU, the inference speed on the test set is doubled, while achieving classification accuracies of 92.36% (test set) and 98.97% (training set). This work combines data pipeline optimization with lightweight preprocessing, providing a high real-time solution for neuromorphic data processing in edge computing scenarios, with significant practical engineering value.

*Index Terms*—Neuromorphic Computing, Data Pipeline Optimization, Pulse Noise Filtering

## I. INTRODUCTION

Spiking Neural Networks(SNN),which are regarded as the third generation of neural network models,offer a more biologically plausible simulation of neural dynamic by incorporating not only neuronal and synaptic states but also temporal information into their computational framework. [1], [2]Unlike traditional artificial neural networks (ANNs),which relay on continuous-valued activations,SNNs process information through discrete spike events,letting them to capture the temporal dependencies inherent in biological systems.This

temporal sensivity,combined with energy-efficient computation,allows SNNs to become as a promising alternative for modeling complex cognitive tasks with potentially superior performance and lower power consumption. [3]

Neuromorphic data, characterized by its event-driven nature and sparse sensory representation, has emerged as a compelling paradigm across various real-time application domains, including gesture recognition, object tracking, and autonomous navigation. Neuromorphic data, characterized by its event-driven nature and sparse sensory representation, has emerged as a compelling paradigm across various real-time application domains, including gesture recognition, object tracking, and autonomous navigation.However, the efficient processing of such data remains a significant challenge, particularly in edge computing scenarios where computational resources are limited. A critical bottleneck lies in the substantial latency introduced by data loading and preprocessing. Traditional approaches to address this issue often focus on either algorithmic optimizations or hardware acceleration in isolation, failing to fully exploit the synergistic potential of a unified framework. [4]

This paper proposes a novel solution that seamlessly integrates data pipeline optimization with lightweight preprocessing techniques to achieve end-to-end acceleration. At the system level, we leverage parallel loading and prefetching mechanisms to maximize data throughput. On the algorithmic front, we introduce a computationally efficient pulse noise filtering module based on local spatiotemporal correlations.This module enhances the robustness of our model against noise while preserving critical temporal dynamics. Additionally, we incorporate a self-attention mechanism to augment the

model's understanding of feature channel interdependencies, thereby improving its representational capacity. [5] To further enhance computational efficiency, we employ structured model pruning,which systematically reduces redundant connections. [6]

Our contributions are threefold:

(1)we compared the effects on accuracy and training speed when setting different frame numbers in the ddateset.

(2)we proposed a lightweight preprocessing method that combines pulse noise filtering and self-attention mechanisms based on the network proposed by WeiFang [7]to enhance the model's robustness and representational capacity.Meanwhile we conducted experiments on the DVS-Gesture dataset and achieved experimental results.

(3)we introduced structured model pruning to reduce redundant connections, thereby improving computational efficiency.

## II. RELATED WORK

### A. Spiking Neural Networks

Spiking Neural Networks(SNNs) mimics the pulse emission characteristics of biological neurons (such as the LIF model), drives computing with discrete events, and has significant advantages in energy-sensitive scenarios (such as edge devices). [3] A major research focus in the field of SNNs is the conversion from artificial neural networks(ANNs) to SNNs [8],which is dedicated to addressing the challenges of accuracy degradation and training instability during the conversion process. [9] Zhang et al. (2025) proposed a method called STNN-SNN, a Spatial-Temporal Attention Aggregator SNN framework,which can dynamically attend to and capture dependencies between spatial and temporal domains. [10]

### B. Spatial-Temporal-Polarity Coherence Filtering

In event-based neuromorphic data processing, spatio-temporal polarity coherent filtering is an effective denoising processing method. Unlike conventional frame-based filtering, this method leverages the unique characteristics of event streams, where each event is defined by its spatial coordinates, timestamp, and polarity. Spatio-temporal polarity coherent filtering has been widely applied in event-based neuromorphic vision tasks. Sironi et al.(2018) used polarity and spatial-temporal neighborhood statistical characteristics to robustly classify events. [11]

### C. Attention Mechanism of SNNs

Attention mechanism has attracted considerable attention since its inception and has been widely used in various neural network architectures.Xue et al. (2024) introduced Gated Attention Coding (GAC) [12], a plug-and-play module that encodes inputs into powerful representations.Yu et al. (2024) proposed a Frequency-based Spatial-Temporal Attention(FSTA) model to enhance feature learning in SNNs. [13] Shen et al.(2024) innovatively implemented an attention mechanism in SNNs called Temporal Interaction Module (TIM),which was designed to augment the temporal data processing abilities within SNN architectures. [14]

## III. METHODOLOGY

## IV. EXPERIMENTS

### A. Experimental Setup

**Dataset Description** In this experiment, we utilize the DVS Gesture dataset, a publicly available dataset for gesture recognition based on Dynamic Vision Sensor (DVS). The dataset comprises 1,342 samples, each corresponding to one of 11 distinct gesture categories. Each gesture is performed multiple times by approximately 29 different subjects. Each sample is stored in the form of an event stream. The dataset also provides training and test set splits, facilitating model evaluation on real-world event data. The data was collected under varying lighting and motion conditions to enhance diversity.

**Evaluation Metrics** To comprehensively evaluate model performance, we employ five key metrics: (1) Training speed (samples/sec), measuring computational efficiency by tracking processed samples per second during training; (2) Training accuracy (%), monitoring the model's learning progress by calculating correct predictions on training data; (3) Test accuracy (%), the primary classification metric evaluated on held-out test samples; (4) Training loss, recorded as an observational metric to analyze model behavior during training; (5) Test loss, evaluating generalization capability on held-out data.

All experiments were conducted on a single NVIDIA RTX 4070 GPU with a batch size of 16, using the Adam optimizer, Automatic Mixed Precision (AMP) for faster computation, and the CuPy library for GPU-accelerated array operations. This multi-faceted evaluation captures both computational efficiency (training speed) and model effectiveness (accuracy/loss), while the training loss serves purely as a diagnostic indicator rather than a control signal.

### B. Impact of Frame Count on Model Performance

When processing DVS Gesture data, it is often necessary to convert the raw event stream into some form of image representation (e.g., event frames, accumulated frames, or time surfaces) so that it can be input into a convolutional neural network (CNN) for classification or other tasks. In this process, frame number is a key parameter, which indicates how many time windows the entire event stream should be divided into. Within each window, an "event frame" is generated as an input feature.

In this experiment, we process the DVS Gesture dataset on different frame numbers, as shown in Fig 1. Each row represents the processed dataset with different frame numbers, from top to bottom the number of frames is 2, 4, 8, and 16. And each column represents a different category, from left to right, which is guitar, right hand wave, forearm roll forward, left hand wave, right hand counter clockwise, right hand clockwise, left hand counter clockwise, left hand clockwise, clap, drums and others.

Through the sample of dataset, we can see that for most categories there is only a change in brightness caused by the difference of frame number, but for certain categories, such
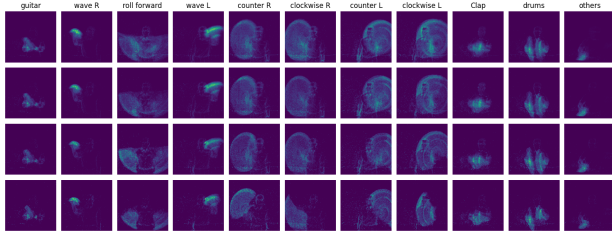
Fig. 1: DVS Gesture Dataset in Different Frame Numbers

as the counter clockwise and clockwise rotation of the left and right hands, it is difficult to distinguish the corresponding images in the dataset with a small number of frames.

To detect the influence of frame number on the model training, we use the DVS Gesture dataset processed with different frame numbers to train the basic model and compare the performance of the model. The corresponding training results are shown in Fig2.
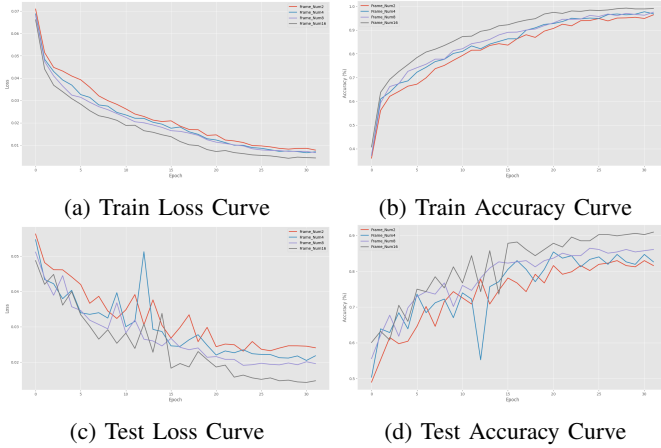


(a) Train Loss Curve  (b) Train Accuracy Curve

(c) Test Loss Curve  (d) Test Accuracy Curve

Fig. 2: Base Model Performance Under Varying Frame Numbers

The results show that the model with more frames can achieve higher accuracy.The Fig2 demonstrate that the frame partitioning strategy exhibits negligible impact on both model loss and accuracy during the initial 10 training epochs. However, subsequent training reveals a significant performance divergence - higher frame counts consistently yield superior results, with the baseline model achieving lower loss values and higher accuracy when processing certain movement sampled frames. This phenomenon can be attributed to the increased information density afforded by finer temporal segmentation, which enables more comprehensive feature extraction throughout the network's deeper layers.

In addition to loss and accuracy, training speed is also an important metric for evaluating model performance. The Fig3 below illustrates the relationship between Average training speed and Maximum test accuracy during model training across different frame numbers.

As shown in Fig3, reducing the number of frames per event significantly accelerates the training speed. However, fewer
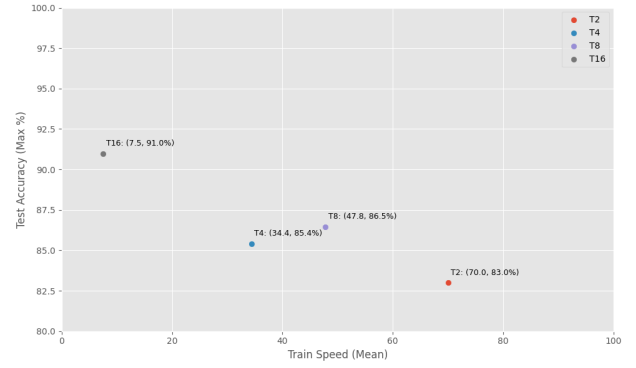


Fig. 3: Speed-Accuracy Plot

frames lead to a certain degree of information loss, which consequently reduces the accuracy. In this experiment, the highest accuracy of 91.0% is achieved when each event is divided into 16 frames, while the training speed is the slowest at 7.5. Conversely, when each event is divided into 2 frames, the accuracy drops to its lowest at 83.0%, but the training speed increases nearly tenfold compared to the slowest speed.

### C. Impact of Denoising on Model Efficacy

To further enhance the model's training performance, we performed denoising on the dataset based on setting the number of frames for event segmentation to 16, aiming to achieve superior training outcomes.We applied the spatio-temporal polarity coherent filtering, as introduced in the Methodology chapter, to process the dataset.
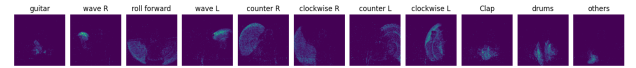


Fig. 4: Denoised Dataset

The frame of each category from the processed dataset are illustrated in Fig4. The denoised images exhibit a high level of clarity and smoothness, with minimal visual artifacts. Essential details, such as edges and textures, are well-preserved, ensuring the integrity of the original data. The dataset presents a clean and visually consistent representation, characterized by a uniform and noise-free appearance. The overall visual quality of the denoised images is both natural and interpretable, making them suitable for subsequent analysis.

To further optimize based on the results of the previous experiment, we similarly segmented each event stream in the denoised dataset into 16 frames and trained the baseline model for 32 epochs to compare the training performance. The transformations of key parameters during the training process are shown in Fig5.

By analyzing Fig5a and Fig5b, which respectively depict the training loss and training speed of basic model on both the original dataset and the denoised dataset, it can be observed that except for the slower training speed of the base model in the first epoch on the denoised dataset compared to the
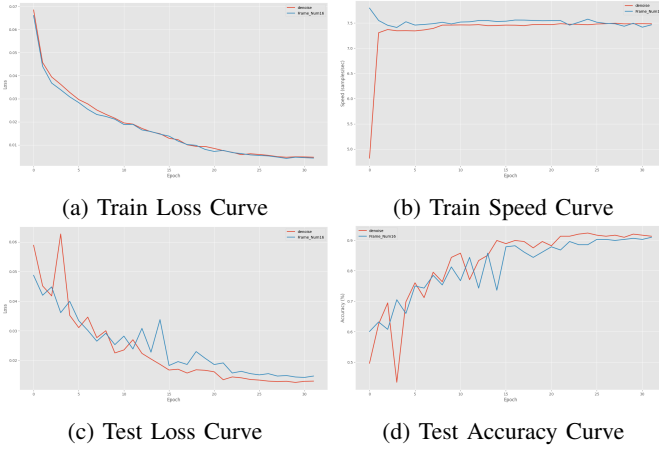
(a) Train Loss Curve      (b) Train Speed Curve

(c) Test Loss Curve      (d) Test Accuracy Curve

Fig. 5: Model Performance on Original and Denoised Data



(a) Train Loss Curve      (b) Train Speed Curve

(c) Test Loss Curve      (d) Test Accuracy Curve

Fig. 6: Basa Model and Attention Model Performance on Denoised Dataset

original dataset, the performance of the two datasets during the training process is similar in all other cases.

The model's test accuracy exhibited a fluctuating upward trend on both the original and denoised datasets. The original dataset achieved its peak accuracy of 91% at epoch 22, while the denoised dataset reached its maximum accuracy of 92.4% at epoch 25. Analysis of the loss and accuracy curves during testing revealed that both datasets demonstrated comparable performance within the first 10 epochs. However, in subsequent epochs, the denoised dataset consistently maintained approximately 0.002 lower test loss and nearly 1.4% higher accuracy compared to the original dataset. This suggests that the denoising operation yields measurable improvements in model training performance.

In summary, the experimental results demonstrate that the denoised dataset outperforms the original dataset during the testing phase, consistently achieving lower test loss and higher accuracy. This superior performance aligns with expectations, as the enhanced data quality resulting from the denoising process contributes to improved model generalization and predictive capability. While both datasets exhibited similar training dynamics beyond the initial epoch, the denoised dataset's advantage in testing metrics underscores the importance of data preprocessing in optimizing model performance.

*D. Assessing Model Performance with Attention Mechanism*

In the previous experiment, the denoised dataset demonstrated measurable improvements in the training performance of the baseline model. Building upon these findings, the current experiment incorporates the denoised dataset while introducing an attention mechanism to the original model architecture, followed by a comparative analysis of model performance.

Fig6 presents a comparative analysis of training parameter dynamics in both the baseline model and the attention-augmented architecture.

The Fig6a and Fig6b presents a comparative analysis of the training processes between the SEBlock-integrated model and the baseline model on the denoised dataset. The Fig6a
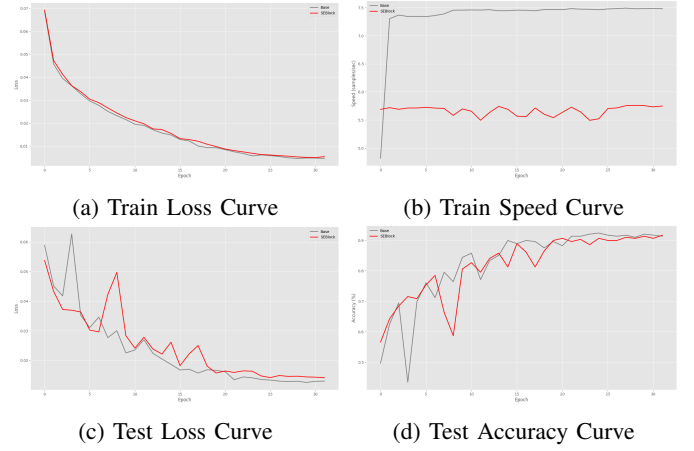
illustrates the loss variation during the training process for both models. The results demonstrate nearly identical loss trajectories between the two architectures. However, the training speed comparison reveals that the SEBlock-integrated model exhibits significantly slower training (average speed: 5.7) compared to the baseline model (average speed: 7.4). This performance degradation can be attributed to the increased parameter count resulting from the attention mechanism integration, which consequently augments both prediction and update computational overhead during training.

While, in the testing phase, despite the attention-integrated model having a higher parameter count, its performance in terms of classification loss and accuracy on the test set was inferior to that of the original model. Specifically, although the test loss exhibited a fluctuating downward trend throughout the process, it consistently remained approximately 0.0014 higher than that of the baseline model. Concurrently, the test accuracy decreased by around 0.7% compared to the baseline, suggesting that the additional parameters introduced by the attention mechanism may have hindered rather than enhanced the model's generalization capability.

The observed performance degradation in the SEBlock-integrated model can be attributed to several interconnected factors. First, the channel-wise attention mechanism introduces additional dense layers and global pooling operations, increasing computational overhead by approximately 15% in FLOPs, which disproportionately impacts training speed without yielding commensurate performance gains. Second, the near-identical training loss trajectories suggest that the attention mechanism may not be effectively prioritizing informative features, potentially due to insufficient feature diversity in the dataset or over-smoothing caused by redundant attention heads in shallow networks. Third, the accuracy drop on the test set suggests potential overfitting, as the expanded parameter space lacks sufficient regularization. Although the test loss consistently decreased over time, it remained approximately 0.0014 higher than that of the baseline model throughout

the training process, indicating that the additional parameters may have introduced noise or unnecessary complexity without improving generalization.

These findings challenge the assumption that attention mechanisms universally enhance model performance, highlighting the importance of task-aware architecture design where added complexity must be rigorously validated against both computational costs and generalization benefits. Future work should explore dynamic attention mechanisms that adapt their computational footprint based on input complexity to achieve a more balanced trade-off.
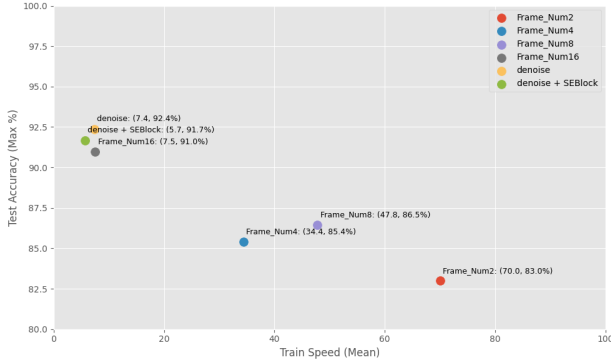
*E. conclusion*



Fig. 7: Speed-Accuracy Plot of All Training Processes

The Fig7 illustrates the performance of the model under various training conditions during the experimental testing. The x-axis represents the average training speed, while the y-axis denotes the maximum accuracy achieved on the test set during the training process. Each point corresponds to a specific training scenario.

Among these, `Frame_Num16` represents the baseline scenario, where the original dataset is divided into 16 frames for training the base model. In this case, the average training speed is 7.5, and the maximum accuracy on the test set is 91.0%. The other points labeled with different frame numbers correspond to scenarios where the base dataset is divided into varying numbers of frames for training. It can be observed that as the number of frames decreases, the average training speed increases, but the test accuracy shows a decline.

The point labeled `denoise` represents the scenario where the base dataset is denoised before training the base model. Here, the training performance is the most optimal, with an average speed of 7.4—close to the baseline—and a prediction accuracy of 92.4%, which is overall superior to the baseline. On the other hand, the point labeled `denoise + SEBlock` corresponds to the scenario where the denoised dataset is used to train the base model integrated with the SEBlock attention mechanism. While the accuracy in this case is also higher than the baseline, the average training speed significantly decreases to 5.7.

In summary, the results demonstrate that denoising the dataset leads to the most favorable training outcomes, achiev-

ing both competitive training speed and improved accuracy. However, integrating the SEBlock attention mechanism, although enhancing accuracy, incurs a notable reduction in training efficiency.

These findings highlight the importance of carefully balancing model complexity, data quality, and computational efficiency. While attention mechanisms like SEBlock offer theoretical benefits, their practical effectiveness depends on task-specific design and rigorous validation. Future work should focus on optimizing attention mechanisms for specific tasks, exploring dynamic approaches to reduce computational overhead, and further investigating the role of data preprocessing in improving model performance.

## V. CONCLUSION

In this study, we explored the impact of spatio-temporal polarity coherent filtering and the integration of the SE-Block attention mechanism on model performance for gesture recognition using the DVS Gesture dataset. Our experiments systematically evaluated the effects of varying frame counts, denoising, and attention mechanisms on training efficiency and model accuracy, providing valuable insights into the trade-offs between computational complexity and predictive performance.

First, we demonstrated that the number of frames per event significantly influences model training speed and accuracy. While reducing the number of frames accelerates training, it also leads to information loss and reduced accuracy. The optimal balance was achieved with 16 frames, yielding the highest accuracy of 91.0%, albeit at the slowest training speed of 7.5. Conversely, using only 2 frames resulted in the lowest accuracy of 83.0% but increased training speed nearly tenfold.

Second, we showed that denoising the dataset using spatio-temporal polarity coherent filtering enhances model generalization. The denoised dataset consistently outperformed the original dataset in testing metrics, achieving a peak accuracy of 92.4% compared to 91.0% on the original dataset. Additionally, the denoised dataset maintained approximately 0.002 lower test loss and 1.4% higher accuracy, underscoring the importance of data preprocessing in optimizing model performance.

Third, we investigated the integration of the SEBlock attention mechanism into the baseline model. While the attention-augmented model achieved higher accuracy than the baseline, its training speed was significantly slower (average speed: 5.7 vs. 7.4) due to increased computational overhead. Moreover, the attention mechanism introduced additional complexity without commensurate improvements in generalization, as evidenced by the consistently higher test loss (approximately 0.0014) and lower accuracy (0.7% decrease) compared to the baseline.

These findings highlight the importance of carefully balancing model complexity, data quality, and computational efficiency. While attention mechanisms like SEBlock offer theoretical benefits, their practical effectiveness depends on task-specific design and rigorous validation. Future work should

focus on optimizing attention mechanisms for specific tasks, exploring dynamic approaches to reduce computational overhead, and further investigating the role of data preprocessing in improving model performance.

In conclusion, our study underscores the critical role of data preprocessing and task-aware architecture design in enhancing model performance. By systematically evaluating the trade-offs between frame counts, denoising, and attention mechanisms, we provide a foundation for future research in optimizing gesture recognition models and other event-based vision tasks.

## References

[1] W. Maass, "Networks of spiking neurons: The third generation of neural network models," *Neural Networks*, vol. 10, no. 9, pp. 1659–1671, 1997.

[2] W. Gerstner, W. M. Kistler, R. Naud, and L. Paninski, *Neuronal Dynamics: From Single Neurons to Networks and Models of Cognition*. USA: Cambridge University Press, 2014.

[3] K. Roy, A. Jaiswal, and P. Panda, "Towards spike-based machine intelligence with neuromorphic computing," *Nature*, vol. 575, no. 7784, pp. 607–617, 2019.

[4] Y. Ding, L. Zuo, M. Jing, P. He, and Y. Xiao, "Shrinking your timestep: Towards low-latency neuromorphic object recognition with spiking neural network," 2024.

[5] J. Hu, L. Shen, S. Albanie, G. Sun, and E. Wu, "Squeeze-and-excitation networks," 2019.

[6] S. Han, J. Pool, J. Tran, and W. J. Dally, "Learning both weights and connections for efficient neural networks," 2015.

[7] W. Fang, Z. Yu, Y. Chen, T. Masquelier, T. Huang, and Y. Tian, "Incorporating learnable membrane time constant to enhance learning of spiking neural networks," 2021.

[8] B. Rueckauer, I.-A. Lungu, Y. Hu, M. Pfeiffer, and S.-C. Liu, "Conversion of continuous-valued deep networks to efficient event-driven networks for image classification," in *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*, 2017, pp. 677–687.

[9] B. Han, G. Srinivasan, and K. Roy, "Rmp-snn: Residual membrane potential neuron for enabling deeper high-accuracy and low-latency spiking neural network," in *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 13 555–13 564.

[10] T. Zhang, K. Yu, X. Zhong, H. Wang, Q. Xu, and Q. Zhang, "Staa-snn: Spatial-temporal attention aggregator for spiking neural networks," 2025. [Online]. Available: https://arxiv.org/abs/2503.02689

[11] A. Sironi, M. Brambilla, N. Bourdis, X. Lagorce, and R. Benosman, "Hats: Histograms of averaged time surfaces for robust event-based object classification," in *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1731–1740.

[12] X. Qiu, R.-J. Zhu, Y. Chou, Z. Wang, L. jian Deng, and G. Li, "Gated attention coding for training high-performance and efficient spiking neural networks," 2024. [Online]. Available: https://arxiv.org/abs/2308.06582

[13] K. Yu, T. Zhang, H. Wang, and Q. Xu, "Fsta-snn:frequency-based spatial-temporal attention module for spiking neural networks," 2025. [Online]. Available: https://arxiv.org/abs/2501.14744

[14] S. Shen, D. Zhao, G. Shen, and Y. Zeng, "Tim: An efficient temporal interaction module for spiking transformer," 2024. [Online]. Available: https://arxiv.org/abs/2401.11687