

SPIKING TRANSFORMER-CNN FOR EVENT-BASED OBJECT DETECTION

Anonymous authors

Paper under double-blind review

ABSTRACT

Spiking Neural Networks (SNNs) enable energy-efficient computation through event-driven computing and multiplication-free inference, making them well-suited for processing sparse events. Recently, deep Spiking Convolutional Neural Networks (CNNs) have shown energy efficiency advantages on event-based object detection. However, spiking CNNs have been limited to local and single-scale features, making it challenging for them to achieve better detection accuracy. To address this challenge, we propose a hierarchical Spiking Transformer-CNN (i.e., Spike-TransCNN) architecture, which is the first attempt to leverage the global information extraction capabilities of Spiking Transformers and the local information capture abilities of Spiking CNNs for event-based object detection. Technically, we first propose using the Spiking Transformer to extract global features and employ a multi-scale local feature extraction CNN module to complement the Spiking Transformers in local feature extraction. Then, we design intra-stage and inter-stage feature fusion modules to integrate global and multi-scale local features within the network architecture. Experimental results demonstrate that our Spike-TransCNN significantly outperforms existing SNN-based object detectors on the Gen1 dataset, achieving higher detection accuracy (mAP 0.336 vs. 0.321) with lower energy consumption (5.49 mJ vs. 7.26 mJ). Our code can be available in the supplementary materials.

1 INTRODUCTION

Object detection is essential in computer vision and robotics applications. Nevertheless, conventional cameras operating at fixed frame rates struggle in challenging conditions, such as fast motion, over-exposure, and low light, leading to a significant decline in object detection performance Liu et al. (2020); Sayed & Brostow (2021). Recently, event cameras like DVS Lichtsteiner et al. (2008) and ATIS Posch et al. (2010) have surpassed RGB cameras in dynamic range, temporal resolution, and energy efficiency. These advanced capabilities make them especially well-suited for object detection Peng et al. (2023a;b); Wang et al. (2023a; 2024); Zubic et al. (2024); Yuan et al. (2024); Hamaguchi et al. (2023); Zubić et al. (2023); Tomy et al. (2022) in these challenging scenarios.

Most current event-based object detectors rely on Artificial Neural Networks (ANNs) Perot et al. (2020); Li et al. (2022a), which deliver high performance but come with high computational complexity and energy consumption. In contrast, Spiking Neural Networks (SNNs) Maass (1997); Zhu et al. (2022) present a novel approach inspired by the brain’s temporal information processing dynamics. SNNs propagate information through binary spike sequences, allowing for energy-effective computing with event-driven computation and multiplication-free inference. This makes SNNs a more efficient and biologically inspired alternative for event-based object detection.

Early SNN-based object detectors are often derived from existing ANNs through conversion, which introduces several limitations. The most critical problem is that most conversion methods are designed for static images and may not be able to handle sparse temporal event data effectively. This is because these methods focus on approximating the activation patterns of ANNs, often neglecting the spatiotemporal information inherent in event data. Furthermore, converted models like Spiking-YOLO Kim et al. (2020) require a large number of time steps to match the performance of the original ANN. Although Spike Calibration Li et al. (2022b) can reduce this to hundreds of time steps, its effectiveness still hinges on the quality of the original ANN model.

054 Directly-trained Spiking Convolutional Neural Networks (CNNs) could be trained with much fewer
055 steps Su et al. (2023), but they primarily focus on local features, which can limit the overall detec-
056 tion performance. For instance, Spiking-DenseNet Cordone et al. (2022) employs a multi-layered
057 approach to process features at various local scales, and SFOD Fan et al. (2024) introduces a fusion
058 mechanism to combine spike features across different local scales. Despite these efforts to opti-
059 mize spike features across multiple local scales to capture the local features, these Spiking CNNs
060 still face challenges in incorporating global and high-level semantic information, constraining their
061 overall performance.

062 Transformer architectures have increasingly been integrated into SNNs, such as Spikformer Zhou
063 et al. (2023), Auto-Spikformer Che et al. (2024), and Attention-free Spikformer Wang et al. (2023b).
064 These models have demonstrated superior performance over spiking CNNs in various tasks, primar-
065 ily due to the Transformer’s capability for global attention and parallel computation. However, most
066 current research focuses on classification, with limited exploration of Spiking Transformers in the
067 regression task of object detection. Additionally, recent ANN studies Fang et al. (2022); Chen et al.
068 (2023) suggest that combining the global modeling strengths of Transformers with the local fea-
069 ture extraction capabilities of CNNs can further enhance network performance. Nevertheless, this
070 potential has yet to be thoroughly investigated within the context of SNNs.

071 To address these challenges, we propose a hierarchical Spiking Transformer-CNN (i.e., Spike-
072 TransCNN), which is a directly-trained deep SNN designed to extract both global and multi-scale
073 local features for event-based object detection. Our model is the first attempt to leverage the global
074 information extraction capabilities of Spiking Transformers and the local information capture abil-
075 ities of Spiking CNNs. We employ spike-driven token selection to selectively capture tokens, and
076 spike self-attention for holistic perception of spike features. Specifically, we first present a multi-
077 scale local feature extraction module to compensate for the limitations of Spiking Transformers in
078 local feature extraction. Furthermore, we design intra-stage and inter-stage feature fusion modules
079 to integrate global and multi-scale local features within the architecture. The results show that our
080 Spike-TransCNN reduces energy consumption by 4.7× compared to the same ANN architecture.
081 Moreover, it significantly outperforms state-of-the-art methods on the Gen1 dataset, achieving a
082 higher mAP and lower energy consumption.

083 In summary, the main contributions of this work are:

- 084 • We propose Spike-TransCNN, a novel *hierarchical Spiking Transformer-CNN*, which is the
085 first hybrid spiking architecture that combines Spiking CNNs and Spiking Transformers to
086 leverage their complementary strengths, yielding both high-accuracy and energy-efficiency
087 in event-based object detection.
- 088 • We present *spike-driven token selection* to select tokens and spike self-attention for global
089 spike feature perception, along with *spiking dilated convolution* for extracting local multi-
090 scale features and optimizing them with *temporal-channel joint attention*.
- 091 • We design *intra-stage spike feature fusion and inter-stage spike feature fusion modules* that
092 effectively aggregate features extracted from different architectures and multi-scale features
093 from the event stream to improve object detection performance.

096 2 RELATED WORK

098 **Event-based Object Detection.** Most event-based object detection methods use ANN approaches,
099 such as RED Perot et al. (2020), ASTMNet Li et al. (2022a), and RVT Gehrig (2023), which demon-
100 strate impressive detection performance but come with high energy consumption. More recently,
101 some works have explored achieving energy-efficient object detection for event data using SNNs.
102 For example, Hybrid-SNN Kugele et al. (2021) combines an SNN backbone for efficient event-
103 based feature extraction with an ANN head for object detection tasks. Spiking-DenseNet Li et al.
104 (2022a) is notable for applying SNNs to event-based object detection using the SSD architecture. A
105 feature pyramid structure Zhang et al. (2023) is introduced to support multi-scale feature extraction.
106 SFOD Fan et al. (2024) introduces a fusion mechanism to combine spike features across different
107 scales. Despite these efforts, Spiking CNNs excel at capturing local features but face challenges in
integrating global information, which limits overall detection performance.

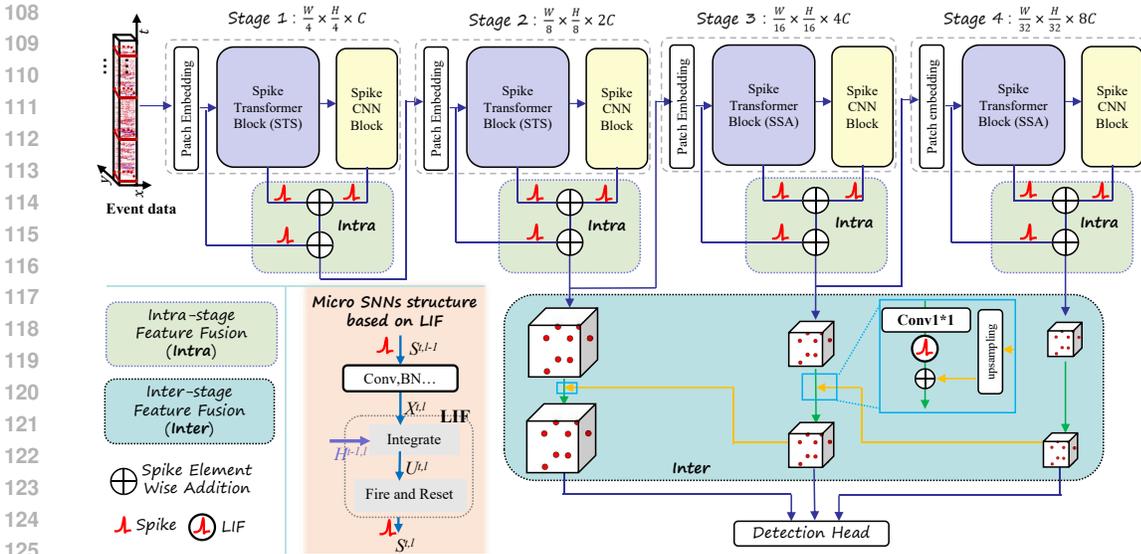


Figure 1: The pipeline of our hierarchical Spiking Transformer-CNNs. Initially, the event stream is processed through a hierarchical hybrid backbone that integrates Spiking Transformer and Spiking CNN blocks. Then, we apply the patch embedding operation across four stages to extract features at four different scales and integrate these features using intra-stage and inter-stage feature fusion modules. Finally, the detection results are predicted on the fused features using the YOLOX head.

Spiking CNNs. Current Spiking CNN-based models for object detection can be broadly categorized into two types. The first type involves ANN-to-SNN methods, which convert pre-trained ANNs into SNNs by replacing continuous activation functions with spiking neurons. For example, Spiking-YOLOv4 Wang et al. (2023c) implements a converted CNN model for fast and accurate object detection from event streams. However, these methods face several limitations, including the need for a large number of time steps to match the performance of the original ANN. The second type refers to directly-trained methods, which use surrogate gradients to train deep and large-scale SNNs for object detection. For instance, EMS-YOLO Su et al. (2023) is a directly-trained SNN that surpasses ANN-to-SNN conversion methods, requiring only a few time steps for real-time inference. Besides, a training scheme for SNNs Caccavella et al. (2023) deployed on the neuromorphic chip achieves low-power face detection. These directly-trained SNN models can achieve comparable performance to ANNs with the same architecture while significantly reducing energy consumption.

Spiking Transformers. Transformers have been integrated into SNN models with notable success across various tasks. For example, Spikformer Zhou et al. (2023) introduces a pioneering Spiking Self-Attention (SSA) version of the self-attention mechanism. Meanwhile, Spike-Driven Self-Attention (SDSA) Yao et al. (2024) employs mask and addition operations to avoid multiplication, significantly reducing computational energy compared to conventional self-attention mechanisms. Besides, a spatiotemporal self-attention mechanism Wang et al. (2023d) has been proposed for SNNs, effectively capturing feature dependencies while preserving the asynchronous transmission property of SNNs. QKFormer Zhou et al. (2024) utilizes sparse matrices to filter tokens and channels. Despite these advancements in Transformer-based SNN models, the exploration of object detection tasks using event data remains relatively limited.

3 METHODS

3.1 OVERVIEW

This work aims at designing a novel hierarchical Spiking Transformer-CNN, termed **Spike-TransCNN**, which combines Spiking CNNs and Spiking Transformers to leverage their complementary strengths for high-accuracy and energy-efficient object detection. As depicted in Fig. 1, our framework consists of four key parts: Spiking Transformer Block (STB), Spiking CNN Block

(SCB), intra-state spike feature fusion module, and inter-stage spike feature fusion module. More precisely, our Spike-TransCNN starts by processing the event stream through a hybrid backbone that combines Spiking Transformers and Spiking CNNs. We then extract features at four different scales using patch embeddings (PE) across four stages in the hybrid architecture, integrating these features with both intra-stage and inter-stage fusion modules. Finally, we use the YOLOX head to predict detection results from the fused features.

3.2 THE BASICS OF SPIKE NEURAL NETWORKS

Spike Neural Networks (SNNs) are inspired by the brain’s functioning, communicating through discrete spikes, enabling efficient information processing for event data. Compared to traditional ANNs, SNNs excel in spatio-temporal processing and energy efficiency. Spiking neurons are the essential components of SNNs, communicating through discrete spikes that mimic the behavior of biological brain neurons. The Leaky Integrate-and-Fire (LIF) neuron model Delorme et al. (1999) is commonly used in SNNs as follows:

$$V_t = V_{t-1} + \frac{1}{\tau}(-(V_{t-1} - V_{rest}) + X_t), \quad (1)$$

$$S_t = Hea(V_t - V_{th}), \quad (2)$$

where V_t denotes the membrane potential after neuronal dynamics at timestep t . X_t represents the input to neuron. τ is the decay factor for leakage. $Hea()$ is the Heaviside step function, which satisfies: $Hea(x) = 1$ when $x \geq 0$, otherwise $Hea(x) = 0$. The generation of output spikes is controlled by the threshold V_{th} , and once the neuron emits a spike at time step $t + 1$, the current membrane potential will be reset to V_{rest} . SNNs Wu et al. (2018) based on spike neurons (LIF) can be developed for network training as:

$$U^{t,l} = H^{t-1,l} + X^{t,l}, S^{t,l} = Hea(U^{t,l} - u_{th}), \quad (3)$$

$$H^{t,l} = V_{rest}S^{t,l} + (\beta U^{t,l}) \odot (1 - S^{t,l}), \quad (4)$$

where t and l respectively represent timestamps and network layers. V represents the membrane potential generated by integrating spatial dimension input $X^{t,l}$ and temporal dimension input $H^{t-1,l}$. u_{th} is the threshold that determines whether the output spike tensor $S^{t,l}$ should be fired or kept as zero. $H^{t,l}$ represents the internal state of neurons propagating over time, where $\beta = e^{-\frac{dt}{\tau}}$ reflects the decay factor, and \odot denotes element-wise multiplication.

3.3 SPIKING TRANSFORMERS FOR GLOBAL FEATURES

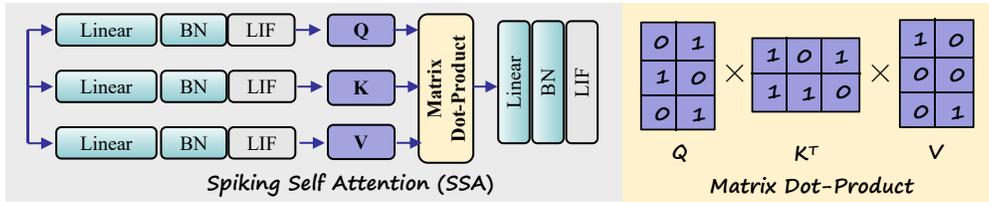


Figure 2: Spiking Self Attention (SSA). SSA is employed for global feature extraction in the latter two stages. The right illustrates the operational processes of ‘Matrix Dot-Product’.

To obtain hierarchical features, the patch embedding operation is implemented for downsampling before each STB. It specifically includes conv1 * 1, max-pooling, and LIF. We use Spike-driven Token Selection (STS) to select tokens for sparse features in the first two STB modules (see Fig. 3) and employ Spiking Self-Attention (SSA) to extract global features in the latter two STB modules (see Fig. 2). SSA utilizes sparse spike-form Q, K, V for computation, without the need for softmax operations and floating-point matrix multiplication. The calculation process of SSA can be formulated as follows:

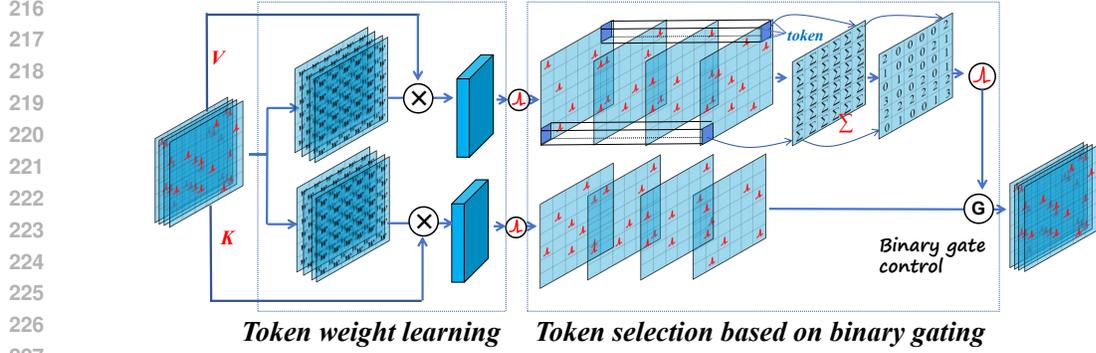


Figure 3: Spike-driven Token Selection (STS). STS is used for token selection in the first two stages.

$$I = SN_I(BN(XW_I)), I \in (Q, K, V), \quad (5)$$

$$SSA(Q, K, V) = SN(BN(Linear(SN(QK^T V * s)))), \quad (6)$$

where $Q, K, V \in R^{T \times N \times D}$, N is the token number, and D is the channel number. The spike-form Q, K, V are computed by learnable linear layers. s represents a scaling factor, and SN refers to the LIF layer.

In order to select important regions within shallow features, we designed Spike-driven Token Selection (STS) in the first two STBs. Firstly, use linear functions to learn the weights W_v and W_k of the value and key domains for each token. Then, utilize LIF to map the key values to (0 or 1) as the gating control for tokens. The STS can be formulated as follows:

$$V = SN_V(BN(XW_V)), K = SN_K(BN(XW_K)), \quad (7)$$

$$G_t = \sum_{i=0}^D SN(K_{i,j}), X' = G_t \otimes V, \quad (8)$$

where G_t is the $N * 1$ token attention vector, which models the binary importance of different tokens. D is the channel number. \otimes is the Hadamard product between spike tensors.

In our backbone network, we utilize STS for token filtering in the first two stages, and SSA in the latter two stages. Experimental results demonstrate that the combination of these two attention mechanisms may enhance the detection performance of the network.

3.4 SPIKING CNNs FOR MULTI-SCALE LOCAL FEATURES

The STB excels at extracting global information but lacks the ability to capture local details. Therefore, we designed a local multi-receptive field Spiking CNN Block (SCB) that leverages the local feature extraction capabilities of CNNs and automatically selects important local multi-scale features using time channel attention (see Fig. 4).

Initially, the SCB utilizes dilated convolutions with dilation rates of [1, 3, 5] to capture local multi-scale information. Subsequently, the multi-scale local spiking features are stacked along the channel dimension to ensure feature binarization, which can be formulated as follows:

$$X_i = SN_i(BN(DC_{dr=i}(X))), i \in (1, 3, 5), \quad (9)$$

$$X_{SDC} = Cat(X, X_1, X_3, X_5), \quad (10)$$

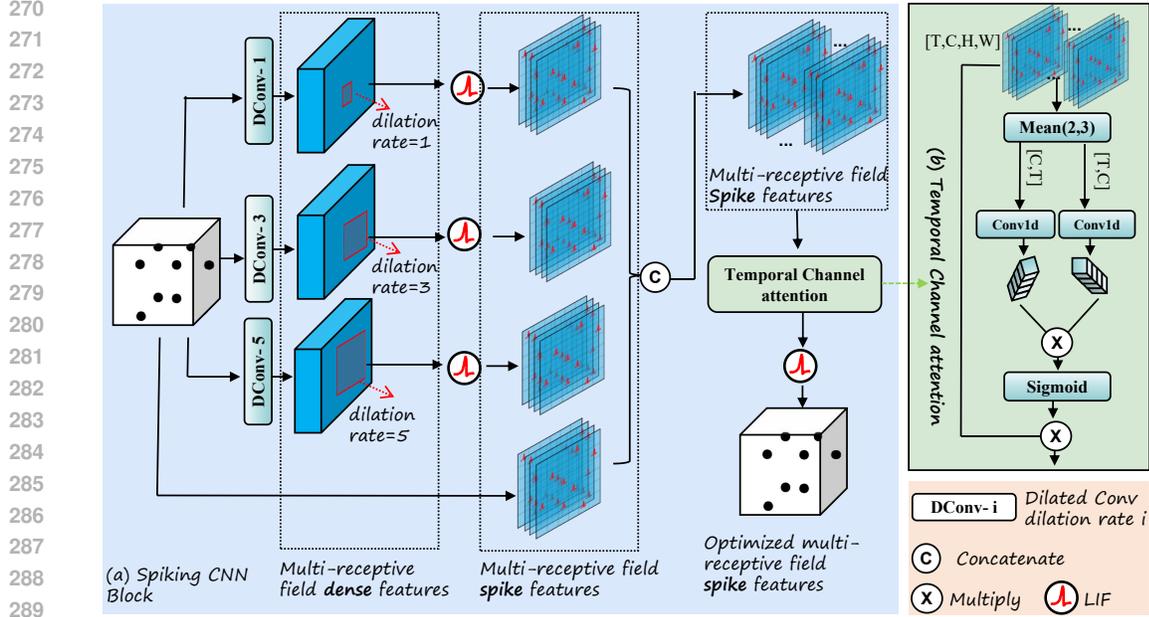


Figure 4: Spiking CNN Block. Extract local multi-receptive field features using dilated convolutions, and select features using temporal channel attention.

where $DC_{dr=i}$ denotes dilated convolution with a dilation rate of i , while Cat represents channel-wise concatenation.

For the $[T,C,H,W]$ dimensional spiking features X_{SDC} , a time channel attention mechanism is employed to automatically select the time and channel of the multi-scale local spiking features. Finally, LIF is used to activate the features, ensuring the binarization of the features input to the next stage and maintaining the binarized characteristics of the SNN network, which can be formulated as:

$$X_{SCB} = SN(X_{SDC} \odot \sigma(\text{Conv1d}(F_t) \cdot \text{Conv1d}(F_c))), \quad (11)$$

where F_t is the mean of X_{SDC} along dimensions H and W , and F_c is obtained by swapping the dimensions (T, C) of F_t .

3.5 INTRA-STAGE SPIKE FEATURE FUSION

To ensure that the features in each stage used for object detection contain both global features and local multi-receptive field features, we use spike-element-wise (SEW) addition to fuse the original input features, STB output features, and SCB output features in the intra-stage spike feature fusion module. SEW-ResNet Fang et al. (2021) has demonstrated that spike-element-wise addition in residual connections can effectively prevent issues of gradient vanishing and gradient explosion. Fig. 5 illustrates the (a) SEW block Fang et al. (2021) and (b) the proposed intra-stage spike feature fusion module. The specific operations of the intra-stage spike feature fusion module are as follows:

$$X_{Intra} = X_{PE} \oplus (X_{STB} \oplus X_{SCB}), \quad (12)$$

where X_{PE} , X_{STB} , and X_{SCB} denote the spike features output by PE, STB, and SCB, respectively. \oplus refers to the spike-element-wise addition operation, which ensures that the fusion operation in-

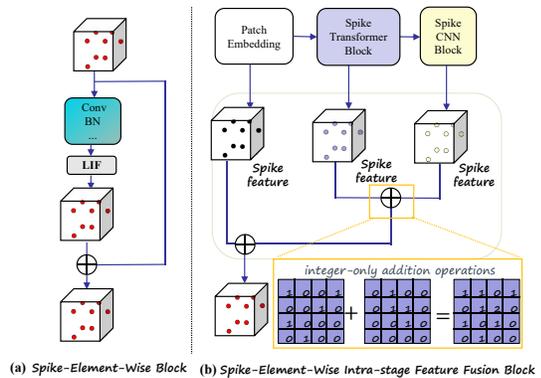


Figure 5: Intra-stage Feature Fusion. (a) SEW Block. (b) Intra-stage spike feature fusion block.

324 involves only integer addition. The intra-stage spike feature fusion not only ensures the integration of
 325 global and local features at each stage, but it also prevents the loss of features extracted by distant
 326 STBs when passing input to the next stage.
 327

3.6 INTER-STAGE SPIKE FEATURE FUSION

330 Intra-stage spike feature fusion can
 331 merge features with the same spatial
 332 resolution, but it lacks interaction
 333 between features of different stages.
 334 Therefore, we design an inter-stage
 335 spike feature fusion module to facilitate
 336 interaction between different
 337 stages. The specific design of the
 338 inter-layer feature fusion module
 339 is shown in Fig. 6. It utilizes a 1×1
 340 convolution to adjust the channel
 341 dimensions of the features at the
 342 current stage, applies up-sampling
 343 to the features from the previous
 344 stage, and then adds them element-wise.
 345 To ensure the binary and sparse nature
 346 of the features, LIF is used for spike
 347 activation after each convolution. The
 348 operations between two adjacent stages
 349 can be formulated as follows:

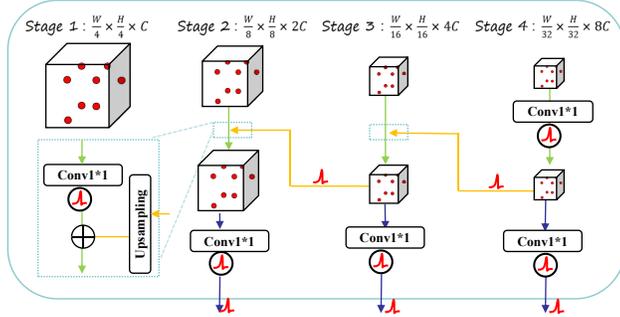


Figure 6: Inter-stage Feature Fusion. This module consists of two main operations: vertical channel modification and horizontal upsampling, enabling interaction between shallow features and deep features over longer distances.

$$X_{S-(j)} = SN(Conv(X_{Intra-(j)})), \tag{13}$$

$$X_{Inter-(j)} = X_{S-(j)} \oplus Up2(X_{S-(j-1)}), \tag{14}$$

$$X_{Intra-(j)} = SN(Conv(X_{Inter-(j)})), \tag{15}$$

357 where $X_{Intra-(j)}$ denotes the output spike features of the j -th stage after intra-stage feature fusion,
 358 and $Up2$ represents 2x up-sampling. This module enables interaction and fusion of features at
 359 different depths and resolutions.

360 Finally, the top three features obtained from the inter-stage feature fusion are fed into the YOLOX
 361 Ge et al. (2021) detection head for classification and regression predictions.
 362

4 EXPERIMENTS

4.1 EXPERIMENT SETTINGS

367 **Datasets.** We evaluate the effectiveness of our Spike-TransCNN using two publicly available anno-
 368 tated datasets: Gen1 De Tournemire et al. (2020) and 1Mpx Perot et al. (2020). Both datasets were
 369 captured by event cameras in real-world driving scenarios. The Gen1 dataset includes 39 hours of
 370 event data with a resolution of 304×240 , providing 228k car and 28k pedestrian bounding boxes
 371 labeled at frequencies of 1, 2, or 4 Hz. The 1Mpx dataset offers recordings with a resolution of 720
 372 $\times 1280$, totaling around 15 hours of event data, and includes 25 million bounding box labels at either
 373 30 or 60 Hz for cars, pedestrians, and two-wheelers. For processing these datasets, we follow the
 374 methodology outlined in RVT-B Gehrig (2023).

375 **Implementation Details.** Our models are trained for 400k iterations using the ADAM optimizer
 376 Kinga et al. (2015) with a OneCycle learning rate schedule Smith & Topin (2019), which includes
 377 2000 warmup iterations followed by linear decay of the maximum learning rate. This training strat-
 egy is consistent across all studies. For the Gen1 dataset, we use a batch size of 4 and a maximum

Table 1: Comparison with state-of-the-art methods and our Spike-TranCNN on the Gen1 dataset and the 1Mpx dataset. Note that, our Spike-TransCNN significantly outperforms existing SNN-based object detectors on the Gen1 dataset.

Method	Type	T	Params	Firing Rate	Energy (mJ)	Gen1	1Mpx
						mAP	mAP
Inception Iacono et al. (2018)	CNN	–	>60M	–	–	0.301	0.340
RRC-Events Chen (2018)	CNN	–	>100M	–	–	0.312	0.343
Matrix Cannici et al. (2019)	RNN+CNN	–	61.5M	–	–	0.310	–
YOLOv3E Jiang et al. (2019)	CNN	–	>60M	–	–	0.312	0.346
RED Perot et al. (2020)	CNN+RNN	–	24.1M	–	> 24	0.400	0.430
ASTMNet Li et al. (2022a)	CNN+RNN	–	>100M	–	–	<u>0.467</u>	0.483
RVT-B Gehrig (2023)	Transformer	–	18.5M	–	–	0.472	<u>0.474</u>
VGG-11Li et al. (2022a)	SNN	5	12.64M	22.22%	11.06	0.174	–
MobileNetLi et al. (2022a)	SNN	5	24.26M	29.44%	5.76	0.147	–
DenseNetLi et al. (2022a)	SNN	5	8.2M	37.2%	3.89	0.189	–
FPDAGNet Zhang et al. (2023)	SNN	5	22M	19.1%	–	0.223	–
SNN-CN Bodden et al. (2024)	SNN	5	12.97M	10.8%	–	0.223	–
KD-CN Bodden et al. (2024)	SNN	5	12.97M	17.4%	–	0.229	–
EMS-YOLO Su et al. (2023)	SNN	5	6.20M	21.15%	–	0.267	–
EMS-YOLO Su et al. (2023)	SNN	5	9.34M	20.09%	–	0.286	–
EMS-YOLO Su et al. (2023)	SNN	5	14.4M	17.80%	–	0.310	–
SFOD Fan et al. (2024)	SNN	5	11.90M	24.04%	7.26	<u>0.321</u>	–
Spike-TransCNN (ours)	SNN	5	24.3M	19.83%	5.49	0.336	0.250

learning rate of $2e - 4$, running the training on a single NVIDIA A100 GPU. For the 1Mpx dataset, we set an effective batch size of 4, a maximum learning rate of $2.45e - 4$, and also train the model using a single NVIDIA A100 GPU.

4.2 MAIN TEST

Quantitative Evaluation. As shown in Table 1, we compare state-of-the-art event-based object detection methods with our Spike-TransCNN on the Gen1 dataset and the 1Mpx dataset. While ANN-based models demonstrate high performance, they come with significant energy consumption. Note that, our Spike-TransCNN significantly outperforms ten state-of-the-art methods. More precisely, our Spike-TransCNN surpasses the best competitor, namely SFOD, achieving a higher mAP (0.336 vs. 0.321) and lower energy consumption (5.49 mJ vs. 7.26 mJ). Additionally, it’s important to note that most event-based object detectors on the 1Mpx dataset are ANN-based models, with very few SNN-based models. However, our Spike-TransCNN has been tested on 1Mpx to establish a benchmark for future comparisons. The limited use of SNNs on 1Mpx may be attributed to the large dataset size, as SNNs typically require more training resources and video memory than ANNs. Fig. 7 provides a visual comparison of the accuracy, spike firing rate, and energy consumption of our approach with other SNN-based methods.

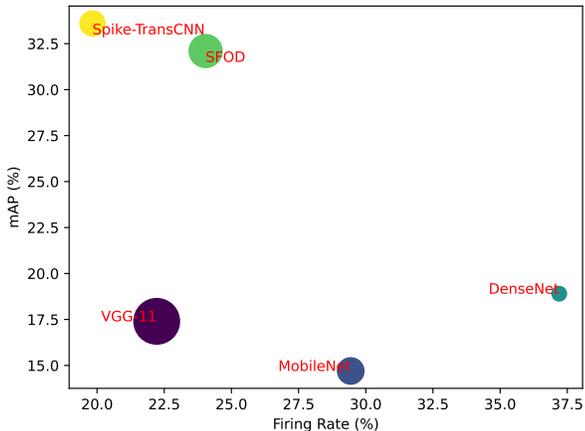


Figure 7: Detection performance vs firing rate of our Spike-TransCNN on the GEN1 dataset. The areas of the circles correspond to the energy.

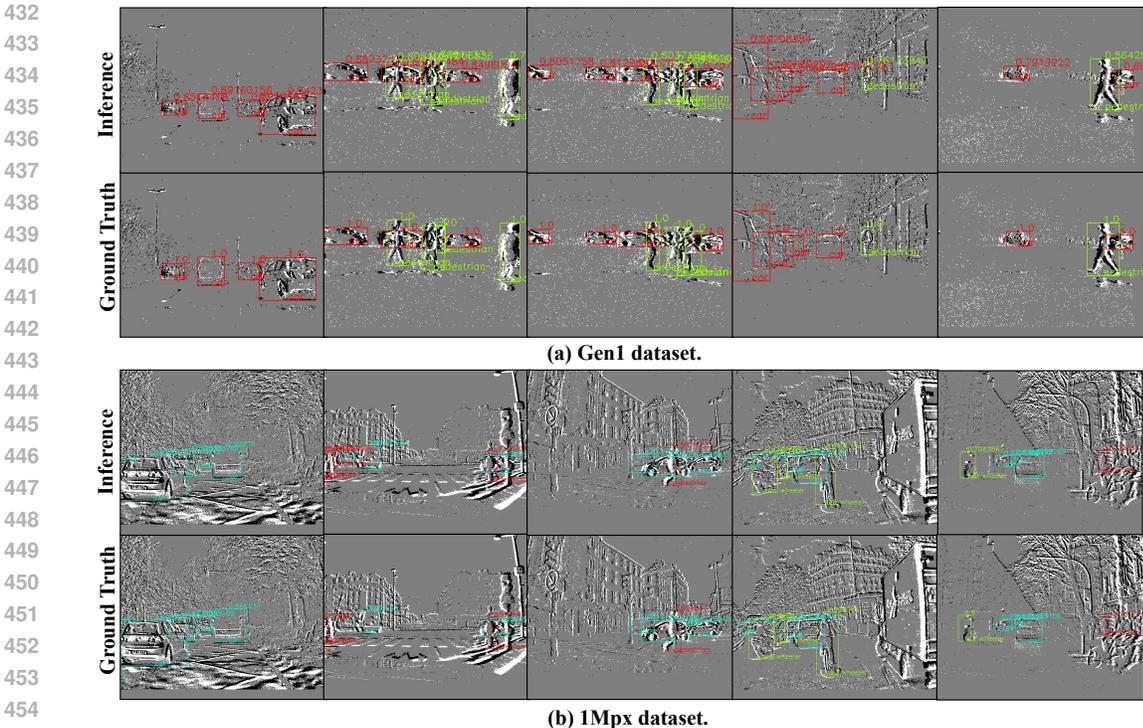


Figure 8: Representative visualization examples of object detection results on the Gen1 dataset and the 1Mpx dataset.

Table 2: The contribution of each component to our Spike-TransCNN on the Gen1 dataset.

Method	STB	SCB	Intra-stage fusion	Inter-stage fusion	Params	mAP _L	mAP _M	mAP _S	mAP
(a)	✓				10.8M	0.208	0.254	0.08	0.184
(b)	✓	✓			21.9M	0.239	0.324	0.134	0.240
(c)	✓	✓	✓		23.5M	0.213	0.326	0.140	0.245
(d)	✓	✓	✓	✓	24.1M	0.251	0.378	0.179	0.285

Visualization Evaluation. We further present representative visualization results on the Gen1 dataset and the 1Mpx dataset (Fig. 8). We select challenging scenarios with occlusions and multi-scale objects. The detection results from our Spike-TransCNN are very close to the ground truth. It indicates that our method performs well in some specific scenarios with single-modal input, particularly when deployed on energy-constrained edge devices.

4.3 ABLATION EXPERIMENTS

We conduct ablation studies on the test set of the Gen1 dataset. We evaluate the impact of various modules and take a deep look at the impact of each design choice as follows.

Contribution of Each Component. To explore the impact of each component on the final performance, we choose Spiking Transformer as the baseline. As shown in Table 2, four methods, namely (a), (b), (d), and (d), utilize Spiking Transformer Block (STB), Spiking CNN Block (SCB), Intra-stage feature fusion module, and Inter-stage feature fusion module, consistently achieve higher performance on the Gen1 dataset than the baseline. Specifically, comparing (b) and (d), the absolute promotion of mAP is 4.5%, which demonstrates that it is feasible to adopt intra-stage spike feature fusion between Spiking Transformer block and Spiking CNN block as well as inter-layer spike feature fusion.

Influence of Spike-driven Token Selection. As shown in Table 3, we explore the influence of the multi-head self-attention operation in our Spike-TransCNN on the Gen1 dataset. We can find that replacing SSA with STS in the first two stages and combining STS with SSA enhanced the overall performance of the model.

4.4 ENERGY CONSUMPTION

The energy consumption of SNNs in neuromorphic hardware are usually assessed based on the number of computational operations Su et al. (2023). In ANNs, each operation involves floating-point multiplications and additions (MAC), and the computation cost is estimated using the number of floating-point operations (FLOPs). SNNs exhibit high energy efficiency in neuromorphic hardware because only neurons involved in spike generation contribute to accumulation calculations (AC), and computations can be performed with roughly the same number of synaptic operations (SyOPs). Hence, we quantify the energy consumption of the original SNN as $E_{SNN} = \sum E_l$, where the energy of the l -th layer can be calculated as:

$$E_l = T \times (S_{fr} \times E_{AC} \times OP_{AC}) + E_{MAC} \times OP_{MAC}, \quad (16)$$

where T represents the time step, S_{fr} denotes the firing rate, and OP_{AC} and OP_{MAC} represent the numbers of AC and MAC operations, respectively. Table 4 presents the energy consumption of our method compared to the theoretical energy consumption of the ANN with the same network architecture. We assume a 32-bit floating-point implementation using 45nm technology, with energy values of $E_{MAC} = 4.6$ pJ and $E_{AC} = 0.9$ pJ Horowitz (2014). Our Spike-TransCNN has an energy consumption of only 5.84 mJ, achieving a $4.4\times$ improvement in energy efficiency compared to the same ANN architecture.

4.5 DISCUSSION

Indeed, our Spike-TransCNN achieves higher detection accuracy and lower power consumption compared to existing pure SNN-based object detection methods, making it suitable for energy-constrained edge devices. Nevertheless, pure SNN models may still exhibit a slight performance gap compared to equivalent ANN architectures or hybrid SNN-ANN models Yu et al. (2024). To further match ANN-level accuracy, we could increase the simulation time steps Luo et al. (2024) or extend the model to handle multiple modalities, such as combining RGB frames with event data.

5 CONCLUSION

This paper proposes a novel hybrid network that takes advantage of both Spiking Transformers and Spiking CNNs for event-based object detection. To the best of our knowledge, this is the first use of a hierarchical hybrid network that includes intra-stage and inter-stage spike feature fusion modules to ensure comprehensive integration of global and multi-scale local information. Experimental results demonstrate that our Spike-TransCNN significantly outperforms existing SNN-based object detectors on the Gen1 dataset, achieving higher mAP and lower energy consumption. We believe our work presents a conceptual hybrid framework that integrates Spiking Transformers and Spiking CNNs, offering potential for various event-based vision applications and feasibility for deployment on neuromorphic hardware.

Table 3: The influence of Spike-driven Token Selection (STS) on the Gen1 dataset.

Operation	Params	mAP _L	mAP _M	mAP _S	mAP
SSA	24.1M	0.251	0.378	0.179	0.285
STS+SSA	24.3M	0.292	0.426	0.237	0.336

Table 4: Energy consumption analysis on the Gen1 dataset.

Model	#OP _{AC}	#OP _{MAC}	Energy	Efficiency
TransCNN (ANN)	1	5.59G	25.70mJ	1×
Spike-TransCNN	5.45G	0.14G	5.49mJ	4.7×

REFERENCES

- 540
541
542 Lennard Bodden, Franziska Schwaiger, Duc Bach Ha, Lars Kreuzberg, and Sven Behnke. Spik-
543 ing centernet: A distillation-boosted spiking neural network for object detection. *arXiv preprint*
544 *arXiv:2402.01287*, 2024.
- 545 Caterina Caccavella, Federico Paredes-Vallés, Marco Cannici, and Lyes Khacef. Low-power event-
546 based face detection with asynchronous neuromorphic hardware. In *arXiv*, 2023.
- 547
548 Marco Cannici, Marco Ciccone, Andrea Romanoni, and Matteo Matteucci. Asynchronous convolu-
549 tional networks for object detection in neuromorphic cameras. In *Proceedings of the IEEE/CVF*
550 *Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0, 2019.
- 551 Kaiwei Che, Zhaokun Zhou, Jun Niu, Zhengyu Ma, Wei Fang, Yanqi Chen, Shuaijie Shen, Li Yuan,
552 and Yonghong Tian. Auto-spikformer: Spikformer architecture search. *Frontiers in Neuroscience*,
553 18:1372257, 2024.
- 554
555 Nicholas FY Chen. Pseudo-labels for supervised learning on dynamic vision sensor data, applied
556 to object detection under ego-motion. In *Proceedings of the IEEE conference on computer vision*
557 *and pattern recognition workshops*, pp. 644–653, 2018.
- 558 Xiang Chen, Jinshan Pan, Jiyang Lu, Zhentao Fan, and Hao Li. Hybrid cnn-transformer feature
559 fusion for single image deraining. In *Proceedings of the AAAI conference on artificial intelligence*,
560 volume 37, pp. 378–386, 2023.
- 561
562 Loïc Cordone, Benoît Miramond, and Philippe Thierion. Object detection with spiking neural net-
563 works on automotive event data. In *2022 International Joint Conference on Neural Networks*
564 *(IJCNN)*, pp. 1–8. IEEE, 2022.
- 565
566 Pierre De Tournemire, Davide Nitti, Etienne Perot, Davide Migliore, and Amos Sironi. A large scale
567 event-based detection dataset for automotive. *arXiv preprint arXiv:2001.08499*, 2020.
- 568
569 Arnaud Delorme, Jacques Gautrais, Rufin Van Rullen, and Simon Thorpe. Spikenet: A simulator
570 for modeling large networks of integrate and fire neurons. *Neurocomputing*, 26:989–996, 1999.
- 571
572 Yimeng Fan, Wei Zhang, Changsong Liu, Mingyang Li, and Wenrui Lu. Sfod: Spiking fusion
573 object detector. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
574 *Recognition*, pp. 17191–17200, 2024.
- 575
576 Jinsheng Fang, Hanjiang Lin, Xinyu Chen, and Kun Zeng. A hybrid network of cnn and transformer
577 for lightweight image super-resolution. In *Proceedings of the IEEE/CVF conference on computer*
578 *vision and pattern recognition*, pp. 1103–1112, 2022.
- 579
580 Wei Fang, Zhaofei Yu, Yanqi Chen, Tiejun Huang, Timothée Masquelier, and Yonghong Tian. Deep
581 residual learning in spiking neural networks. *Advances in Neural Information Processing Systems*,
582 34:21056–21069, 2021.
- 583
584 Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in
585 2021, 2021. URL <https://arxiv.org/abs/2107.08430>.
- 586
587 Mathias Gehrig. Recurrent vision transformers for object detection with event cameras. In *Proceed-*
588 *ings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 13884–13893,
589 2023.
- 590
591 Ryuhei Hamaguchi, Yasutaka Furukawa, Masaki Onishi, and Ken Sakurada. Hierarchical neural
592 memory network for low latency event processing. In *Proceedings of the IEEE/CVF Conference*
593 *on Computer Vision and Pattern Recognition*, pp. 22867–22876, 2023.
- 589
590 Mark Horowitz. 1.1 computing’s energy problem (and what we can do about it). In *IEEE interna-*
591 *tional solid-state circuits conference digest of technical papers*, pp. 10–14, 2014.
- 592
593 Massimiliano Iacono, Stefan Weber, Arren Glover, and Chiara Bartolozzi. Towards event-driven
object detection with off-the-shelf deep learning. In *IEEE/RSJ International Conference on Intel-*
ligent Robots and Systems, pp. 1–9, 2018.

- 594 Zhuangyi Jiang, Pengfei Xia, Kai Huang, Walter Stechele, Guang Chen, Zhenshan Bing, and Alois
595 Knoll. Mixed frame-/event-driven fast pedestrian detection. In *IEEE International Conference on*
596 *Robotics and Automation*, pp. 8332–8338, 2019.
- 597
- 598 Seijoon Kim, Seongsik Park, Byunggook Na, and Sungroh Yoon. Spiking-yolo: spiking neural
599 network for energy-efficient object detection. In *Proceedings of the AAAI conference on artificial*
600 *intelligence*, volume 34, pp. 11270–11277, 2020.
- 601
- 602 D Kinga, Jimmy Ba Adam, et al. A method for stochastic optimization. In *Proceedings of Interna-*
603 *tional conference on learning representations*, volume 5, pp. 6. San Diego, California;, 2015.
- 604
- 605 Alexander Kugele, Thomas Pfeil, Michael Pfeiffer, and Elisabetta Chicca. Hybrid snn-ann: Energy-
606 efficient classification and object detection for event-based vision. In *DAGM German Conference*
607 *on Pattern Recognition*, pp. 297–312. Springer, 2021.
- 608
- 609 Jianing Li, Jia Li, Lin Zhu, Xijie Xiang, Tiejun Huang, and Yonghong Tian. Asynchronous spatio-
610 temporal memory network for continuous event-based object detection. *IEEE Transactions on*
611 *Image Processing*, 31:2975–2987, 2022a.
- 612
- 613 Yang Li, Xiang He, Yiting Dong, Qingqun Kong, and Yi Zeng. Spike calibration: Fast and accu-
614 rate conversion of spiking neural network for object detection and segmentation. *arXiv preprint*
615 *arXiv:2207.02702*, 2022b.
- 616
- 617 Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A 128×128 120 db 15 μ s latency
618 asynchronous temporal contrast vision sensor. *IEEE journal of solid-state circuits*, 43(2):566–
619 576, 2008.
- 620
- 621 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
622 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
623 *Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014,*
624 *Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- 625
- 626 Li Liu, Wanli Ouyang, Xiaogang Wang, Paul Fieguth, Jie Chen, Xinwang Liu, and Matti Pietikäinen.
627 Deep learning for generic object detection: A survey. *International journal of computer vision*,
628 128:261–318, 2020.
- 629
- 630 Xinhao Luo, Man Yao, Yuhong Chou, Bo Xu, and Guoqi Li. Integer-valued training and spike-
631 driven inference spiking neural network for high-performance and energy-efficient object detec-
632 tion. *arXiv preprint arXiv:2407.20708*, 2024.
- 633
- 634 Wolfgang Maass. Networks of spiking neurons: the third generation of neural network models.
635 *Neural networks*, 10(9):1659–1671, 1997.
- 636
- 637 Yansong Peng, Yueyi Zhang, Peilin Xiao, Xiaoyan Sun, and Feng Wu. Better and faster: Adaptive
638 event conversion for event-based object detection. In *Proceedings of the AAAI Conference on*
639 *Artificial Intelligence*, volume 37, pp. 2056–2064, 2023a.
- 640
- 641 Yansong Peng, Yueyi Zhang, Zhiwei Xiong, Xiaoyan Sun, and Feng Wu. Get: Group event trans-
642 former for event-based vision. In *Proceedings of the IEEE/CVF International Conference on*
643 *Computer Vision*, pp. 6038–6048, 2023b.
- 644
- 645 Etienne Perot, Pierre De Tournemire, Davide Nitti, Jonathan Masci, and Amos Sironi. Learning
646 to detect objects with a 1 megapixel event camera. *Advances in Neural Information Processing*
647 *Systems*, 33:16639–16652, 2020.
- 648
- 649 Christoph Posch, Daniel Matolin, and Rainer Wohlgenannt. A qvga 143 db dynamic range frame-
650 free pwm image sensor with lossless pixel-level video compression and time-domain cds. *IEEE*
651 *Journal of Solid-State Circuits*, 46(1):259–275, 2010.
- 652
- 653 Mohamed Sayed and Gabriel Brostow. Improved handling of motion blur in online object detection.
654 In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp.
655 1706–1716, 2021.

- 648 Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using
649 large learning rates. In *Artificial intelligence and machine learning for multi-domain operations*
650 *applications*, volume 11006, pp. 369–386. SPIE, 2019.
- 651 Qiaoyi Su, Yuhong Chou, Yifan Hu, Jianing Li, Shijie Mei, Ziyang Zhang, and Guoqi Li. Deep
652 directly-trained spiking neural networks for object detection. In *Proceedings of the IEEE/CVF*
653 *International Conference on Computer Vision*, pp. 6555–6565, 2023.
- 654 Abhishek Tomy, Anshul Paigwar, Khushdeep S Mann, Alessandro Renzaglia, and Christian Laugier.
655 Fusing event-based and rgb camera for robust object detection in adverse conditions. In *Proceed-*
656 *ings of the International Conference on Robotics and Automation*, pp. 933–939, 2022.
- 657 Dongsheng Wang, Xu Jia, Yang Zhang, Xinyu Zhang, Yaoyuan Wang, Ziyang Zhang, Dong Wang,
658 and Huchuan Lu. Dual memory aggregation network for event-based object detection with learn-
659 able representation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 37,
660 pp. 2492–2500, 2023a.
- 661 Qingyu Wang, Duzhen Zhang, Tielin Zhang, and Bo Xu. Attention-free spikformer: Mixing spike
662 sequences with simple linear transforms. *arXiv preprint arXiv:2308.02557*, 2023b.
- 663 Yuan-Kai Wang, Shao-En Wang, and Ping-Hsien Wu. Spike-event object detection for neuromorphic
664 vision. *IEEE Access*, 11:5215–5230, 2023c.
- 665 Yuchen Wang, Kexin Shi, Chengzhuo Lu, Yuguo Liu, Malu Zhang, and Hong Qu. Spatial-temporal
666 self-attention for asynchronous spiking neural networks. In *Proceedings of International Joint*
667 *Conference on Artificial Intelligence*, pp. 3085–3093, 2023d.
- 668 Ziming Wang, Ziling Wang, Huaning Li, Lang Qin, Runhao Jiang, De Ma, and Huajin Tang. Eas-
669 snn: End-to-end adaptive sampling and representation for event-based detection with recurrent
670 spiking neural networks. In *arXiv*, 2024.
- 671 Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for
672 training high-performance spiking neural networks. *Frontiers in neuroscience*, 12:331, 2018.
- 673 Man Yao, Jiakui Hu, Zhaokun Zhou, Li Yuan, Yonghong Tian, Bo Xu, and Guoqi Li. Spike-driven
674 transformer. *Advances in neural information processing systems*, 36, 2024.
- 675 Lixing Yu, Hanqi Chen, Ziming Wang, Shaojie Zhan, Jiankun Shao, Qingjie Liu, and Shu Xu.
676 Spikingvit: a multi-scale spiking vision transformer model for event-based object detection. *IEEE*
677 *Transactions on Cognitive and Developmental Systems*, 2024.
- 678 Mengwen Yuan, Chengjun Zhang, Ziming Wang, Huixiang Liu, Gang Pan, and Huajin Tang. Train-
679 able spiking-yolo for low-latency and high-performance object detection. *Neural Networks*, 172:
680 106092, 2024.
- 681 Hu Zhang, Luziwei Leng, Kaiwei Che, Qian Liu, Jie Cheng, Qinghai Guo, Jiangxing Liao, and Ran
682 Cheng. Automotive object detection via learning sparse events by temporal dynamics of spiking
683 neurons. *arXiv preprint arXiv:2307.12900*, 2023.
- 684 Chenlin Zhou, Han Zhang, Zhaokun Zhou, Liutao Yu, Liwei Huang, Xiaopeng Fan, Li Yuan,
685 Zhengyu Ma, Huihui Zhou, and Yonghong Tian. Qkformer: Hierarchical spiking transformer
686 using qk attention. *arXiv preprint arXiv:2403.16552*, 2024.
- 687 Zhaokun Zhou, Yuesheng Zhu, Chao He, Yaowei Wang, YAN Shuicheng, Yonghong Tian, and
688 Li Yuan. Spikformer: When spiking neural network meets transformer. In *The Eleventh Interna-*
689 *tional Conference on Learning Representations*, 2023.
- 690 Lin Zhu, Xiao Wang, Yi Chang, Jianing Li, Tiejun Huang, and Yonghong Tian. Event-based video
691 reconstruction via potential-assisted spiking neural network. In *Proceedings of the IEEE/CVF*
692 *Conference on Computer Vision and Pattern Recognition*, pp. 3594–3604, 2022.
- 693 Nikola Zubić, Daniel Gehrig, Mathias Gehrig, and Davide Scaramuzza. From chaos comes or-
694 der: Ordering event representations for object recognition and detection. In *Proceedings of the*
695 *IEEE/CVF International Conference on Computer Vision*, pp. 12846–12856, 2023.

702 Nikola Zubic, Mathias Gehrig, and Davide Scaramuzza. State space models for event cameras.
703 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.
704 5819–5828, 2024.
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755

A APPENDIX

A.1 EVENT DATA PREPROCESSING

The output of an event camera with a resolution of $H \times W$ can be represented as an event sequence, denoted as $E = e_i 1^N$, where $e_i = (x_i, y_i, t_i, p_i)$. Here, $p_i \in \{-1, 1\}$ represents the polarity of a brightness change that occurs at time t_i and pixel position (x_i, y_i) . The change is triggered for the pixel (x_n, y_n) at timestamp t_n when the log-intensity $\ln L$ changes beyond the pre-defined threshold θ . This dynamic visual sensing mechanism is depicted by the inequality:

$$\ln L(x_n, y_n, t_n) - \ln L(x_n, y_n, t_n - \Delta t_n) p_n \theta, \quad (17)$$

Here, the polarity $p_n \in \{1, -1\}$ indicates whether the brightness is increasing or decreasing, and Δt_n represents the temporal sampling interval of DVS at a pixel.

In our investigation, we have implemented the technique presented in Gehrig (2023) for preprocessing event data. Our preprocessing generates a 4-dimensional tensor E from discrete event data. The first dimension consists of T components associated with T discretization steps of time. The second dimension includes two components signifying polarity. The third and fourth dimensions represent the height and width of the event camera. We process a set of events ε within the time duration $[t_a, t_b)$ in the following manner:

$$E(\tau, p, x, y) = \sum_{\varepsilon} \delta(p - p_n) \delta(x - x_n) \delta(\tau - \tau_n), \quad (18)$$

$$\tau_n = \left\lfloor \frac{t_n - t_a}{t_b - t_a} \cdot T \right\rfloor$$

The provided equation describes the handling of a set of events ε over a time interval $[t_a, t_b)$, with each t_n falling between t_a and t_b . Here, $\delta(\cdot)$ represents the Dirac delta function, where $\delta(t)$ is 0 for all $t \neq 0$, and $\int \delta(t) dt$ equals 1. The value of T is determined by the selected number of discrete time steps. Following this procedure, a four-dimensional tensor $E \in [T, 2, H, W]$ is obtained, where T , H , and W denote the aggregation time, preprocessed height, and width, respectively.

The datasets used are Gen1¹ and 1Mpx².

A.2 EVALUATION CRITERIA

Table 5: Summary of common metrics for event-based object detection.

Metric	Unit	Description
mAP	-	mAP averaged over ten IoUs: $\{0.5:0.05:0.95\}$.
mAP _{0.5}	-	mAP at a fixed IoU=0.50.
mAP _{0.75}	-	mAP at a fixed IoU=0.75.
mAP _S	-	mAP for small objects of area smaller than 32^2 .
mAP _M	-	mAP for objects of area between 32^2 and 96^2 .
mAP _L	-	mAP for large objects of area bigger than 96^2 .
Model size	MB	The number of parameters for the learning-based model.
Power consum.	mJ	The energy consumption of the SNN model through AC and MAC operations in a neuromorphic chip.

The evaluation metrics for object detection based on event data are summarized in Table 5. For both datasets, the primary metric we use is mean average precision (mAP) Lin et al. (2014). Additionally, to demonstrate the capability of our method for detecting objects at multiple scales, we also use mAP_L, mAP_M, and mAP_S. The AP is derived from precision and recall using the following formulas:

$$AP = \int_0^1 \max\{p(r' | r' \geq r)\} dr, \quad (19)$$

¹<https://www.prophesee.ai/2020/01/24/prophesee-gen1-automotive-detection-dataset>

²<https://www.prophesee.ai/2020/11/24/automotive-megapixel-event-based-dataset/>

where r denotes the recall, $p(r)$ is the precision-recall curve.

Thus, the mAP is calculated as the average of AP values across all object categories as follows:

$$\text{mAP} = \frac{\sum_{i=1}^{C_b} AP(i)}{C_b}, \quad (20)$$

where C_b represents the number of object classes. At present, MS COCO³ is the most widely used benchmark for evaluating event-based object detection methods. Instead of using a fixed IoU threshold, MS COCO Perot et al. (2020) provides a few metrics with various IoUs (i.e., mAP, mAP_{0.5}, and mAP_{0.75}) and AP across different scales (i.e., mAP_S, mAP_M, and mAP_L). As neuromorphic cameras offer continuous visual streams, object detection labels are annotated on the RGB image or reconstructed image at a fixed frame rate. Thus, many existing methods assess the detection accuracy at the timestamp when the label is provided, lacking the ability to evaluate the entire event stream continuously.

We calculate the spike Firing Rate by counting the number of spikes and the number of neurons in each layer of the network, which can be computed using the following formula:

$$\text{Firing Rate} = \frac{\text{Numbers of spikes}}{\text{Numbers of neurons}}. \quad (21)$$

A.3 THE IMPLEMENTATION OF LIF NEURONS

The firing and membrane potential updates are two main modules in the LIF neuron Delorme et al. (1999), for which we provided the implementation functions in Algorithm 1 and Algorithm 2. In Algorithm 1, the forward propagation function during the firing process is presented, and a gradient substitution function is defined in the backward propagation function. We use the functions defined in Algorithm 2 as the LIF neuron for this paper, where the input x has dimensions of [T, 2, H, W]. Table 6 provides specific numerical values of the hyperparameters used in the LIF model in this work.

Table 6: Values of LIF parameters.

Parameter	thresh	lens	decay	time_window
Value	0.5	0.5	0.25	5

A.4 PATCH EMBEDDING

In order to obtain hierarchical features, we implement the patch embedding operation for downsampling before each STB. This operation consists of conv1 * 1, max-pooling, and LIF. We downsample the spatial resolution using max pooling, which helps preserve the binary nature of the elements.

A.5 SPIKE TRANSFORMER BLOCK

As shown in Fig. 10, detailed design details of the Spike Transformer Block(STB) are provided. The main difference between the first two stages and the last two stages is the multi-head attention. In the first two stages, we used Spike-driven token selection (STS) in the multi-head attention, and Spike self-attention (SSA) mechanism in the last two stages.

A.6 THE ORDER OF STB AND SCB

Spiking Transformer Block (STB) and Spiking CNN Block (SCB) serve as the core modules for feature extraction in the backbone. While STB focuses on extracting global features, SCB specializes in capturing local multi-scale features. In each stage of our proposed method, we first employ STB to extract global features, followed by utilizing SCB to capture local multi-scale features. Our analysis of the sequential order of these two modules in each stage, illustrated in Fig. 11, revealed

³<https://github.com/cocodataset/cocoapi>

Algorithm 1: Approximate Firing Function

Input: *input***Output:** *output*

```

1 Function ActFun (ctx, input) :
  Data: input
  Result: output
2   Function forward (ctx, input) :
  |   Save input for backward pass in ctx;
  |   return float(input > thresh);
3   Function backward (ctx, grad_output) :
  |   Data: grad_output
  |   Result: grad_input
  |   input  $\leftarrow$  Retrieve saved tensors from ctx;
  |   grad_input  $\leftarrow$  grad_output.clone();
  |   temp  $\leftarrow$  abs(input - thresh) < lens;
  |   temp  $\leftarrow$  temp / (2  $\times$  lens);
  |   return grad_input  $\times$  float(temp);
  Invoke forward and backward functions as needed;
  return forward(ctx, input);

```

Algorithm 2: Membrane Potential Update

Input: *x***Output:** *output*

```

1 Function mem_update (x) :
  Data: x
  Result: output
2   mem  $\leftarrow$  torch.zeros_like(x[0]).to(device);
3   spike  $\leftarrow$  torch.zeros_like(x[0]).to(device);
4   output  $\leftarrow$  torch.zeros_like(x);
5   mem_old  $\leftarrow$  0;
6   for i = 0 to time_window - 1 do
7     if i  $\geq$  1 then
8       | mem  $\leftarrow$  mem_old  $\times$  decay  $\times$  (1 - spike.detach()) + x[i];
9     else
10      | mem  $\leftarrow$  x[i];
11      | spike  $\leftarrow$  ActFun(mem);
12      | mem_old  $\leftarrow$  mem.clone();
13      | output[i]  $\leftarrow$  spike;
14  return output;

```

that the model is minimally impacted by the order in which they are applied. As depicted in Table 7, regardless of whether SCB is used before STB (as in A) or vice versa (as in B), the overall effect on the model is negligible. We attribute this to the fusion module proposed in this paper, which comprehensively integrates global and local features. Notably, the information from the preceding module is retained, irrespective of the order in which STB and SCB are employed.

A.7 SPIKE-ELEMENT-WISE ADDITION

In this study, \oplus operations are employed in both intra-stage and inter-stage contexts. Herein, we aim to demonstrate how this operation effectively mitigates issues such as gradient vanishing or gradient explosion.

Taking the SEW Block Fang et al. (2021) as an example, upon the application of \oplus :

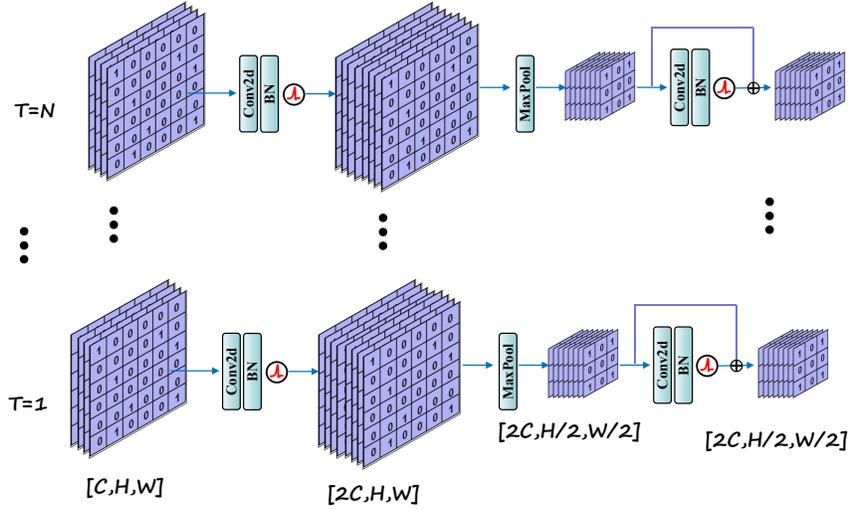


Figure 9: Patch Embedding.

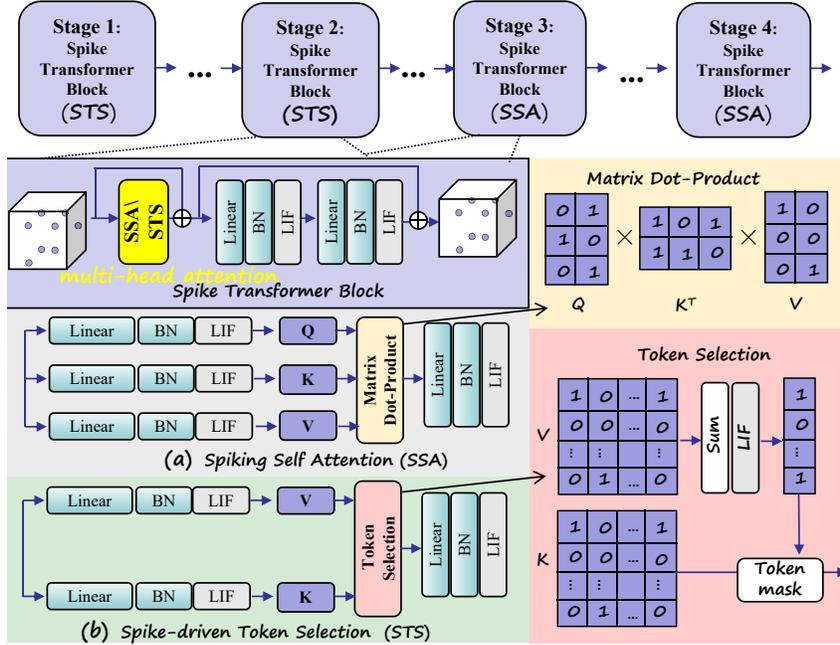


Figure 10: Spike Transformer Block. Spike-driven Token Selection (STS) is used for token selection in the first two stages, and Spiking Self Attention (SSA) is employed for global feature extraction in the latter two stages. (a) and (b) present the specific details of SSA and QKA. The bottom right corner illustrates the operational processes of ‘Matrix Dot-Product’ and ‘Token Selection’ in SSA and QKA, respectively.

$$O^l[t] = A^l[t] \oplus S^l[t], \quad (22)$$

the gradient from the output of the $(l+k-1)$ -th SEW block to the input of the l -th SEW block can be computed in a layer-by-layer manner:

$$\frac{\partial O_j^{l+k-1}[t]}{\partial S_j^l[t]} = \prod_{i=0}^{k-1} \frac{\partial (A_j^{l+i}[t] \oplus S_j^{l+i}[t])}{\partial S_j^{l+i}[t]} = \prod_{i=0}^{k-1} \frac{\partial (0 + S_j^{l+i}[t])}{\partial S_j^{l+i}[t]}. \quad (23)$$

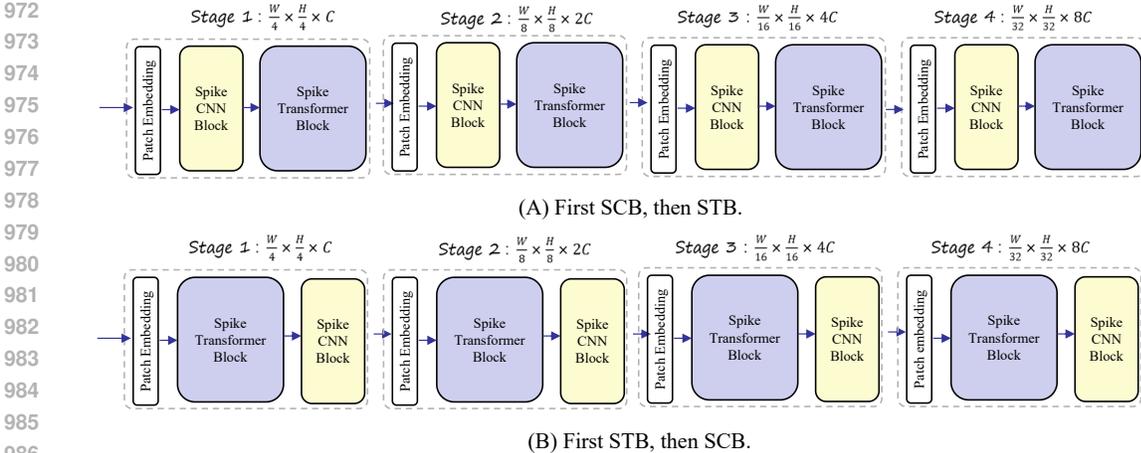


Figure 11: The order of Spiking Transformer Block (STB) and Spiking CNN Block (SCB) in the backbone.

Table 7: Performance comparison based on the order of Spiking Transformer Block (STB) and Spiking CNN Block (SCB) in the backbone.

Backbone	mAP_L	mAP_M	mAP_S	mAP	mAP_{50}
A	0.329	0.420	0.177	0.328	0.587
B	0.336	0.426	0.267	0.336	0.604

The equality holds as identity mapping is achieved by setting $A^{l+i}[t] \equiv 0$. Since the gradient is a constant ($= 1$), the SEWA (\oplus) can overcome gradient vanishing or gradient explosion.

A.8 TRAINING DETAILS

In this research, we present an in-depth exploration of the specific training process values for the Gen1 and 1Mpx datasets depicted in Fig. 12 and Fig. 13. These figures not only offer a comprehensive overview of the learning rates and loss values at each iteration but also provide valuable insights into the stability and evolution of the training procedure. Furthermore, the visual representations encapsulate the mAP_{50} , mAP , mAP_{large} , mAP_{middle} , and mAP_{small} metrics, affording a holistic comprehension of the network’s performance across diverse evaluation criteria throughout the duration of the training process. It is noteworthy that the evaluation metrics consistently demonstrate a gradual enhancement in performance on the training set, underscoring the efficacy and robustness of the training approach adopted in this study. This detailed analysis serves to enrich our understanding of the intricate dynamics and progressive refinement observed during the training phase, thus contributing to a more nuanced interpretation of the network’s learning process.

A.9 MORE VISUALIZATION RESULTS

In addition to the visualization results previously presented, we offer further insights into the characteristics of the Gen1 and 1Mpx datasets in Fig. 14 and Fig. 15, respectively. These visualizations provide a deeper understanding of the dataset attributes and distribution, shedding light on the diverse features and patterns captured within the data.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

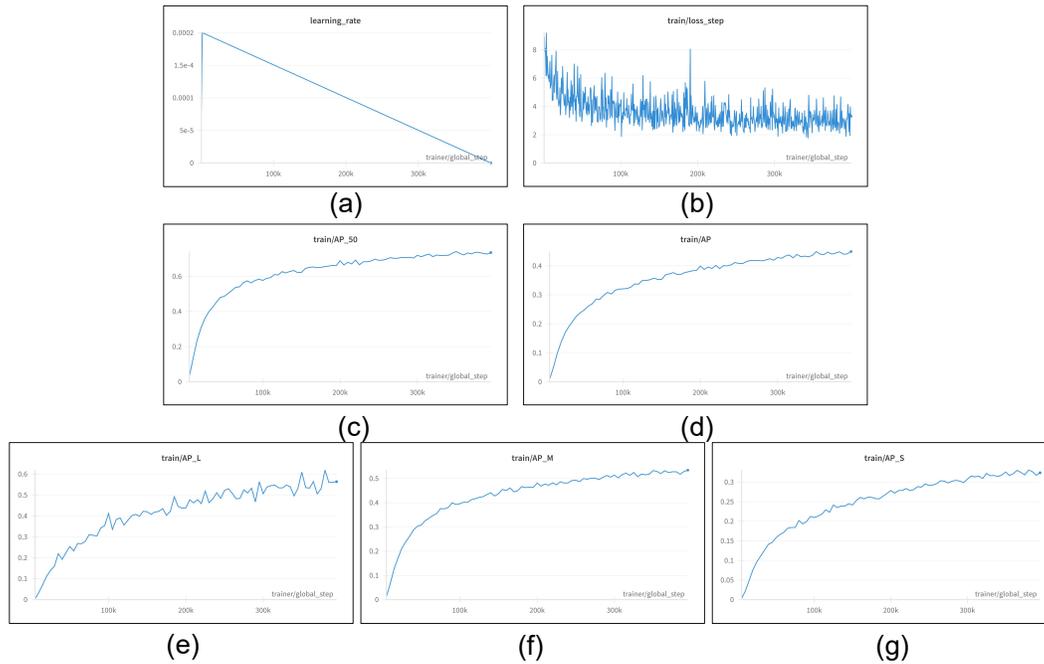


Figure 12: Training Details on the Gen1. (a) Learning rate, (b) loss reduction, (c) mAP_{50} , (d) mAP , (e) mAP_{large} , (f) mAP_{middle} , (g) mAP_{small} .

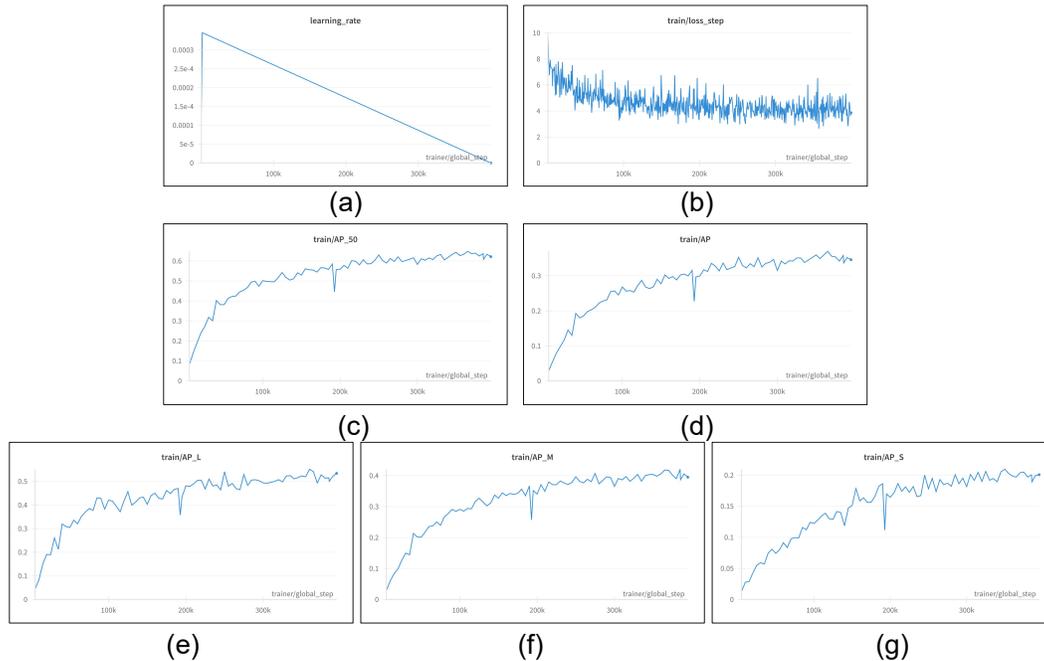


Figure 13: Training Details on the 1Mpx. (a) Learning rate, (b) loss reduction, (c) mAP_{50} , (d) mAP , (e) mAP_{large} , (f) mAP_{middle} , (g) mAP_{small} .

1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106
1107
1108
1109
1110
1111
1112
1113
1114
1115
1116
1117
1118
1119
1120
1121
1122
1123
1124
1125
1126
1127
1128
1129
1130
1131
1132
1133

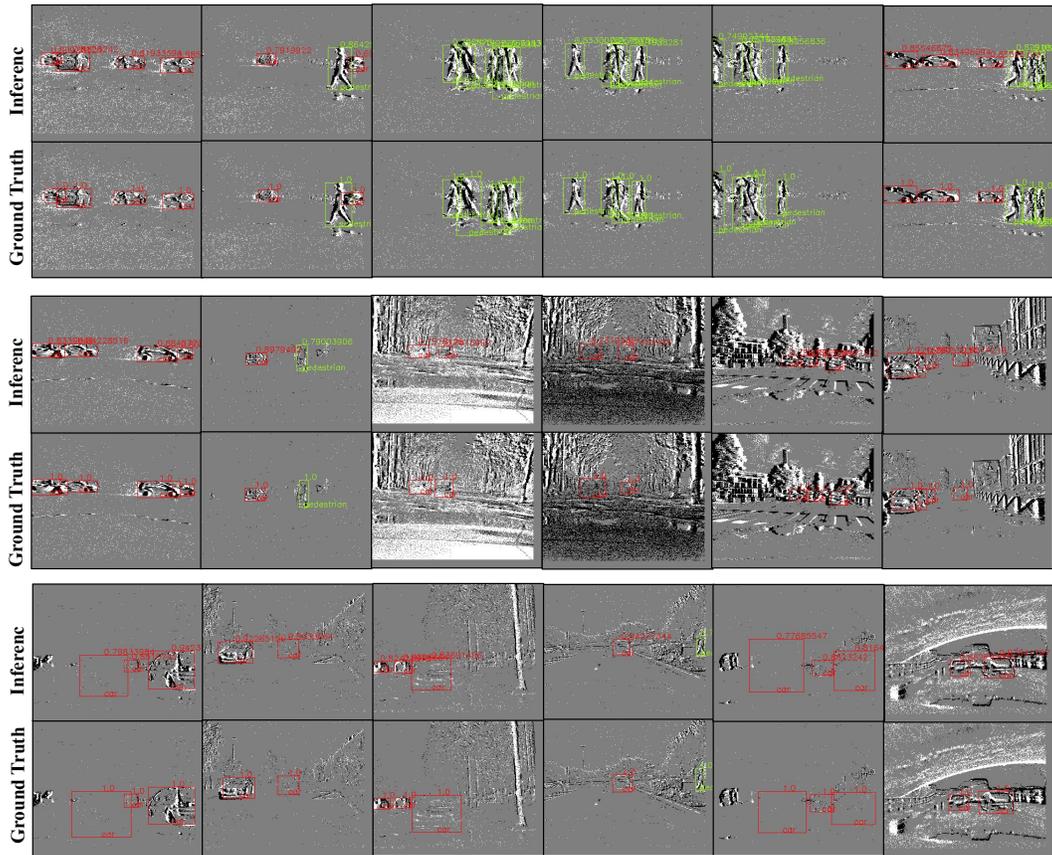


Figure 14: More visualization results on the Gen1 dataset.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187

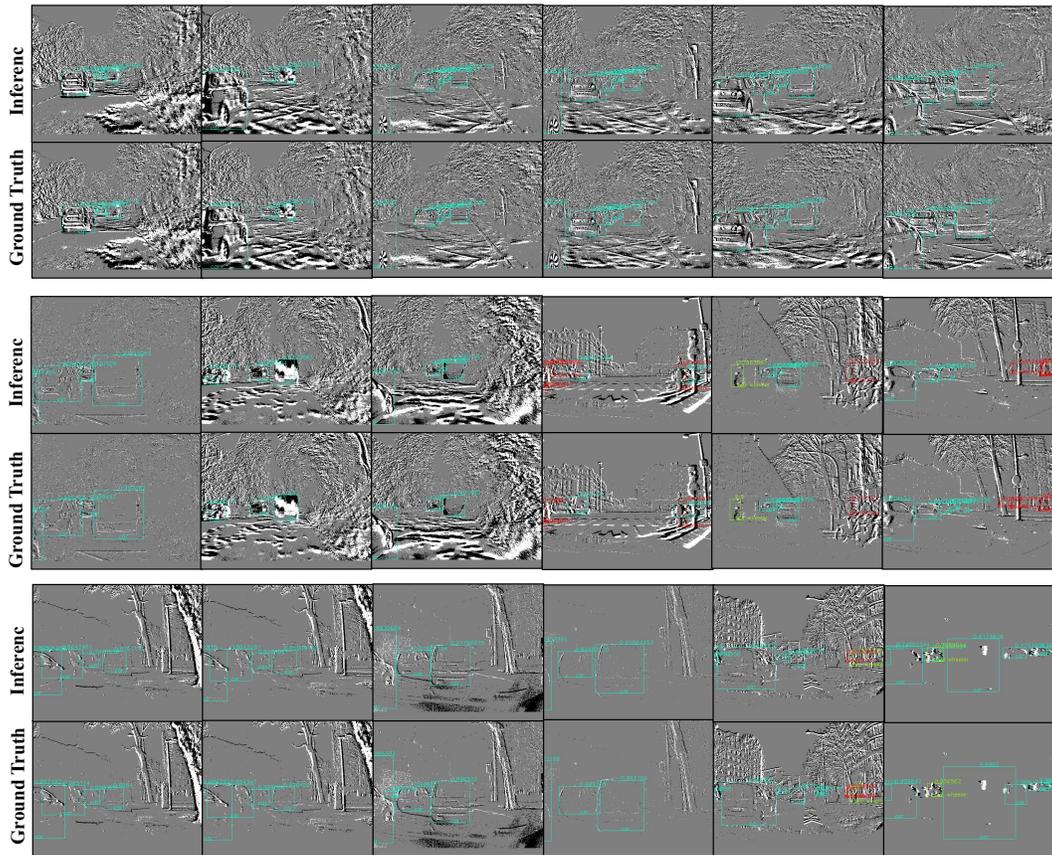


Figure 15: More visualization results on the 1Mpx dataset.