

## State Space Models for Event Cameras

Nikola Zubić, Mathias Gehrig, Davide Scaramuzza  
 Robotics and Perception Group, University of Zurich, Switzerland

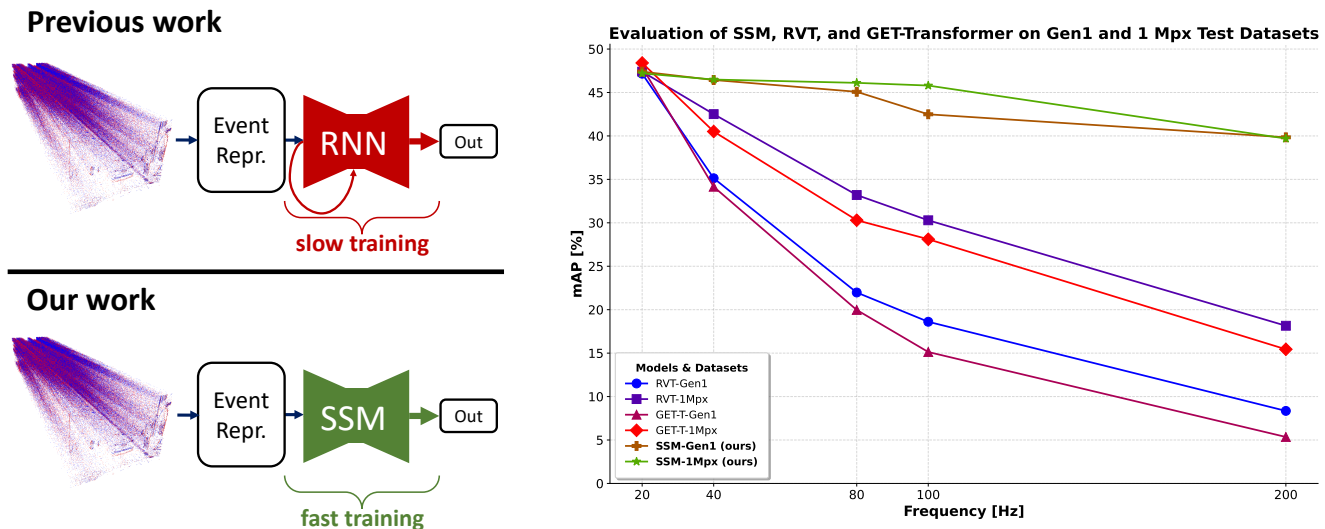


Figure 1. **Top-Left** Previous works [12, 29] use RNN architectures with convolutional or attention mechanisms to train models that have superior performance on downstream tasks. However, the use of RNNs leads to slower training and the learned weights only generalize well to data deployed at the same frequency as that used at training time. **Bottom-Left** We solve this problem by utilizing SSMs for temporal aggregation, which enables faster training by either utilizing the S4 model [16] or S5 [36] parallel scans. By their nature, these models allow deployment at different frequencies than those used at training time since they have a learnable timescale parameter. **Right** Our SSM-based models achieve an average performance drop between training and testing frequencies of 3.31 mAP averaged on both Gen1 [7] and 1 Mpx [29] datasets, while RVT [12] and GET [28] have a drop of 21.25 and 24.53 mAP, respectively.

### Abstract

Today, state-of-the-art deep neural networks that process event-camera data first convert a temporal window of events into dense, grid-like input representations. As such, they exhibit poor generalizability when deployed at higher inference frequencies (i.e., smaller temporal windows) than the ones they were trained on. We address this challenge by introducing state-space models (SSMs) with learnable timescale parameters to event-based vision. This design adapts to varying frequencies without the need to retrain the network at different frequencies. Additionally, we investigate two strategies to counteract aliasing effects when deploying the model at higher frequencies. We comprehensively evaluate our approach against existing methods based on RNN and Transformer architectures across various benchmarks, including Gen1 and 1 Mpx event camera datasets. Our results demonstrate that SSM-based models train 33% faster and also exhibit minimal performance degradation when tested at higher frequencies than the

training input. Traditional RNN and Transformer models exhibit performance drops of more than 20 mAP, with SSMs having a drop of 3.31 mAP, highlighting the effectiveness of SSMs in event-based vision tasks.

### 1. Introduction

Event cameras emerged as a class of sensor technologies that noticeably deviate from the operational mechanics of conventional frame-based cameras [9]. While standard frame-based cameras [42] capture full-frame luminance levels at fixed intervals, event cameras record per-pixel relative brightness changes in a scene at the time they occur. The output is, therefore, an asynchronous stream of events in a four-dimensional spatio-temporal space. Each event is represented by its spatial coordinates on the pixel array, the temporal instance of the occurrence, and its polarity, which denotes the direction of the brightness change and encapsulates the increase or decrease in luminance. The

power of event cameras primarily lies in their ability to capture dynamic scenes with unparalleled temporal resolution (microseconds). This property becomes invaluable in rapidly changing environments or applications requiring very fast response times [8, 39]. However, the richness of the spatio-temporal data they generate introduces complexities in data interpretation and processing. Sophisticated algorithms are required to efficiently parse and make sense of the high-dimensional data space. As such, while event cameras represent a promising frontier in visual sensor technologies, their pervasive utilization depends upon solving these inherent computational challenges.

Current methodologies for addressing problems with event cameras fall predominantly into two categories. The first involves converting the raw stream of spatio-temporal events into dense representations akin to multi-channel images [12, 22, 40, 43]. This transformation allows leveraging conventional computer vision techniques designed for frame-based data. The second category employs sparse computational paradigms, often utilizing spiking neural networks or graph-based architectures [4, 11]. While promising, these methods are not without limitations; they frequently encounter hardware incompatibility issues and compromised accuracy. In this work, we utilize dense representations for their advantages in computational efficiency.

Despite the advances in both paradigms, models trained on event representations at a specific frequency exhibit poor generalizability when deployed in settings with higher frequencies which is crucial for high-speed, dynamical visual scenarios. Additionally, to achieve high performance, it is necessary to include recurrent memory, thereby sacrificing computational efficiency during the training phase. An ideal model would seamlessly merge the training speed of convolutional architectures with the temporal sensitivity and memory benefits inherent to recurrent models.

While recent advancements have introduced efficient recurrent vision transformer architectures [12, 24] to achieve better performance, they face several limitations. Specifically, these architectures suffer from longer training cycles due to their reliance on conventional recurrent mechanisms.

The issue of slow training and generalization at higher event representation frequency than the lower one we trained on remains unresolved as conventional recurrent training methodologies are predominantly utilized in event-based vision. These methods do not incorporate learnable timescale parameters, thus inhibiting the model’s capacity to generalize across varying inference frequencies.

In this work, we address this limitation by introducing structured variations of state-space models [16, 36] as layers within our neural network framework.

State-space models [14] function as CNN during training and are converted to an efficient RNN at test time. To achieve this, S4 [16] employs a technique where the SSM,

which is not practical for training on modern hardware due to its sequential nature, is transformed into a CNN by unrolling. S5 [36] uses parallel scans during training, and is employed as RNN during inference. State-space models [14] can be deployed at different frequencies at inference time because they are Linear Time-Invariant (LTI) continuous-time systems that can be transformed into a discrete-time system with an arbitrary step size. This feature permits inference at arbitrary frequencies by globally adjusting the timescale parameter based on the ratio between the new and old frequency sampling rates. Consequently, we tackle the longstanding issue in event-based vision, which requires multiple training cycles with different frequencies to adapt the neural network for various frequencies during inference.

For the task of object detection, we find that incorporating state-space layers accelerates the training by up to 33% relative to existing recurrent vision transformer approaches [12, 24]. This is achieved while maintaining performance comparable to existing methods. Notably, our approach demonstrates superior generalization to higher temporal frequencies with a drop of only 3.31 mAP, while previous methods experience a drop of 21.25 mAP or more. Also, we achieve comparable performance to RVT although we use a linear state-space model rather than a non-linear LSTM model. This also shows that the complexity of the LSTM might not be needed. For this to work, we introduce two strategies (frequency-selective masking and  $H_2$  norm) to counteract the aliasing effect encountered with increased temporal frequencies in event cameras. First one is a low-pass bandlimiting modification to SSM that encourages the learned convolutional kernels to be smooth. Second one mitigates the aliasing problem by attenuating the frequency response after a chosen frequency. We argue that state-space modeling offers a significant new direction for research in event-based vision, offering promising solutions to the challenges inherent in processing event-based data effectively and efficiently.

Our contributions are concisely outlined as follows:

- We introduce state-space models for event cameras to address two key challenges in event-based vision: (i) model performance degradation when operating event cameras at temporal frequencies different from their training conditions and (ii) training efficiency.
- Our experimental results outperform existing methods at higher frequencies by 20 mAP on average and show 33% increase in the training speed.
- We introduce two strategies (bandlimiting &  $H_2$  norm) designed to alleviate aliasing issues.

## 2. Related Work

### 2.1. Object detection with Event Cameras

Approaches in event camera literature, thus in object detection, can be broadly classified into two branches.

The first research branch investigates dense feed-forward methods. Early attempts in this direction relied on a constrained temporal window for generating event representations [3, 17, 19]. The resultant models were deficient in tracking slow-moving or stationary objects, as they failed to incorporate data beyond this limited time-frame. To mitigate these drawbacks, later studies introduced recurrent neural network (RNN) layers [29] into the architecture. The RNN component enhanced the model’s capacity for temporal understanding, thereby improving its object detection capabilities. The work of Zubic et al. [43] takes this a step further by optimizing event representations and incorporating cutting-edge swin transformer architecture. Nonetheless, their approach was limited in its ability to re-detect slowly moving objects following extended periods of absence. Subsequent research [12, 24], merged the transformer and RNN architectures, pushing the performance further. However, this significantly increased computational demands during the training phase. Importantly, all methodologies examined to date suffer from an inability to adapt when deployed at variable frequencies.

The second research branch investigates the use of Graph Neural Networks (GNNs) or Spiking Neural Networks (SNNs). GNNs dynamically construct a spatio-temporal graph where new nodes and edges are instantiated by selectively sub-sampling events and identifying pre-existing nodes that are proximate in the space-time continuum [11, 18, 38]. A pivotal challenge lies in architecting the network such that information can disseminate effectively across extensive regions of the space-time volume. This becomes particularly important when dealing with large objects that exhibit slow relative motion to the camera. Moreover, while aggressive sub-sampling is often necessary to achieve low-latency inference, it introduces the risk of omitting important information from the event stream. On the other hand, SNNs [1, 6, 34] transmit information sparsely within the system. Unlike RNNs, SNNs emit spikes only when a threshold is met, making them hard to optimize due to the non-differentiable spike mechanism. Some solutions [27] bypass the threshold, but this sacrifices sparsity in deeper layers. Overall, SNNs remain a challenging area needing more foundational research for optimized performance.

### 2.2. Continuous-time Models

Gu et al. [16] introduced the S4 model as an alternative to CNNs and Transformers for capturing long-range dependencies through LTI systems. This was followed by

the S4D model [15], designed for easier understanding and implementation, offering similar performance to S4. The S5 model [36] improved efficiency by avoiding frequency domain computations and utilizing time-domain parallel scans. However, these models have not been thoroughly evaluated on complex, high-dimensional visual data with significant temporal resolution.

In our study, we empirically show that S4, S4D and S5 models achieve results on-par with state-of-the-art when combined with attention mechanisms on complex data. We also identify and address aliasing issues in these models, proposing two corrective strategies. Our work extends the range and robustness of continuous-time models for complex sequence modeling tasks such as object detection for event cameras.

## 3. Method

In this section, we firstly formalize the operating mechanism of event cameras and provide a notation used for describing state-space models in the preliminaries (Sec. 3.1). Secondly, we describe our approach of incorporating variants of state-space models as layers within our block (Sec. 3.2). This innovative design solves the problems associated with slow training and the variable frequency inference for event cameras.

### 3.1. Preliminaries

**Event cameras.** Event cameras are bio-inspired vision sensors that capture changes in log intensity per pixel asynchronously, rather than capturing entire frames at fixed intervals. Formally, let  $I(x, y, t)$  denote the log intensity at pixel coordinates  $(x, y)$  and time  $t$ . An event  $e$  is generated at  $(x, y, t)$  whenever the change in log intensity  $\Delta I$  exceeds a certain threshold  $C$ :

$$\Delta I(x, y, t) = I(x, y, t) - I(x, y, t - \Delta t) \geq C \quad (1)$$

Each event  $e$  is a tuple  $(x, y, t, p)$ , where  $(x, y)$  are the pixel coordinates,  $t$  is the timestamp, and  $p = \{-1, 1\}$  is the polarity of the event, indicating the direction of the intensity change.

**State-Space Models (SSMs).** Linear State-Space Models (SSMs) form the crucial part of the backbone in our architecture, where we compare S4 [16], S4D [15] and S5 [36] layer variants. Given an input vector  $\mathbf{u}(t) \in \mathbb{R}^U$ , a latent state vector  $\mathbf{x}(t) \in \mathbb{R}^P$ , and an output vector  $\mathbf{y}(t) \in \mathbb{R}^M$ , the governing equations of a continuous-time linear SSM can be mathematically represented as:

$$\frac{d\mathbf{x}(t)}{dt} = \mathbf{A}\mathbf{x}(t) + \mathbf{B}\mathbf{u}(t), \quad \mathbf{y}(t) = \mathbf{C}\mathbf{x}(t) + \mathbf{D}\mathbf{u}(t), \quad (2)$$

The model is parameterized by a state transition matrix  $\mathbf{A} \in \mathbb{R}^{P \times P}$ , an input matrix  $\mathbf{B} \in \mathbb{R}^{P \times U}$ , an output matrix  $\mathbf{C} \in \mathbb{R}^{M \times P}$ , and a direct transmission matrix  $\mathbf{D} \in \mathbb{R}^{M \times U}$ .

Given a fixed step size  $\Delta$ , this continuous-time model can be *discretized* into a linear recurrence using various methods such as Euler, bilinear, or zero-order hold (ZOH):

$$\mathbf{x}_k = \bar{\mathbf{A}}\mathbf{x}_{k-1} + \bar{\mathbf{B}}\mathbf{u}_k, \quad \mathbf{y}_k = \bar{\mathbf{C}}\mathbf{x}_k + \bar{\mathbf{D}}\mathbf{u}_k, \quad (3)$$

The parameters in the discrete-time model are functions of the continuous-time parameters, defined by the chosen discretization method. Details on SSMs are available in the appendix.

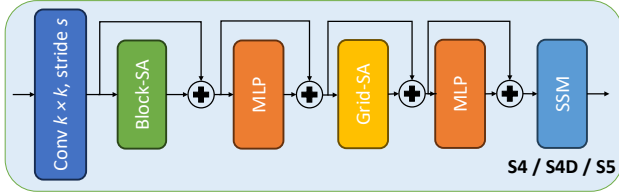


Figure 2. SSM-ViT block structure

### 3.2. SSM-ViT block

In this section, we introduce the SSM-ViT block, a novel block depicted in Figure 2, which showcases the structured flow of the proposed block structure, designed for efficient event-based information processing.

We build a 4-stage hierarchical backbone like in [12] where in each stage we utilize our proposed SSM-ViT block. Events are processed into a tensor representation like in [12] before they are used as input to the first stage. Each stage takes the previous features as input and reuses the SSM state from the last timestep to compute features for the next stage. By saving SSM states, each stage retains temporal information for the whole feature map, while also being able to generalize to different frequencies since we use SSM instead of RNN that is used in [12].

Regarding the block structure, initially, the input undergoes convolution with a defined kernel of size  $k \times k$  with a stride  $s$ . This operation effectively captures the essential spatial features of the input. Following the convolution operation, the structure introduces a 'Block-SA' module. This module is pivotal in implementing self-attention mechanisms, but it does so specifically within local windows. The localized nature of the attention in this block ensures a focus on immediate spatial relations, allowing for a detailed representation of close-proximity features.

Subsequent to 'Block-SA', the 'Grid-SA' module comes into play. In contrast to the localized approach of the previous block, 'Grid-SA' employs dilated attention, ensuring a global perspective. This module, by considering a broader scope, encapsulates a more comprehensive understanding of the spatial relations and the overall structure of the input.

The final and crucial component of the block structure is the State-Space Model (SSM) block, that is designed to

compute both the output and the state in parallel, either by using S4 [16], S4D [15] or S5 [36] model variant. This ensures temporal consistency and a seamless transition of information between consecutive time steps. Efficient computation is crucial since it allows for the faster training than RNN, and timescale parameter on temporal aggregation is important since we can rescale it during inference and deploy at any frequency we want.

### 3.3. Low-pass bandlimiting

However, the capability to deploy state space models at resolutions higher than those encountered during training is not without its drawbacks. It gives rise to the well-documented issue of aliasing [32, 33], which occurs when the kernel bandwidth surpasses the Nyquist frequency. We address this challenge in following two subsections - 3.3.1 & 3.3.2.

#### 3.3.1 Output Masking

In the realm of signal processing, frequency content beyond the Nyquist rate can lead to aliasing, adversely affecting model performance. To address this, we integrate a frequency-selective masking strategy into the training and inference processes. This bandlimiting method has been empirically validated to be crucial for generalizing across different frequencies [32, 33], with ablation studies indicating a decrease in accuracy by as much as 20% in its absence.

Let the SSM be governed by a matrix  $\mathbf{A}$ , with its diagonal elements  $a_n$  influencing the temporal evolution of the system states. The kernel's basis function for the  $n$ -th state is given by  $K_n(t) = e^{ta_n} B_n$ , where the frequency characteristics are primarily dictated by the imaginary part of  $a_n$  -  $\Im(a_n)$  and  $B_n$  represents the  $B$  matrix.

To modulate the frequency spectrum of the model, we define a hyperparameter  $\alpha$ , which is used for masking of the computed effective frequency  $f_n$  for each state:

$$f_n = \frac{\Delta t}{r} \cdot \frac{|\Im(a_n)|}{2\pi}, \quad (4)$$

where  $\Delta t$  denotes the discrete time-step, and  $r$  is the rate at which we train the model, by default it is 1, rate gets halved when deploying at twice the trained frequency etc.

The bandlimiting mask is then applied as follows:

$$C_n = \begin{cases} C_n & \text{if } f_n \leq \frac{\alpha}{2}, \\ 0 & \text{otherwise,} \end{cases} \quad (5)$$

where  $C_n$  represents the coefficients in the state representation. The selection of  $\alpha$  is critical, with  $\alpha = 1.0$  representing the Nyquist limit under idealized conditions. However, practical considerations and model constraints typically necessitate a lower empirical threshold. Our findings suggest that setting  $\alpha = 0.5$  yields optimal outcomes for systems

with diagonal state matrices, such as in S4D and S5 configurations.

This frequency-selective masking is proven to be a cornerstone for the adaptability of our SSMs, significantly contributing to their generalization across differing frequencies.

### 3.3.2 $H_2$ Norm

This section introduces another strategy to mitigate the aliasing problem by suppressing the frequency response of the continuous-time state-space model beyond a selected frequency  $\omega_{\min}$ .

This approach makes use of the  $H_2$  norm of a continuous-time linear system which measures the power (or steady-state covariance) of the output response to unit white-noise input.

Given a continuous-time system described by matrices  $\mathbf{A}$ ,  $\mathbf{B}$ , and  $\mathbf{C}$ , the transfer function  $\mathbf{G}(s)$  from the Laplace transform can be defined in the frequency domain as:

$$\mathbf{G}(j\omega) = \mathbf{C}(j\omega\mathbf{I} - \mathbf{A})^{-1}\mathbf{B}, \quad (6)$$

where  $\omega$  denotes the frequency, and  $\mathbf{I}$  is the identity matrix. The  $H_2$  norm is computed as the integral of the squared magnitude of the frequency response over the range of interest, typically the entire frequency spectrum. However, in our case, we would like to suppress the frequency response of the system to frequencies beyond  $\omega_{\min}$  which can be done by minimizing the following:

$$\|\mathbf{G}\|_{H_2(\omega_{\min}, \infty)} = \sqrt{\frac{1}{\pi} \int_{\omega_{\min}}^{\infty} \|\mathbf{G}(j\omega_k)\|_F^2 d\omega}, \quad (7)$$

as part of the loss function, where  $\|\cdot\|_F$  is the Frobenius norm. In practice, to numerically estimate this integral, we choose a maximum frequency  $\omega_{\max}$  and discretize the frequency range  $[\omega_{\min}, \omega_{\max}]$  into  $N$  points and apply numerical integration methods to the squared Frobenius norm of  $\mathbf{G}(j\omega)$ . This yields an approximate  $H_2$  norm of the system in the desired frequency range.

## 4. Experiments

We perform ablation studies and systematic evaluations of our proposed models utilizing both the Gen1 [7] and 1 Mpx [29] event-based camera datasets. We assess the model’s proficiency in adapting to unseen inference frequencies during training time on both datasets. We employed two variants for training across both datasets: the base model ViT-SSM-B, alongside its scaled-down derivative ViT-SSM-S (small). Additionally, to study the robustness and generalization capabilities of our architecture, we subject it to empirical testing on the DSEC dataset [13], which was not part of the original training corpus and the visual results are provided in the appendix.

## 4.1. Setup

**Implementation details.** Our models are trained using 32-bit precision arithmetic across 400k iterations, making use of the ADAM optimizer [20] and a OneCycle learning rate schedule [37] that decays linearly from its highest value. We adopt a mixed batching technique that applies standard Backpropagation Through Time (BPTT) to half of the batch samples and Truncated BPTT (TBPTT) to the other half. This technique was introduced in the RVT [12] and ablated in their supplementary material. We found it to be equally effective also for SSMs. For the S4(D) method, details can be found in [15, 16]. As for S5 [36], we utilize an efficient implementation from scratch in PyTorch; original public version of S5 is available in JAX. Data augmentation is carried out through random horizontal flips and zooming operations, both inward and outward. Event representations are formed based on 50 ms time windows that correspond to 20 Hz sampling frequency, divided into  $T = 10$  discrete bins. We incorporate a YOLOX detection head [10] that integrates IOU, class, and regression losses, averaged across both batch and sequence lengths during each optimization phase. Bandlimiting is implemented during training and inference of the SSM model by masking the output matrix  $\mathbf{C}$  as explained in 3.3.1. In case of  $H_2$  norm, we add it as term to the loss function of the model. We conduct training on the GEN1 dataset using A100 GPU, employing a batch size of 8 and a sequence length of 21. We use a global learning rate of  $2e-4$ , which is propagated to the SSM components. For the 1 Mpx dataset, we train with a batch size of 12 and a sequence length of 10, using a learning rate of  $3.5e-4$  across two A100 GPUs.

**Datasets.** The Gen1 Automotive Detection dataset [7] comprises 39 hours of event camera footage with a resolution of  $304 \times 240$  pixels. It includes 228k bounding boxes for cars and 28k for pedestrians. Using original evaluation criteria [7], we discard bounding boxes having a side length shorter than 10 pixels or a diagonal less than 30 pixels.

The 1 Mpx dataset [29] similarly focuses on driving environments but offers recordings at a higher resolution of  $720 \times 1280$  pixels, captured over multiple months during both day-time and night-time. It contains roughly 15 hours of event data and accumulates around 25 million bounding box labels, spread across three classes: cars, pedestrians, and two-wheelers. Following original evaluation criteria, we eliminate bounding boxes with a side length under 20 pixels and a diagonal less than 60 pixels, while also reducing the input resolution to  $640 \times 360$  pixels.

## 4.2. Benchmark comparisons

In this section, we present a detailed comparison of our proposed methods, S4D-ViT-B and S5-ViT-B, with the SotA approaches in the domain of event-based vision for object detection, which can be seen in Table 1. The evalua-

| Method                  | Backbone          | Detection Head | Gen1        |             | 1 Mpx       |             | Params (M) |
|-------------------------|-------------------|----------------|-------------|-------------|-------------|-------------|------------|
|                         |                   |                | mAP         | Time (ms)   | mAP         | Time (ms)   |            |
| Asynet [27]             | Sparse CNN        | YOLOv1 [31]    | 14.5        | -           | -           | -           | 11.4       |
| AEGNN [35]              | GNN               | YOLOv1         | 16.3        | -           | -           | -           | 20.0       |
| Spiking DenseNet [5]    | SNN               | SSD [23]       | 18.9        | -           | -           | -           | 8.2        |
| Inception + SSD [17]    | CNN               | SSD            | 30.1        | 19.4        | 34.0        | 45.2        | > 60*      |
| RRC-Events [3]          | CNN               | YOLOv3 [30]    | 30.7        | 21.5        | 34.3        | 46.4        | > 100*     |
| MatrixLSTM [2]          | RNN + CNN         | YOLOv3         | 31.0        | -           | -           | -           | 61.5       |
| YOLOv3 Events [19]      | CNN               | YOLOv3         | 31.2        | 22.3        | 34.6        | 49.4        | > 60*      |
| RED [29]                | CNN + RNN         | SSD            | 40.0        | 16.7        | 43.0        | 39.3        | 24.1       |
| ERGO-12 [43]            | Transformer       | YOLOv6 [21]    | <b>50.4</b> | 69.9        | 40.6        | 100.0       | 59.6       |
| RVT-B [12]              | Transformer + RNN | YOLOX [10]     | 47.2        | 10.2        | <u>47.4</u> | 11.9        | 18.5       |
| Swin-T v2 [25]          | Transformer + RNN | YOLOX          | 45.5        | 26.6        | 46.4        | 34.5        | 21.1       |
| Nested-T [28, 41]       | Transformer + RNN | YOLOX          | 46.3        | 25.9        | 46.0        | 33.5        | 22.2       |
| GET-T [28]              | Transformer + RNN | YOLOX          | <u>47.9</u> | 16.8        | <b>48.4</b> | 18.2        | 21.9       |
| <b>S4D-ViT-B (ours)</b> | Transformer + SSM | YOLOX          | 46.2        | <u>9.40</u> | 46.8        | <u>10.9</u> | 16.5       |
| <b>S5-ViT-B (ours)</b>  | Transformer + SSM | YOLOX          | 47.4        | <b>8.16</b> | 47.2        | <b>9.57</b> | 18.2       |

Table 1. **Comparisons on test sets of Gen1 and 1 Mpx datasets (20 Hz)**. Best results in **bold** and second best underlined. A star \* suggests that this information was not directly available and estimated based on the publications. Runtime is measured in milliseconds for a batch size of 1. We used a T4 GPU for SSM-ViT to compare against indicated timings in prior work [29] on comparable GPUs (Titan Xp). Ours achieves comparable results with other SotA approaches and generalizes well to other frequencies where others fail.

tion is conducted on two distinct datasets: Gen1 and 1 Mpx, to assess the performance and robustness of these methods under varying conditions.

The comparative analysis, as summarized in Table 1, encompasses different backbones and detection heads. These include Sparse CNNs, GNNs, SNNs, RNNs, and various implementations of Transformers.

Our S4D-ViT-B and S5-ViT-B models, which integrate SSMs with Attention (3.2), demonstrate competitive performance across both datasets. Specifically, on the 1 Mpx dataset, our models achieve mAP scores of 46.8 and 47.2, respectively, which are competitive with the leading scores in this benchmark. While our models do not outperform the top-performing ERGO-12 [43] and GET-T [28] in terms of mAP, they exhibit a notable balance between accuracy and model complexity, as indicated by their parameter counts of 16.5M and 17.5M, respectively.

Notably, our models perform consistently well across different frequencies as it can be seen in Table 2, highlighting their robustness and generalizability in comparison to RVT and GET-T, which tend to exhibit large performance drops. We achieve an average performance drop between training and testing frequencies of 3.31 mAP averaged on both Gen1 [7] and 1 Mpx [29] datasets, while RVT [12] and GET [28] have a drop of 21.25 and 24.53 mAP, respectively.

In summary, our proposed S4D-ViT-B and S5-ViT-B models establish themselves as effective and efficient contenders in the field of event-based vision for object detec-

tion. Their competitive performance metrics, coupled with their generalization capabilities across various inference frequencies, make them valuable contributions to event-based vision community.

| Model | Dataset | Frequency Evaluation (Hz) |       |       |        |        | Perf. Drop  |
|-------|---------|---------------------------|-------|-------|--------|--------|-------------|
|       |         | 20 Hz                     | 40 Hz | 80 Hz | 100 Hz | 200 Hz |             |
| RVT   | Gen1    | 47.16                     | 35.13 | 21.98 | 18.61  | 8.35   | 26.14       |
|       | 1Mpx    | 47.40                     | 42.51 | 33.20 | 30.29  | 18.15  | 16.36       |
| S5    | Gen1    | 47.40                     | 46.44 | 45.08 | 42.49  | 39.84  | <u>3.94</u> |
|       | 1Mpx    | 47.20                     | 46.49 | 46.11 | 45.80  | 39.70  | <b>2.68</b> |
| GET   | Gen1    | 47.90                     | 34.15 | 19.97 | 15.13  | 5.35   | 29.25       |
|       | 1Mpx    | 48.40                     | 40.51 | 30.30 | 28.11  | 15.44  | 19.81       |

Table 2. Evaluation of RVT [12], S5 [36], and GET [28] across different frequencies on test datasets

### 4.3. Ablation study

In this section we investigate different SSMs and their initializations with and without bandlimiting parameter for event-based vision to address the problem of inference at different frequencies (4.3.1). After that, we study the importance of SSM layers in various stages during training of the model (4.3.2). Finally, we evaluate models on differing frequencies with two proposed strategies (4.3.3).

#### 4.3.1 SSMs: initializations & bandlimiting

This section presents an ablation study focusing on the performance impact of various SSM variants and their initialization strategies in conjunction with the bandlimiting parameter  $\alpha$ . The SSM variants under consideration are S4

| Model    | Mean Average Precision - mAP <sub>val</sub> |                |              |              |
|----------|---|----------------|--------------|--------------|
|          | $\alpha = 0$                                | $\alpha = 0.5$ | $\alpha = 1$ | Average      |
| S4-legS  | 46.66                                       | -              | -            | 46.66        |
| S4D-legS | 46.93                                       | 47.33          | 46.50        | 46.92        |
| S4D-inv  | 46.15                                       | 46.23          | 46.11        | 46.16        |
| S4D-lin  | 44.82                                       | 46.02          | 45.04        | 45.29        |
| S5-legS  | 48.33                                       | <b>48.48</b>   | 48.00        | <b>48.27</b> |
| S5-inv   | 47.26                                       | <u>47.43</u>   | 46.98        | <u>47.22</u> |
| S5-lin   | 46.12                                       | 46.40          | 45.59        | 46.04        |

Table 3. Performance comparison between the S4 [16], S4D [15] and S5 [36] models for different values of  $\alpha$  and initializations on Gen1 [7] validation dataset.

[16], S4D [15], and the more recent S5 [36]. The analysis is conducted on the Gen1 [7] validation dataset, with the mean Average Precision (mAP) serving as the performance metric. Table 3 provides a comprehensive view of how different initialization strategies (legS, inv, lin) introduced in [15] and values of  $\alpha \in \{0, 0.5, 1\}$  influence model performance. The S4-legS model, not equipped for bandlimiting in non-diagonal matrix scenarios, achieves a baseline mAP of 46.66 at  $\alpha = 0$ . The S4D variants demonstrate a diverse performance spectrum. The S4D-legS variant, particularly at  $\alpha = 0.5$ , achieves the highest mAP of 47.33 among the S4D models, also maintaining the best average mAP of 46.92. The S4D-inv and S4D-lin variants show less favorable outcomes, with mAPs peaking at 46.23 and 46.02, respectively, for  $\alpha = 0.5$ .

Notably, the S5 model variants exhibit an improvement over their predecessors. The S5-legS variant achieves the highest mAP of 48.48 at  $\alpha = 0.5$ , and also records the best average mAP of 48.27. This result is not only the best in its category but also the best overall. The S5-inv and S5-lin models also demonstrate comparable performance, with the former reaching its peak mAP of 47.43 at  $\alpha = 0.5$ .

This study emphasizes the critical role of initialization strategies and bandlimiting in optimizing SSM-based neural networks for event-based vision tasks. The distinct performance variations across different models and configurations underscore the importance of selecting appropriate initializations and  $\alpha$  values, as these choices impact the efficacy of the models in handling the dynamic and complex nature of event-based vision data. Higher values of parameter  $\alpha$  encourage SSM to learn smoother kernels and discard more complex and higher-frequency convolution kernels.

### 4.3.2 SSM Utilization Analysis

In this ablation study, we examine the impact of employing temporal recurrence exclusively in a subset of network stages, or not using it at all. Our methodology involves manipulating the SSMs within the network by resetting their

states at predetermined stages during each timestep. This approach enables us to isolate and evaluate the impact of SSM layers’ presence or absence while maintaining a consistent parameter count across different model configurations.

| S1 | S2 | S3 | S4 | mAP <sub>RVT</sub> | mAP <sub>S4D.5</sub> | mAP <sub>S5.5</sub> |
|----|----|----|----|--------------------|----------------------|---------------------|
|    |    |    |    | 33.90              | 39.99                | 43.67               |
|    |    |    | ✓  | 41.68              | 43.11                | 46.10               |
|    |    | ✓  | ✓  | 46.10              | 45.33                | 47.52               |
|    | ✓  | ✓  | ✓  | 48.82              | 47.02                | 48.41               |
| ✓  | ✓  | ✓  | ✓  | 49.52              | 47.33                | 48.48               |

Table 4. SSM contribution in various stages on the Gen1 dataset.

Table 4 shows the outcomes of these manipulations on the Gen1 validation dataset.  $S4D_{.5}$  represents S4D model, and  $S5_{.5}$  is S5 model with  $\alpha = 0.5$ . The data clearly indicate that the complete removal of SSMs from the network leads to a highest decrease in detection performance. This underscores the pivotal role of SSMs in enhancing the model’s capability. On the other hand, initiating the use of SSMs from the fourth stage onwards consistently enhances performance, suggesting a critical threshold for the impact of temporal information processing in the later stages of the network. Another intriguing observation is the performance boost obtained by incorporating an SSM at the very initial stage. This suggests that the early integration of temporal information is also beneficial for the overall detection performance. As a result of these findings, our preferred configuration includes the SSM component right from the initial stage, thereby leveraging the advantages of temporal information processing throughout the entire network. Noteworthy is the observation that the performance drop when employing SSMs is less pronounced than with RVT, suggesting our approach’s superior robustness to temporal aggregation from certain stages compared to RVT.

### 4.3.3 Evaluation at different frequencies

In this section, we delve into the comparative performance analysis of three models: RVT [12], S4D [15], and S5 [36], across various frequency ranges. This evaluation is crucial in determining the robustness and adaptability of these models under diverse inference frequencies, specifically at 20 Hz, 40 Hz, 80 Hz, 100 Hz, and 200 Hz. Our model is trained with event representations formed based on 50 ms time windows that correspond to 20 Hz sampling frequency. Evaluation is done at this frequency, and others above mentioned. With RVT, there are no changes when evaluating at different frequencies, while with SSMs rate  $r$  is halved for the double the inference frequency etc.

Each model in Table 5 is assessed across datasets and model sizes, including Gen1<sub>Base</sub>, Gen1<sub>Small</sub>, and 1 Mpx<sub>Base</sub>. *Base* and *Small* represent two variants of mod-

| Model    | Dataset/Size          | Strategy       | Frequency Evaluation (Hz) |       |       |        |        | Performance Drop |
|----------|-----------------------|----------------|---------------------------|-------|-------|--------|--------|------------------|
|          |                       |                | 20 Hz                     | 40 Hz | 80 Hz | 100 Hz | 200 Hz |                  |
| RVT [12] | Gen1 <sub>Base</sub>  | -              | 49.52                     | 37.16 | 23.25 | 19.36  | 7.83   | 27.62            |
|          | Gen1 <sub>Small</sub> | -              | 48.68                     | 35.28 | 19.95 | 16.05  | 5.75   | 29.42            |
|          | 1Mpx <sub>Base</sub>  | -              | 45.95                     | 40.93 | 31.70 | 29.00  | 18.16  | 16.00            |
| S4D [15] | Gen1 <sub>Base</sub>  | $H_2$ norm     | 46.83                     | 45.98 | 43.91 | 40.10  | 36.11  | 5.31             |
|          |                       | $\alpha = 0.5$ | 47.33                     | 46.36 | 44.51 | 40.02  | 35.98  | 5.61             |
|          | Gen1 <sub>Small</sub> | $H_2$ norm     | 45.88                     | 45.11 | 41.05 | 38.00  | 34.05  | 6.33             |
|          |                       | $\alpha = 0.5$ | 46.30                     | 45.21 | 42.11 | 38.61  | 33.00  | 6.57             |
|          | 1Mpx <sub>Base</sub>  | $H_2$ norm     | 46.66                     | 45.85 | 43.33 | 41.80  | 37.01  | 4.66             |
|          |                       | $\alpha = 0.5$ | 47.93                     | 46.78 | 44.56 | 41.11  | 36.18  | 5.77             |
| S5 [36]  | Gen1 <sub>Base</sub>  | $H_2$ norm     | 48.60                     | 47.11 | 46.06 | 43.80  | 40.51  | 4.23             |
|          |                       | $\alpha = 0.5$ | 48.48                     | 47.34 | 46.11 | 43.23  | 40.03  | 4.30             |
|          | Gen1 <sub>Small</sub> | $H_2$ norm     | 47.33                     | 46.32 | 44.03 | 41.12  | 38.98  | 4.72             |
|          |                       | $\alpha = 0.5$ | 47.83                     | 46.58 | 44.46 | 41.11  | 38.95  | 5.06             |
|          | 1Mpx <sub>Base</sub>  | $H_2$ norm     | 48.65                     | 47.53 | 47.11 | 46.63  | 40.91  | 3.11             |
|          |                       | $\alpha = 0.5$ | 48.35                     | 47.60 | 47.21 | 46.50  | 40.80  | 2.82             |

Table 5. Evaluation of RVT, S4D, and S5 across different frequencies on validation datasets

els on the Gen1 and 1 Mpx datasets, the base one being the larger one with 16.5M parameters for the S4D model and 18.2M parameters for the S5 model (as presented in Table 1 parameters’ column), and small one with 8.8M parameters for S4D and 9.7M parameters for S5 model. The focal point of this analysis is the performance drop, calculated as the average difference between the original performance at 20 Hz and performances at higher frequencies. A notable aspect of this study is the inherent advantage of the S4D and S5 models due to their incorporation of a learnable timescale parameter. This significantly enhances their adaptability, allowing them to dynamically adjust to varying frequencies. This feature is particularly salient in the S4D and S5 models, which are further analyzed based on different operational strategies: the  $H_2$  norm and bandlimiting with  $\alpha = 0.5$ . The inclusion of the learnable timescale parameter in these models underscores their capability to maintain performance across a wide range of frequencies.

The analysis reveals that both the  $H_2$  norm and bandlimiting strategies with  $\alpha = 0.5$  offer comparable performance across the assessed frequency ranges. However, a slight edge is observed with the  $H_2$  norm, particularly at very high frequencies. This marginal superiority can be attributed to the fact that  $H_2$  norm approach does not explicitly mask output’s ( $C$ ) matrix columns.

In the Appendix, we study pure-SSM and SSM models in combination with ConvNext [26].

## 5. Conclusion

In this paper, we presented a novel approach for enhancing the adaptability and training efficiency of models designed for event-based vision, particularly in object detec-

tion tasks. Our methodology leverages the integration of SSMs with a Vision Transformer (ViT) architecture, creating a hybrid SSM-ViT model. This integration not only addresses the long-standing challenge of performance degradation at varying temporal frequencies but also significantly accelerates the training process.

The key innovation of our work lies in the use of learnable timescale parameters within the SSMs, enabling the model to adapt dynamically to different inference frequencies without necessitating multiple training cycles. This feature represents a substantial advancement over existing methods, which require retraining for different frequencies.

The SSM-ViT model outperforms existing methods by 20 mAP at higher frequencies and exhibits a 33% increase in training speed. Furthermore, our introduction of two novel strategies to counteract the aliasing effect (a crucial consideration in high-frequency deployment) further reinforces the model’s suitability for real-world applications. These strategies, involving frequency-selective masking and  $H_2$  norm adjustments, effectively mitigate the adverse effects of aliasing, ensuring the model’s reliability across a spectrum of temporal resolutions. We believe that our approach opens new avenues for research and application in high-speed, dynamic visual environments, setting a new benchmark in the domain of event-based vision.

## 6. Acknowledgment

This work was supported by the European Research Council (ERC) under grant agreement No. 864042 (AG-ILEFLIGHT).



## References

- [1] Julian Büchel, Gregor Lenz, Yalun Hu, Sadique Sheik, and Martino Sorbaro. Adversarial attacks on spiking convolutional neural networks for event-based vision. *Front. Neurosci.*, 16, 2022. 3
- [2] Marco Cannici, Marco Ciccone, Andrea Romanoni, and Matteo Matteucci. Asynchronous convolutional networks for object detection in neuromorphic cameras. In *CVPRW*, 2019. 6
- [3] Nicholas F. Y. Chen. Pseudo-labels for supervised learning on dynamic vision sensor data, applied to object detection under ego-motion. In *CVPRW*, 2018. 3, 6
- [4] Loic Cordone, Benoît Miramond, and Sonia Ferrante. Learning from event cameras with sparse spiking convolutional neural networks. In *IEEE IJCNN*, 2021. 2
- [5] Loic Cordone, Benoît Miramond, and Phillipe Thierion. Object detection with spiking neural networks on automotive event data. In *IEEE IJCNN*, 2022. 6
- [6] Javier Cuadrado, Ulysse Rançon, Benoit R. Cottreau, Francisco Barranco, and Timothée Masquelier. Optical flow estimation from event-based cameras and spiking neural networks. *Front. Neurosci.*, 17, 2023. 3
- [7] Pierre de Tournemire, Davide Nitti, Etienne Perot, Davide Migliore, and Amos Sironi. A large scale event-based detection dataset for automotive. *arXiv*, abs/2001.08499, 2020. 1, 5, 6, 7
- [8] Davide Falanga, Kevin Kleber, and Davide Scaramuzza. Dynamic obstacle avoidance for quadrotors with event cameras. *Science Robotics*, 2020. 2
- [9] Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi, Stefan Leutenegger, Andrew J. Davison, Jörg Conradt, Kostas Daniilidis, and Davide Scaramuzza. Event-based vision: A survey. *IEEE TPAMI*, 2020. 1
- [10] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430*, 2021. 5, 6
- [11] Daniel Gehrig and Davide Scaramuzza. Pushing the limits of asynchronous graph-based object detection with event cameras, 2022. 2, 3
- [12] Mathias Gehrig and Davide Scaramuzza. Recurrent vision transformers for object detection with event cameras. In *CVPR*, pages 13884–13893, 2023. 1, 2, 3, 4, 5, 6, 7, 8
- [13] Mathias Gehrig, Willem Aarents, Daniel Gehrig, and Davide Scaramuzza. DSEC: A stereo event camera dataset for driving scenarios. *IEEE RA-L*, 2021. 5
- [14] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. *NeurIPS*, 33, 2020. 2
- [15] Albert Gu, Karan Goel, Ankit Gupta, and Christopher Ré. On the parameterization and initialization of diagonal state space models. In *NeurIPS*, pages 35971–35983. Curran Associates, Inc., 2022. 3, 4, 5, 7, 8
- [16] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *ICLR*, 2022. 1, 2, 3, 4, 5, 7
- [17] Massimiliano Iacono, Stefan Weber, Arren Glover, and Chiara Bartolozzi. Towards event-driven object detection with off-the-shelf deep learning. In *IEEE/RSJ Int. Conf. Intell. Robot. Syst.*, pages 1–9, 2018. 3, 6
- [18] Kamil Jeziorek, Andrea Pinna, and Tomasz Kryjak. Memory-efficient graph convolutional networks for object classification and detection with event cameras. In *Sign. Proc.: Algo., Arch., Arrang., and Appl.*, pages 160–165, 2023. 3
- [19] Zhuangyi Jiang, Pengfei Xia, Kai Huang, Walter Stechele, Guang Chen, Zhenshan Bing, and Alois Knoll. Mixed frame-/event-driven fast pedestrian detection. In *IEEE Int. Conf. Robot. Autom.*, pages 8332–8338, 2019. 3, 6
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015. 5
- [21] Chuyi Li, Lulu Li, Hongliang Jiang, Kaiheng Weng, Yifei Geng, Liang Li, Zaidan Ke, Qingyuan Li, Meng Cheng, Weiqiang Nie, Yiduo Li, Bo Zhang, Yufei Liang, Linyuan Zhou, Xiaoming Xu, Xiangxiang Chu, Xiaoming Wei, and Xiaolin Wei. Yolov6: A single-stage object detection framework for industrial applications, 2022. 6
- [22] Dianze Li, Yonghong Tian, and Jianing Li. Sodformer: Streaming object detection with transformer using events and frames. *IEEE TPAMI*, 45(11):14020–14037, 2023. 2
- [23] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C. Berg. Ssd: Single shot multibox detector. In *ECCV*, pages 21–37, 2016. 6
- [24] Xu Liu, Jianing Li, Xiaopeng Fan, and Yonghong Tian. Event-based monocular dense depth estimation with recurrent transformers. *arxiv*, abs/2212.02791, 2022. 2, 3
- [25] Ze Liu, Han Hu, Yutong Lin, Zhuliang Yao, Zhenda Xie, Yixuan Wei, Jia Ning, Yue Cao, Zheng Zhang, Li Dong, Furu Wei, and Baining Guo. Swin transformer v2: Scaling up capacity and resolution. In *CVPR*, pages 12009–12019, 2022. 6
- [26] Zhuang Liu, Hanzi Mao, Chao-Yuan Wu, Christoph Feichtenhofer, Trevor Darrell, and Saining Xie. A convnet for the 2020s. *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 8
- [27] Nico Messikommer, Daniel Gehrig, Antonio Loquercio, and Davide Scaramuzza. Event-based asynchronous sparse convolutional networks. In *ECCV*, 2020. 3, 6
- [28] Yansong Peng, Yueyi Zhang, Zhiwei Xiong, Xiaoyan Sun, and Feng Wu. Get: Group event transformer for event-based vision. In *ICCV*, pages 6038–6048, 2023. 1, 6
- [29] Etienne Perot, Pierre de Tournemire, Davide Nitti, Jonathan Masci, and Amos Sironi. Learning to detect objects with a 1 megapixel event camera. In *NeurIPS*, pages 16639–16652. Curran Associates, Inc., 2020. 1, 3, 5, 6
- [30] Joseph Redmon and Ali Farhadi. Yolov3: An incremental improvement. In *CVPR*, 2018. 6
- [31] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 6
- [32] David W. Romero, Robert-Jan Brintjes, Jakub Mikolaj Tomczak, Erik J Bekkers, Mark Hoogendoorn, and Jan van

- Gemert. Flexconv: Continuous kernel convolutions with differentiable kernel sizes. In *ICLR*, 2022. 4
- [33] David W. Romero, Anna Kuzina, Erik J Bekkers, Jakub Mikołaj Tomczak, and Mark Hoogendoorn. CKConv: Continuous kernel convolution for sequential data. In *ICLR*, 2022. 4
- [34] Nikolaus Salvatore and Justin Fletcher. Dynamic vision-based satellite detection: A time-based encoding approach with spiking neural networks. In *Int. Conf. on Comput. Vis. Syst.*, pages 285–298, 2023. 3
- [35] Simon Schaefer, Daniel Gehrig, and Davide Scaramuzza. Aegnn: Asynchronous event-based graph neural networks. In *CVPR*, pages 12371–12381, 2022. 6
- [36] Jimmy T.H. Smith, Andrew Warrington, and Scott Linderman. Simplified state space layers for sequence modeling. In *ICLR*, 2023. 1, 2, 3, 4, 5, 6, 7, 8
- [37] Leslie N. Smith and Nicholay Topin. Super-convergence: Very fast training of residual networks using large learning rates. In *ICLR*, 2018. 5
- [38] Daobo Sun and Haibo Ji. Event-based object detection using graph neural networks. In *IEEE Data Driven Contr. and Learn. Syst. Conf.*, pages 1895–1900, 2023. 3
- [39] Stepan Tulyakov, Daniel Gehrig, Stamatios Georgoulis, Julius Erbach, Mathias Gehrig, Yuanyou Li, and Davide Scaramuzza. Time lens: Event-based video frame interpolation. In *CVPR*, pages 16155–16164, 2021. 2
- [40] Ziyi Wu, Xudong Liu, and Igor Gilitschenski. Eventclip: Adapting clip for event-based object recognition. *ArXiv*, 2023. 2
- [41] Zizhao Zhang, Han Zhang, Long Zhao, Ting Chen, , Serkan Ö. Arik, and Tomas Pfister. Nested hierarchical transformer: Towards accurate, data-efficient and interpretable visual understanding. In *AAAI*, 2022. 6
- [42] Nikola Zubić and Pietro Liò. An effective loss function for generating 3d models from single 2d image without rendering. pages 309–322. Springer International Publishing, 2021. 1
- [43] Nikola Zubić, Daniel Gehrig, Mathias Gehrig, and Davide Scaramuzza. From chaos comes order: Ordering event representations for object recognition and detection. In *ICCV*, pages 12846–12856, 2023. 2, 3, 6