



南京航空航天大学  
Nanjing University of Aeronautics and Astronautics

---

# Event Transformer

---



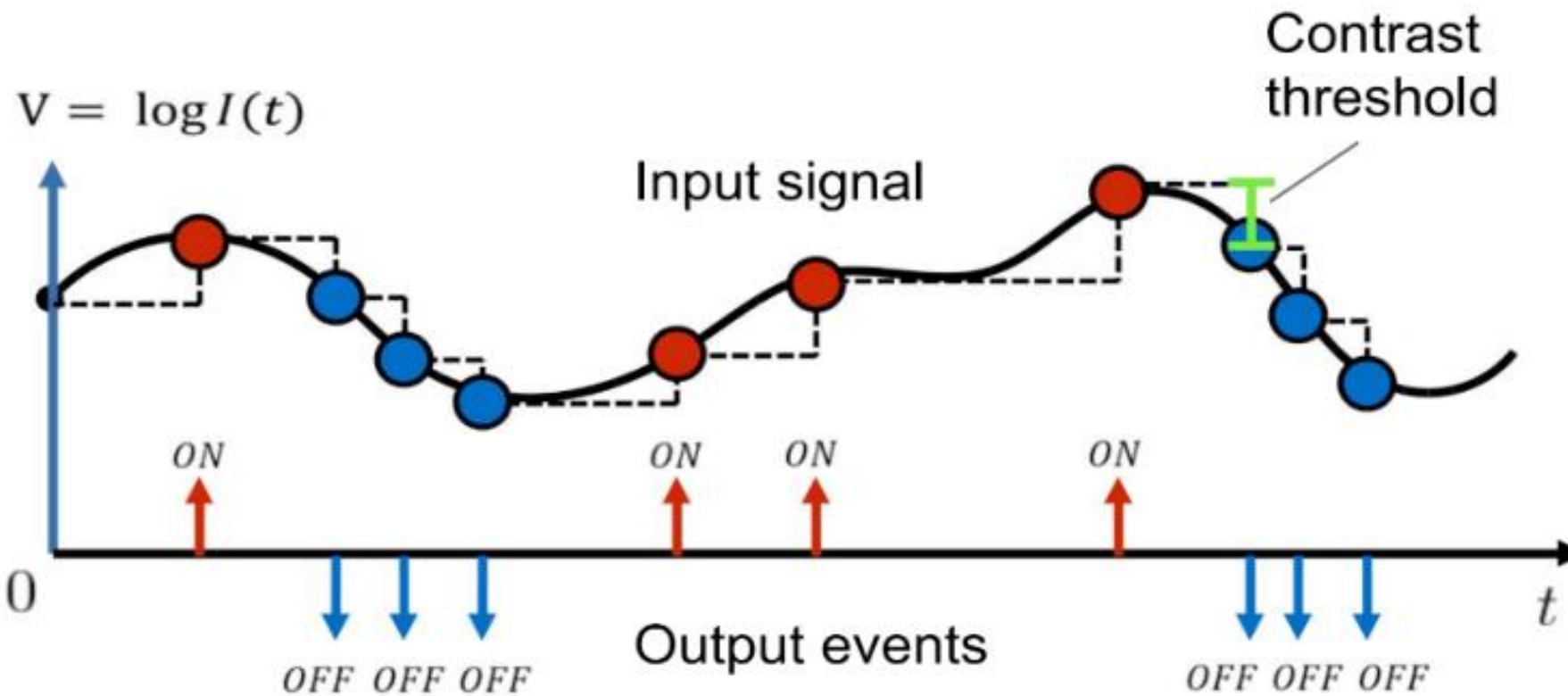
# Event Transformer

Zhihao Li<sup>1</sup>, M. Salman Asif<sup>2</sup>, and Zhan Ma<sup>1</sup>

<sup>1</sup> Vision Lab, Nanjing University, China

<sup>2</sup> University of California at Riverside, USA

`lizhihao6@smail.nju.edu.cn`, `sasif@ece.ucr.edu`,  
`mazhan@nju.edu.cn`

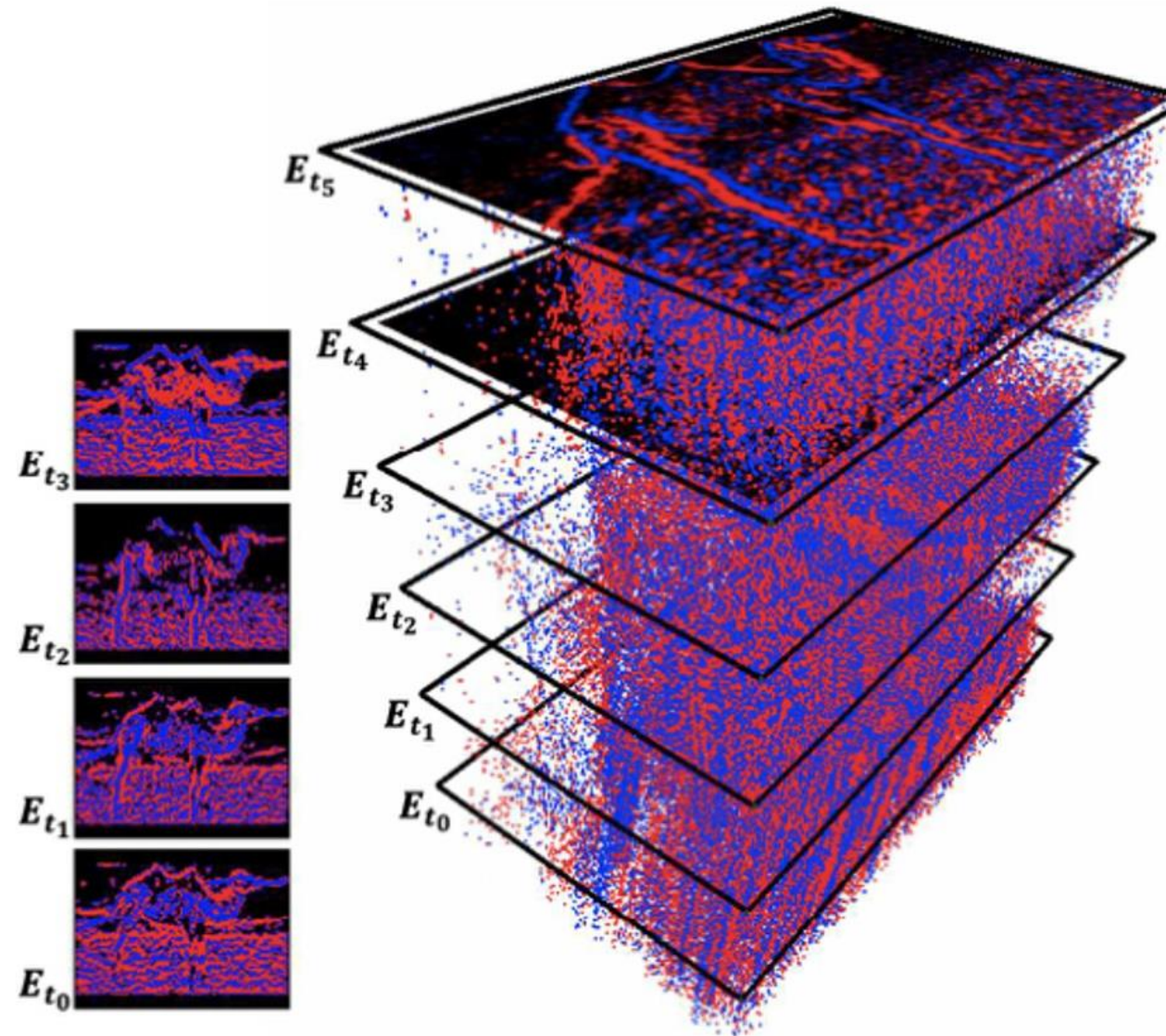


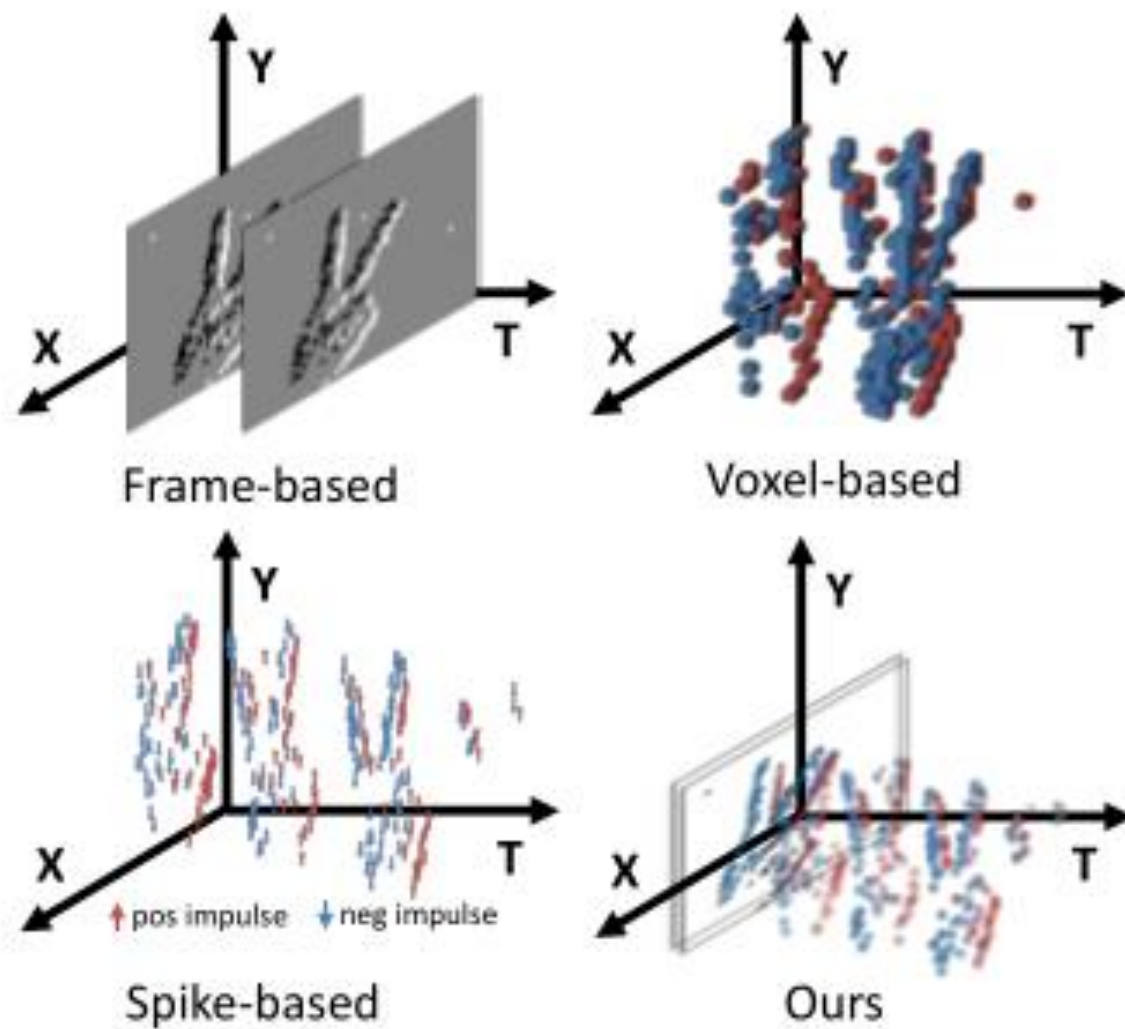
$$\log I(x, y, t) - \log I(x, y, t - \Delta t) = \pm C$$

$$\mathbf{e}_k \triangleq (\mathbf{x}_k, t_k, p_k)$$

Event set:

$$\mathcal{E} = \{\mathbf{e}_k\}_{k=1}^N$$





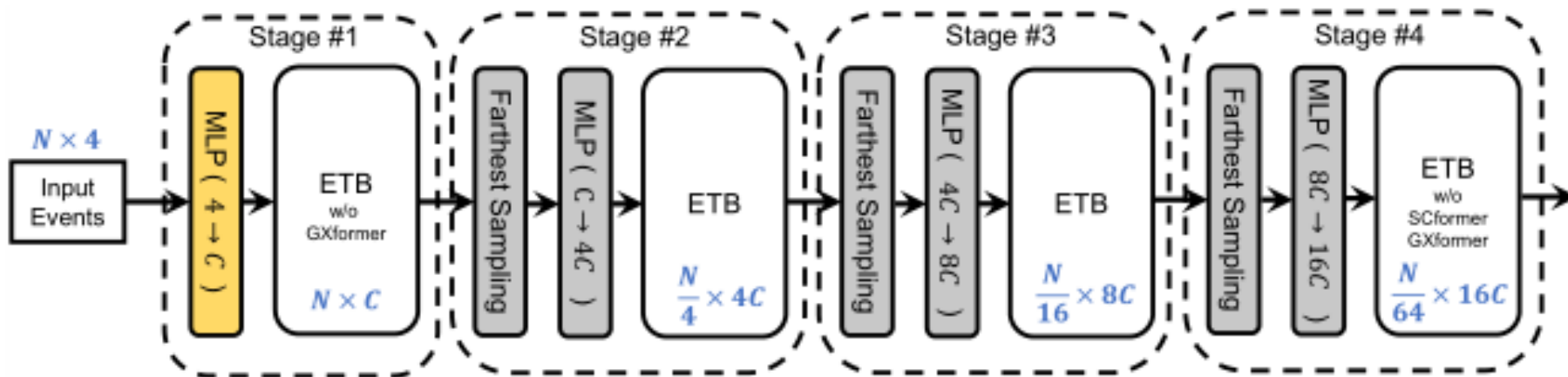
(a) Representation Forms

## merit

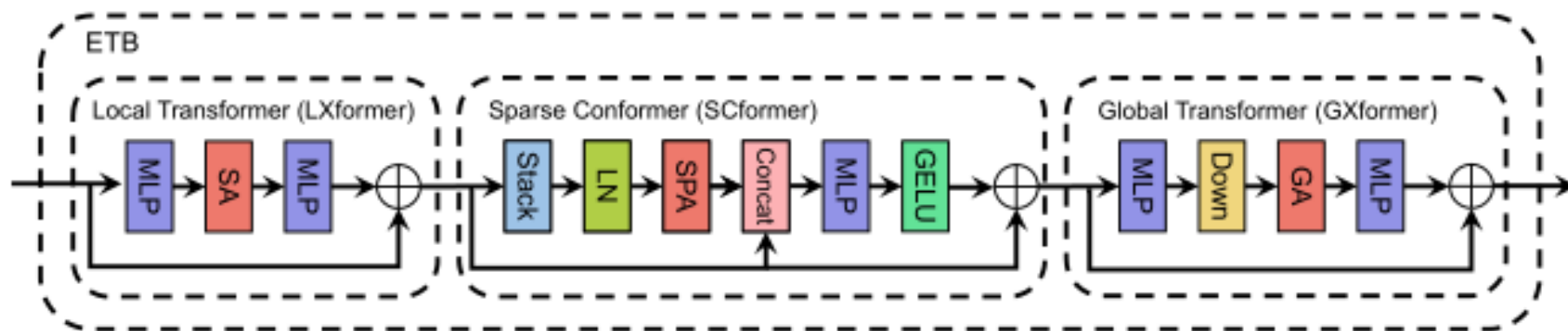
1. Asynchronous
  - a. Pixel-wise independent
  - b. Brightness change
  - c. Low bandwidth
2. Low latency
  - a. No motion blur
3. High dynamic range (120dB)
4. Low power (mW)

## shortcoming

1. Data is asynchronous and not easy to process
2. Lack of texture information about the object



(a) Event Transformer Backbone



(b) Event Transformer Block

1. The input set has  $N$  points, select a point  $p_0$  from the set as the starting point, and obtain the sampling point set  $S=\{P_0\}$ ;
2. Calculate the distance from all points to  $p_0$ , form an  $N$ -dimensional array  $L$ , select the point corresponding to the maximum value as  $p_1$ , and update the sampling point set  $S=\{P_0, P_1\}$ ;
3. calculate the distance from all points to  $p_1$ , for each point  $p_i$ , if its distance from  $p_1$  is less than  $L[i]$ , then update  $L[i] = d(P_i, P_1)$ , therefore, the array  $L$  is stored in the nearest distance from each point to the sample point set  $S$ ;
4. Select the point corresponding to the maximum value in  $L$  as  $P_2$ , and update the sampling point set  $S=\{P_0, P_1, P_2\}$ ;
5. Repeat steps 2-4 until  $N$  target sampling points.

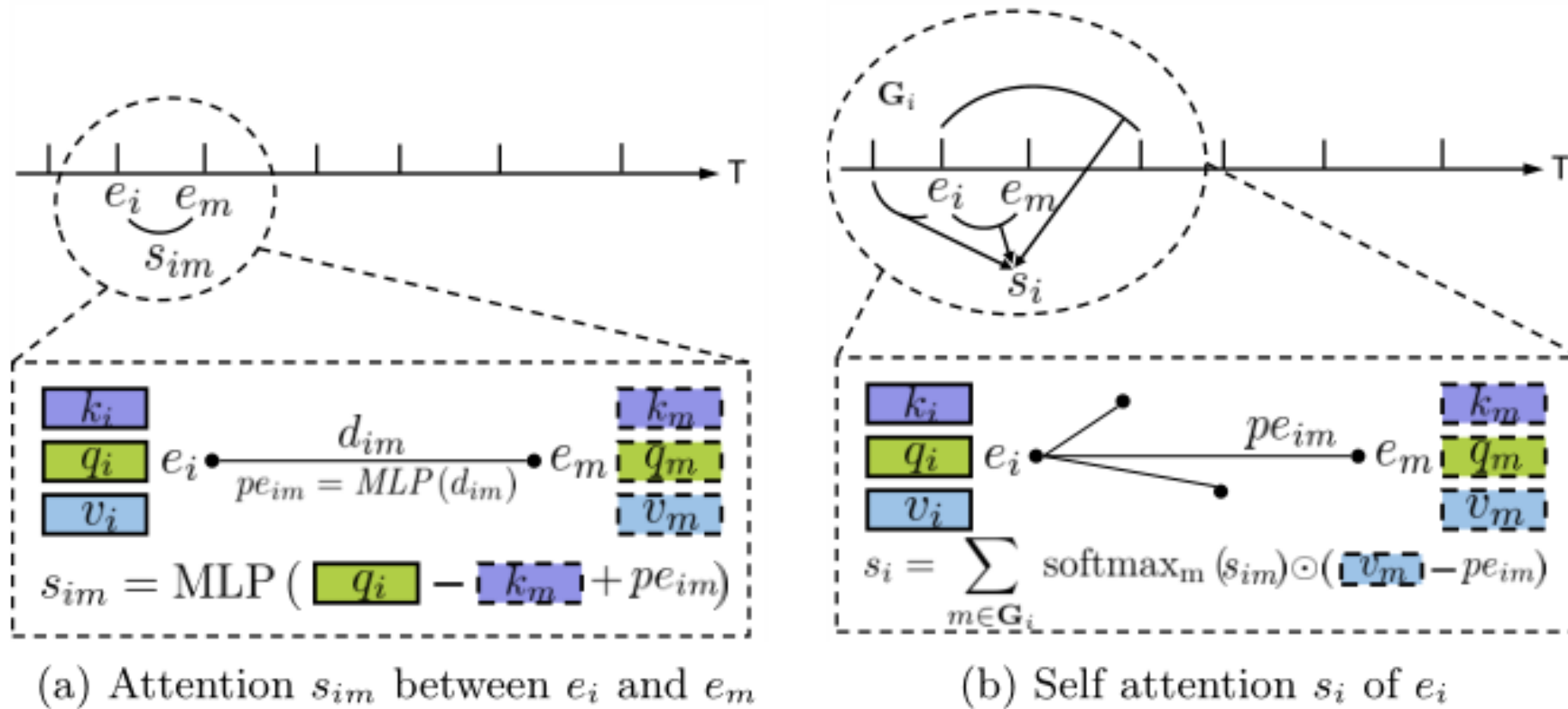


Fig. 4: **Self attention in LXformer.** (a) The attention  $s_{im}$  between  $e_i$  and  $e_m$  is calculated using query  $q_i$  of  $e_i$ , key  $k_m$  of  $e_m$  and the distance  $d_{im}$  between  $e_i$  and  $e_j$ . (b) The self-attention  $s_i$  is calculated using the set  $G_i$  of events that are temporally close to event  $e_i$ .

$$\mathbf{E} = \{e_1, \dots, e_N\} \in R^{N \times 4} \quad \mathbf{F} = \{f_1, \dots, f_N\} \in R^{N \times C} \quad \text{Q, K, and V}$$

$$\mathbf{D} = \{d_{11}, \dots, d_{ij}, d_{NN}\} \in R^{N \times N \times 4} \quad \text{PE} = \text{MLP}(\mathbf{D})$$

$$s_{im} = \text{MLP}(q_i - k_m + pe_{im}), \quad (2)$$

$$s_i = \sum_{m \in \mathbf{G}_i} \text{softmax}_m(s_{im}) \odot (v_m + pe_{im}). \quad (3)$$

$$y_i = f_i + \text{MLP}(s_i). \quad (4)$$

The sequence of events  $\{(e_i, f_i)\}_{i=1}^N$  is converted to a 2D frame  $I \in \mathbb{R}^{H \times W \times (2+C)}$

$$L = \text{LN}(I) \quad Q = \text{SPConv}_q(L), K = \text{SPConv}_k(L), V = \text{SPConv}_v(L)$$

$$\text{event } e_i \text{ and event } e_j \in W \quad s_{ij} = \text{MLP}(q_{y_i, x_i} - k_{y_j, x_j} + pe_{ij}), \quad (5)$$

$$pe_{ij} = \text{MLP}((y_i, x_i, p_i) - (y_j, x_j, p_j))$$

$$s_i = \sum_{j \in W_i} \text{softmax}_j(s_{ij}) \odot (v_j + pe_{ij}). \quad (6)$$

$$y_i = f_i + \phi(\text{concat}(f_i, \text{MLP}(s_i))). \quad (7)$$

The farthest point down-sampling with down-sample rate  $r \rightarrow \hat{E}$ , for each  $e_j \in \hat{E}$ , its corresponding feature  $\hat{f}_j \in \hat{F}$  is calculated as

$$\hat{f}_j = \max_k (\text{MLP} (\text{concat} (e_k, f_k))), \quad (8)$$

$$s_i = \sum_{j \in \hat{E}} \text{softmax}_j (s_{ij}) \odot (\hat{v}_j + p e_{ij}), \quad (9)$$

$$y_i = f_i + \text{MLP} (s_i). \quad (10)$$

Table 2: Accuracy of Event-based Classification Using Various Datasets. N-M refers to N-MNIST, N-Cal to N-Caltech101, CIF10 to CIFAR10-DVS, N-C to N-Cars, and ASL to ASL-DVS.

Method	Type	N-M [21]	N-Cal [33]	CIF10 [27]	N-C [39]	ASL [3]
H-First [34]	frame	0.712	0.054	0.077	0.561	-
HOTS [24]	frame	0.808	0.210	0.271	0.624	-
HATS [39]	frame	0.991	0.642	0.524	0.902	-
EST [18]	frame	0.990	0.753	0.634	0.919	0.979
AMAE [11]	frame	0.983	0.694	0.620	0.936	0.984
M-LSTM [5]	frame	0.986	0.738	0.631	0.927	0.980
MVF-Net [10]	frame	0.981	0.687	0.599	0.927	0.971
ResNet-50 [20]	frame	0.984	0.637	0.558	0.903	0.886
EventNet [37]	voxel	0.752	0.425	0.171	0.750	0.833
RG-CNNs [3]	voxel	0.990	0.657	0.540	0.914	0.901
EV-VGCNN [9]	voxel	0.994	0.748	0.651	0.953	0.983
VMV-GCN [43]	voxel	0.995	0.778	0.663	0.932	0.989
Gabor-SNN [26]	spike	0.837	0.196	0.245	0.789	-
PLIF [15]	spike	0.996	-	0.697	-	-
<b>Ours</b>	tensor	<b>0.999</b>	<b>0.789</b>	<b>0.712</b>	<b>0.954</b>	<b>0.999</b>

Table 3: Model Complexity of Different Methods

Method	Type	Parameters	FLOPs	Top-1 Accuracy
EST [18]	frame	21.38M	4.28G	52.4
M-LSTM [5]	frame	21.43M	4.82G	63.1
MVF-Net [10]	frame	33.62M	5.62G	59.9
ResNet-50 [20]	frame	25.61M	3.87G	55.8
EventNet [37]	voxel	2.81M	0.91G	17.1
RG-CNNs [3]	voxel	19.46M	0.79G	54.0
EV-VGCNN [9]	voxel	0.84M	0.70G	65.1
VMV-GCN [43]	voxel	0.86M	1.30G	66.3
<b>Ours</b>	tensor	15.87M	0.51G	71.2

Table 6: Ablation Study of Closest Events in LXformer

<i>M</i> Closest Events in LXformer	4	8	<b>16</b>	32	64
Top-1 Accuracy	65.8	69.0	<b>71.2</b>	70.9	69.6

Table 7: Ablation Study of Window Size in SCformer

Size of <i>W</i> in SCformer	$1 \times 1$	$3 \times 3$	$5 \times 5$	$7 \times 7$	$9 \times 9$
Top-1 Accuracy	69.9	<b>71.2</b>	71.1	71.0	70.6

Table 8: Ablation Study of Down-sample Rate in GXformer

Down-sample rate <i>r</i> in GXformer	1	8	<b>32</b>	64	128
Top-1 Accuracy	69.8	70.3	<b>71.2</b>	70.6	70.1



南京航空航天大学

Nanjing University of Aeronautics and Astronautics

# Thanks for Listening

