

# SFOD: Spiking Fusion Object Detector

Yimeng Fan<sup>1</sup>, Wei Zhang<sup>1,2,\*</sup>, Changsong Liu<sup>1</sup>, Mingyang Li<sup>1</sup>, Wenrui Lu<sup>1</sup>

<sup>1</sup>School of Microelectronics, Tianjin University, China

<sup>2</sup>Tianjin Key Laboratory of Low-dimensional Electronic Materials and Advanced Instrumentation



## Introduction

Motivation: Why do we need Spiking Fusion in SNNs?

- The combination of deeper and shallower feature maps in the spatial domain
- Enhancing connections between features of different scales in the temporal domain

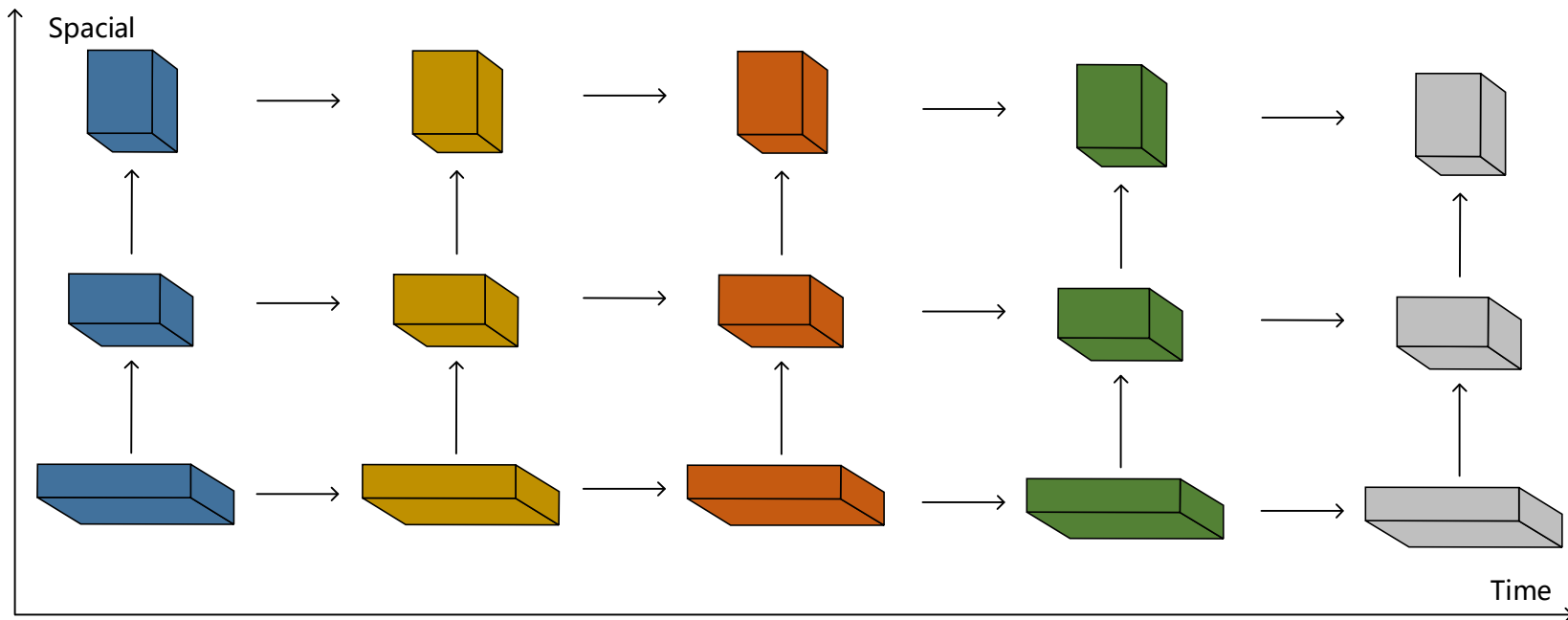


Fig1. Spatiotemporal feature extraction

# Introduction

## Related Work: Lack of corresponding research on Spiking Fusion

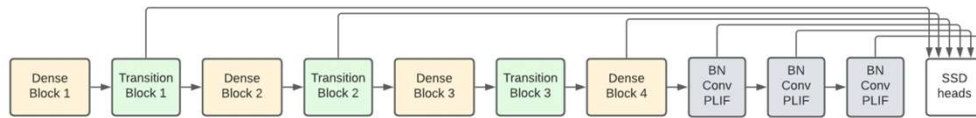


Fig2. Overview of Spiking DenseNet + SSD architecture [1].

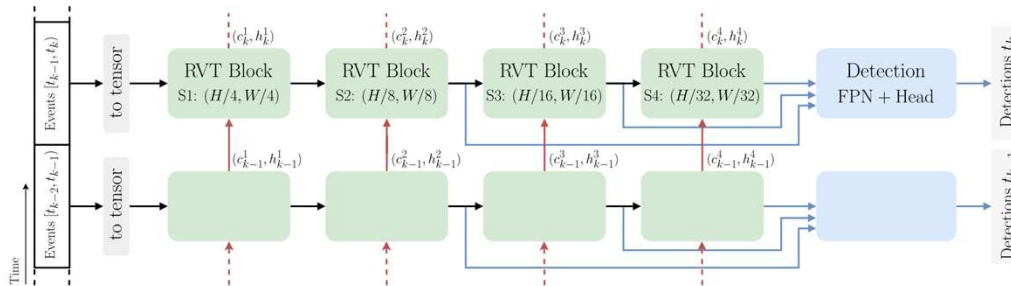


Fig3. Overview of RVT model [2].

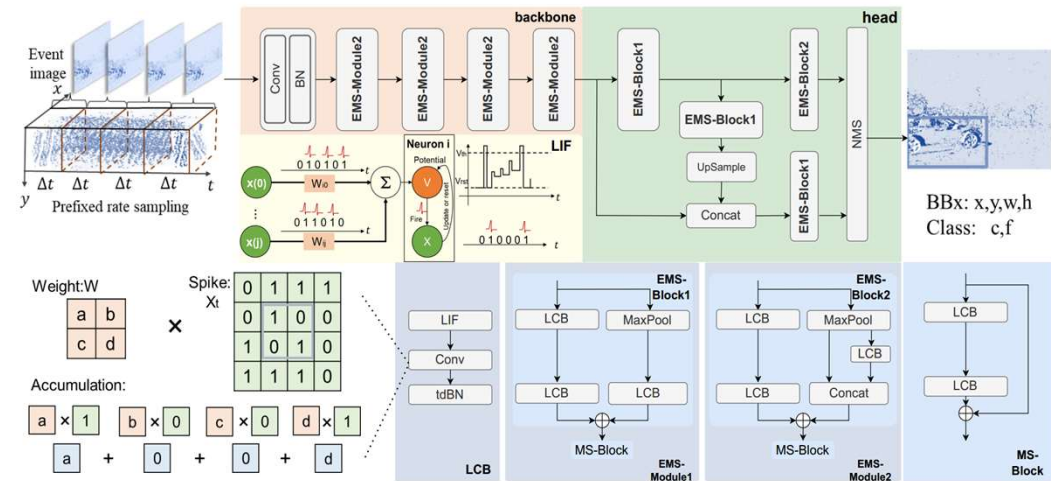


Fig4. Overview of EMS-YOLO model [3].

[1] Cordone L, Miramond B, Thierion P. Object detection with spiking neural networks on automotive event data[C]//2022 International Joint Conference on Neural Networks (IJCNN). IEEE, 2022: 1-8.

[2] Gehrig M, Scaramuzza D. Recurrent vision transformers for object detection with event cameras[C]//Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2023: 13884-13893.

[3] Su Q, Chou Y, Hu Y, et al. Deep directly-trained spiking neural networks for object detection[C]//Proceedings of the IEEE/CVF International Conference on Computer Vision. 2023: 6555-6565.

# Method

## Simple Fusion Model: SFOD(Spiking Fusion Object Detector)

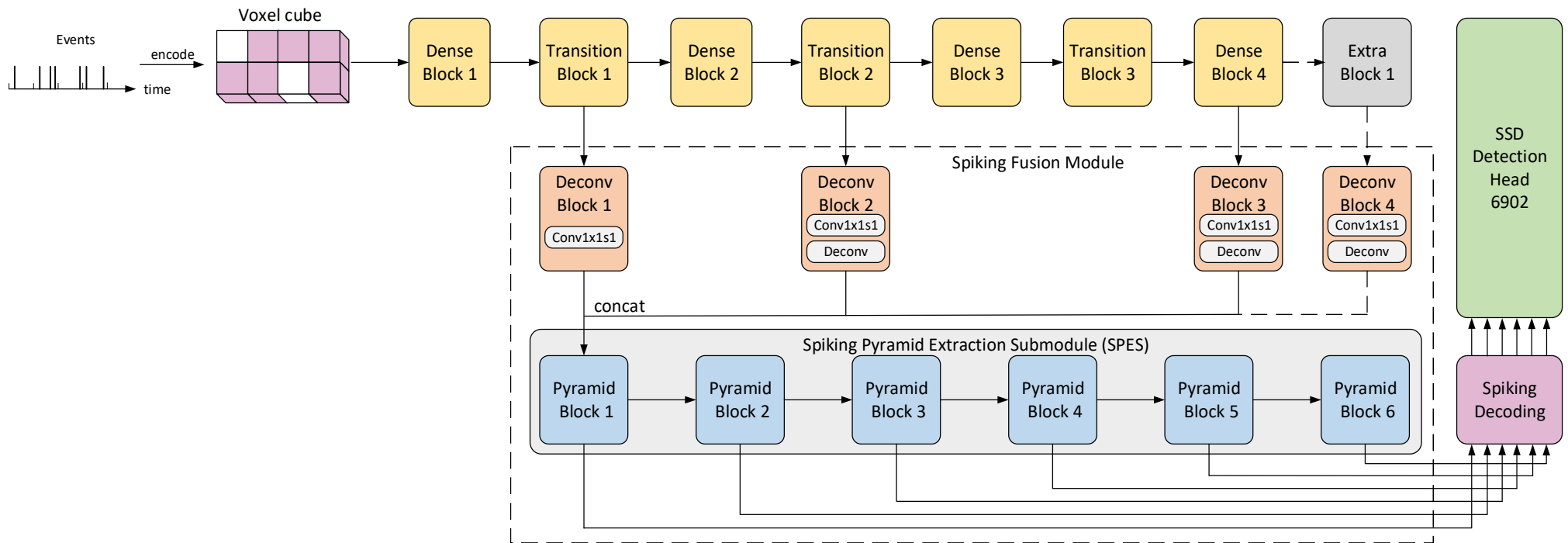


Fig5. Overview of SFOD.

## Method

### The architectures of SPES

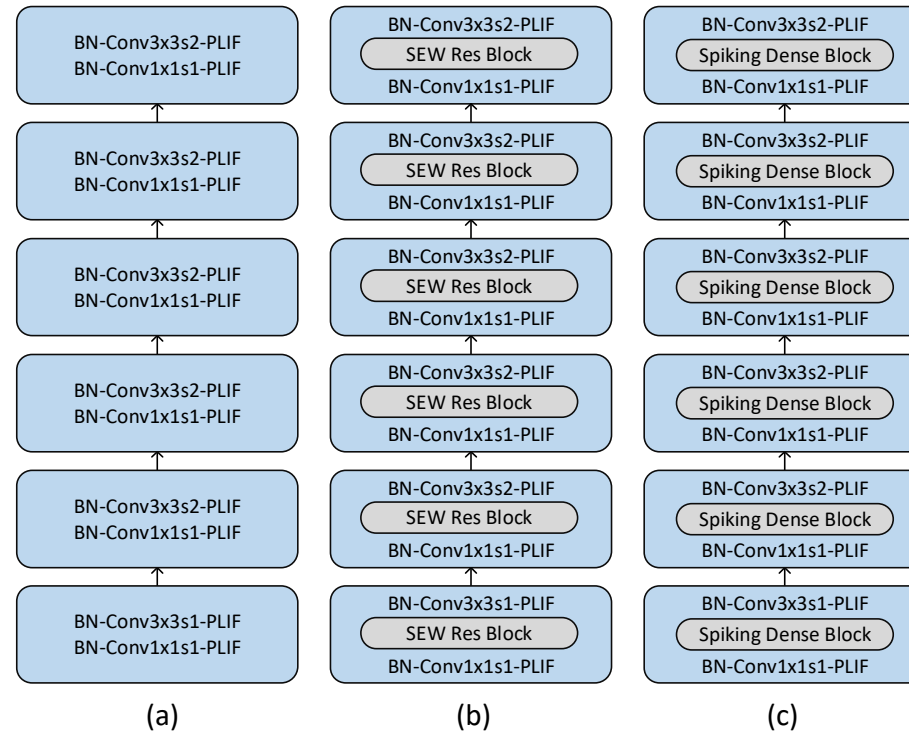


Fig6. The architecture of SPES.

## Experiments



### Ablation Results:

| Models             | Dec. |       | Fusion Layers |   |   | Params       | mAP@0.5:0.95 | mAP@0.5      | Firing Rate   |
|--------------------|------|-------|---------------|---|---|--------------|--------------|--------------|---------------|
|                    | Rate | Count | None          | 3 | 4 |              |              |              |               |
| DenseNet121-16-SSD | ✓    |       | ✓             |   |   | 5.0M         | 0.262        | 0.517        | 21.01%        |
| DenseNet121_24-SSD |      | ✓     | ✓             |   |   | 8.2M         | 0.235        | 0.445        | 22.02%        |
| DenseNet121-24-SSD | ✓    |       | ✓             |   |   | 8.2M         | 0.288        | 0.553        | 22.29%        |
| DenseNet169-16-SSD | ✓    |       | ✓             |   |   | 7.7M         | 0.257        | 0.507        | 22.82%        |
| SFOD-B             | ✓    |       |               |   | ✓ | 15.0M        | 0.294        | 0.570        | 21.13%        |
| SFOD-B             | ✓    |       |               | ✓ |   | 9.9M         | 0.299        | 0.575        | 24.41%        |
| SFOD-D             | ✓    |       |               | ✓ |   | 11.3M        | 0.286        | 0.558        | 26.37%        |
| <b>SFOD-R</b>      | ✓    |       |               | ✓ |   | <b>11.9M</b> | <b>0.321</b> | <b>0.593</b> | <b>24.04%</b> |

Tab1. Results of the ablation study on the GEN1 dataset. We first study the performance differences across object detection models using various backbone networks. Based on this, we select the best backbone and further analyze the impact of different fusion layers. Finally, we compare the performance of various SPES variants. We name the models using the basic, Spiking Dense Block-enhanced, and SEW Res Block-enhanced SPESs as SFOD-B, SFOD-D, and SFOD-R, respectively.

## Experiments



### Benchmark Comparisons:

| Method                 | Networks                | Detection Head    | Params       | mAP<br>@0.5:0.95 | Firing<br>Rate | Time<br>(ms)     | Energy<br>(mJ) |
|------------------------|-------------------------|-------------------|--------------|------------------|----------------|------------------|----------------|
| Asynet [30]            | Sparse CNNs             | YOLOv1 [34]       | 11.4M        | 0.145            | -              | -                | > 4.83         |
| AEGNN [39]             | GNNs                    | YOLOv1            | 20.0M        | 0.163            | -              | -                | -              |
| Inception+SSD [17]     | CNNs                    | SSD [25]          | -            | 0.301            | -              | 19.4             | -              |
| MatrixLSTM [5]         | RNNs+CNNs               | YOLOv3 [33]       | 61.5M        | 0.310            | -              | -                | -              |
| RED [32]               | RNNs+CNNs               | SSD               | 24.1M        | 0.400            | -              | 16.7             | > 24.08        |
| <b>RVT [15]</b>        | <b>Transformer+RNNs</b> | <b>YOLOX [14]</b> | <b>18.5M</b> | <b>0.472</b>     | -              | <b>10.2</b>      | -              |
| MobileNet-64+SSD [9]   | SNNs                    | SSD               | 24.3M        | 0.147            | 29.44%         | 1.7 <sup>†</sup> | 5.76           |
| VGG-11+SDD [9]         | SNNs                    | SSD               | 12.6M        | 0.174            | 22.22%         | 4.4 <sup>†</sup> | 11.06          |
| DenseNet121-24+SSD [9] | SNNs                    | SSD               | 8.2M         | 0.189            | 37.20%         | 4.1 <sup>†</sup> | 3.89           |
| <b>EMS-YOLO [41]</b>   | <b>SNNs</b>             | <b>YOLOv3</b>     | <b>14.4M</b> | <b>0.310</b>     | <b>17.80%</b>  | -                | -              |
| <b>SFOD</b>            | <b>SNNs</b>             | <b>SSD</b>        | <b>11.9M</b> | <b>0.321</b>     | <b>24.04%</b>  | <b>6.7</b>       | <b>7.26</b>    |

Tab2. Comparison with state-of-the-art models on the GEN1 dataset. We present a comparison of our model with other state-of-the-art approaches on the GEN1 dataset. Remarkably, our model achieves a state-of-the-art mAP of 32.1% at the same level of firing rate and parameters compared to other SNN-based methods.



## Experiments

Inference results of the model on the GEN1 dataset:

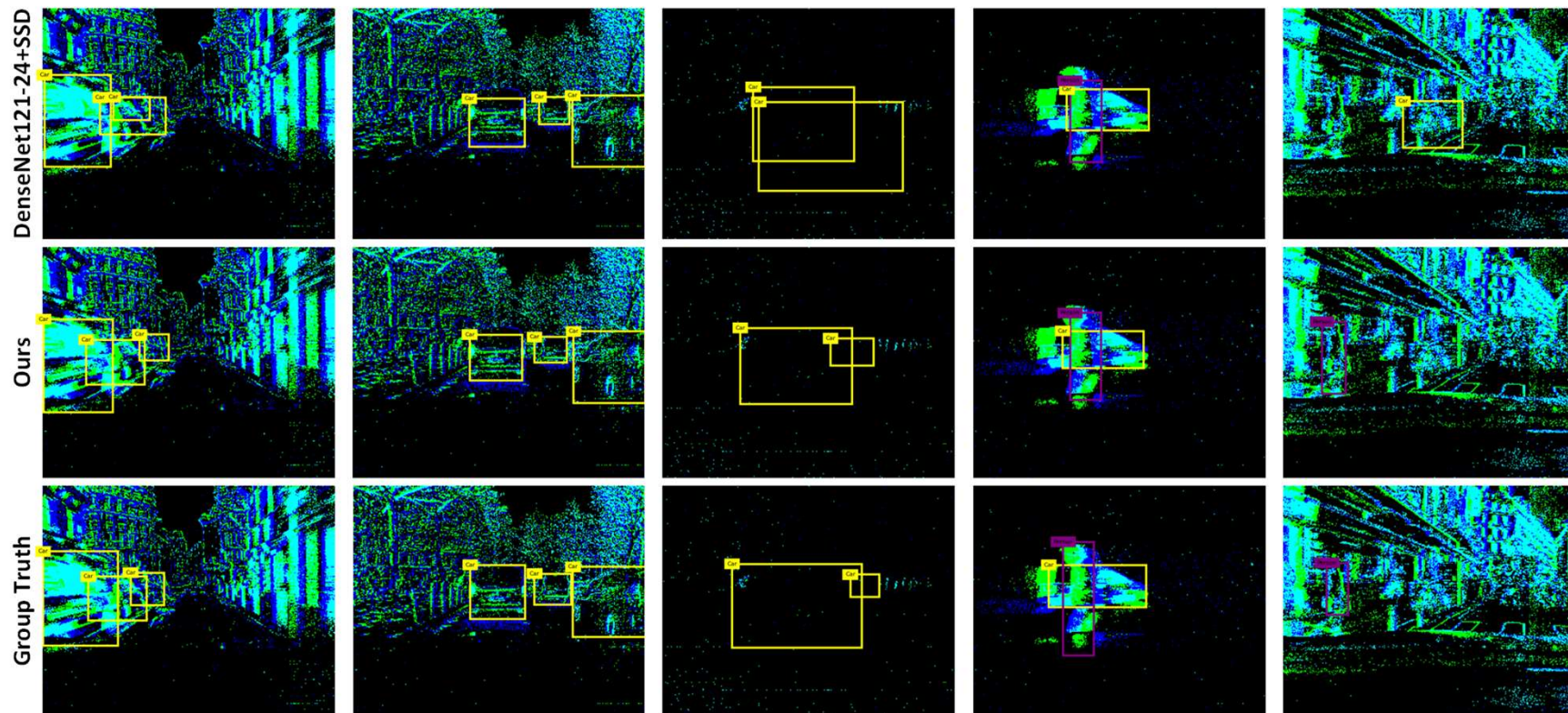


Fig6. Inference results of the model on the GEN1 dataset.



## Conclusion



### Summary and Future Outlook:

- **We propose Spiking Fusion Module, which is the first to implement spiking feature fusion in SNNs for event cameras.**
- **For the first time in SNNs applied to event cameras, we conduct a thorough study of different spiking decoding strategies and classification loss functions to determine their impact on model performance.**
- **On the GEN1 dataset, our SFOD achieves the state-of-the-art object detection performance of 32.1% mAP for SNN-based models.**
- **In the future, we believe that the performance of SFOD is expected to be further improved by adopting a more effective data augmentation strategy. It undeniably represents a promising research direction.**