

TRAINING ROBUST SPIKING NEURAL NETWORKS WITH VIEWPOINT TRANSFORM AND SPATIOTEMPORAL STRETCHING

Haibo Shen¹, Juyu Xiao¹, Yihao Luo^{2,1}, Xiang Cao^{3,1}, Liangqi Zhang¹, Tianjiang Wang¹

School of Huazhong University of Science and Technology¹
Yichang Testing Technique Research Institute²
Changsha University³

ABSTRACT

Neuromorphic vision sensors (event cameras) simulate biological visual perception systems and have the advantages of high temporal resolution, less data redundancy, low power consumption, and large dynamic range. Since both events and spikes are modeled from neural signals, event cameras are inherently suitable for spiking neural networks (SNNs), which are considered promising models for artificial intelligence (AI) and theoretical neuroscience. However, the unconventional visual signals of these cameras pose a great challenge to the robustness of spiking neural networks. In this paper, we propose a novel data augmentation method, ViewPoint Transform and SpatioTemporal Stretching (VPT-STs). It improves the robustness of SNNs by transforming the rotation centers and angles in the spatiotemporal domain to generate samples from different viewpoints. Furthermore, we introduce the spatiotemporal stretching to avoid potential information loss in viewpoint transformation. Extensive experiments on prevailing neuromorphic datasets demonstrate that VPT-STs is broadly effective on multi-event representations and significantly outperforms pure spatial geometric transformations. Notably, the SNNs model with VPT-STs achieves a state-of-the-art accuracy of 84.4% on the DVS-CIFAR10 dataset.

Index Terms— Spiking Neural Networks, Neuromorphic Data, Data Augmentation, ViewPoint Transform and SpatioTemporal Stretching

1. INTRODUCTION

Inspired by the primate visual system, neuromorphic vision cameras generate events by sampling the brightness of objects. For example, the Dynamic Vision Sensor (DVS) [1] camera and the Vidar [2] camera are inspired by the outer three-layer structure of the retina and the foveal three-layer structure, respectively. Both of them have the advantages of

high temporal resolution, less data redundancy, low power consumption, and large dynamic range [3]. In addition, spiking neural networks (SNNs) are similarly inspired by the learning mechanisms of the mammalian brain and are considered a promising model for artificial intelligence (AI) and theoretical neuroscience [4]. In theory, as the third generation of neural networks, SNNs are computationally more powerful than traditional convolutional neural networks (CNNs) [4]. Therefore, event cameras are inherently suitable for SNNs.

However, the unconventional visual signals of these cameras also pose a great challenge to the robustness of SNNs. Most existing data augmentations are fundamentally designed for RGB data and lack exploration of neuromorphic events. For example, Cutout [5] artificially impedes a rectangular block in the image to simulate the impact of occlusion on the image. Random erasing [6] further optimizes the erased pixel value by adding noise. Mixup [7] uses the weighted sum of two images as training samples to smooth the transition line between classes. Since neuromorphic data have an additional temporal dimension and differ widely in imaging principles, novel data augmentations are required to process the spatiotemporal visual signals of these cameras.

In this paper, we propose a novel data augmentation method suitable for events, ViewPoint Transformation and SpatioTemporal Stretching (VPT-STs). Viewpoint transformation solves the spatiotemporal scale mismatch of samples by introducing a balance coefficient, and generates samples from different viewpoints by transforming the rotation centers and angles in the spatiotemporal domain. Furthermore, we introduce spatiotemporal stretching to avoid potential information loss in viewpoint transformation. Extensive experiments are performed on prevailing neuromorphic datasets. It turns out that VPT-STs is broadly effective on multiple event representations and significantly outperforms pure spatial geometric transformations. Insightful analysis shows that VPT-STs improves the robustness of SNNs against different spatial locations. In particular, the SNNs model with VPT-STs achieves a state-of-the-art accuracy of 84.4% on the DVS-CIFAR10 dataset.

Furthermore, while this work is related to EventDrop [8],

This work was supported in part by the National Natural Science Foundation of China under Grant 61572214 and Seed Foundation of Huazhong University of Science and Technology (2020kfyXGYJ114). (Corresponding author: Tianjiang Wang.)

NDA [9], there are some notable differences. For example, NDA is a pure global geometric transformation, while VPT-STS changes the viewpoint of samples in the spatiotemporal domain. EventDrop is only experimented on CNNs, it introduces noise by dropping events, but may cause problems with dead neurons on SNNs. VPT-STS is applicable to both CNNs and SNNs, maintaining the continuity of samples. In addition, EventDrop transforms both temporal and spatial domains, but as two independent strategies, it does not combine the spatiotemporal information of the samples. To our knowledge, VPT-STS is the first event data augmentation that simultaneously incorporates spatiotemporal transformations.

2. METHOD

2.1. Event Generation Model

The event generation model [3, 4] is abstracted from dynamic vision sensors [1]. Each pixel of the event camera responds to changes in its logarithmic photocurrent $L = \log(I)$. Specifically, in a noise-free scenario, an event $e_k = (x_k, y_k, t_k, p_k)$ is triggered at pixel $X_k = (y_k, x_k)$ and at time t_k as soon as the brightness variation $|\Delta L|$ reaches a temporal contrast threshold C since the last event at the pixel. The event generation model can be expressed by the following formula:

$$\Delta L(X_k, t_k) = L(X_k, t_k) - L(X_k, t_k - \Delta t_k) = p_k C \quad (1)$$

where $C > 0$, Δt_k is the time elapsed since the last event at the same pixel, and the polarity $p_k \in \{+1, -1\}$ is the sign of the brightness change. During a period, the event camera triggers event stream \mathcal{E} :

$$\mathcal{E} = \{e_k\}_{k=1}^N = \{(X_k, t_k, p_k)\}_{k=1}^N \quad (2)$$

where N represents the number of events in the set \mathcal{E} .

As shown in Figure 1, an event is generated each time the brightness variances reach the threshold, and then $|\Delta L|$ is cleared. The event stream can be represented as a matrix:

$$M_{\mathcal{E}} = \begin{pmatrix} y_1 & x_1 & t_1 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ y_N & x_N & t_N & 1 \end{pmatrix}_{4 \times N} \quad (3)$$

For convenience, we omit the unconverted polarity p .

2.2. Motivation

This work stems from the observation that it is difficult to maintain absolute frontal view between the sample and cameras, which easily leads to a slight shift of the viewpoint. Considering this small offset distance, we use viewpoint rotation to approximate the deformation of samples in space and time. In addition, since events record the brightness change of samples, especially changes of the edge, variations of the illumination angle will also cause the effect of viewpoint transformation, which suggests that we can enhance the robustness of SNNs by generating viewpoint-transformed samples.

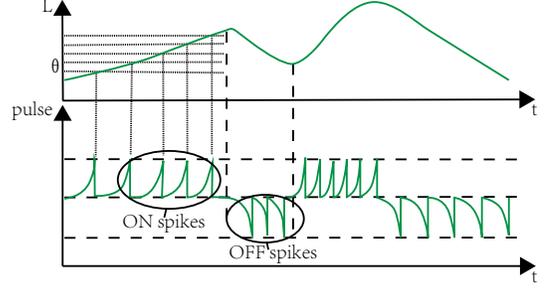


Fig. 1. Event generation model.

2.3. The Proposed Method.

To generate viewpoint-transformed samples, we draw on the idea of spatio-temporal rotation. For viewpoint transformation (VPT), we introduce translation matrices T_b , T_a , which represent the translation to the rotation center (x_c, y_c, t_c) and the translation back to the original position, respectively.

$$T_b = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ -y_c & -x_c & -t_c & 1 \end{pmatrix}, T_a = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ y_c & x_c & t_c & 1 \end{pmatrix} \quad (4)$$

Suppose that rotate along the y and t planes with x as the axis, we can easily derive the rotation matrix R_r^{YT} :

$$R_r^{YT} = \begin{pmatrix} \cos\theta & 0 & \sin\theta & 0 \\ 0 & 1 & 0 & 0 \\ -\sin\theta & 0 & \cos\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (5)$$

where θ is the rotation angle. In practice, Eq 5 is an unbalanced matrix due to the mismatch between the time and space dimensions in the $M_{\mathcal{E}}$ matrix. Therefore, we introduce a balance coefficient τ to scale the space and time dimension, which results in a better visual effects. The balanced matrix R_{br}^{YT} can be formulated as:

$$R_{br}^{YT} = \begin{pmatrix} \cos\theta & 0 & \tau \sin\theta & 0 \\ 0 & 1 & 0 & 0 \\ -\frac{1}{\tau} \sin\theta & 0 & \cos\theta & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} \quad (6)$$

Set $x_c = 0$, the viewpoint transformation matrix M_{br}^{YT} can be formulated by calculating $T_b R_{br}^{YT} T_a$:

$$\begin{pmatrix} \cos\theta & 0 & \tau \sin\theta & 0 \\ 0 & 1 & 0 & 0 \\ -\frac{1}{\tau} \sin\theta & 0 & \cos\theta & 0 \\ -x_c \cos\theta + \frac{1}{\tau} t_c \sin\theta + x_c & 0 & -\tau x_c \sin\theta - t_c \cos\theta + t_c & 1 \end{pmatrix} \quad (7)$$

Similarly, the viewpoint transformation matrix M_{br}^{XT} in

Table 1. Performance of VPT-STs on SNNs and CNNs with various representations.

Datasets	Method	Accuracy (%)				
		SNNs	EventFrame	EventCount	VoxelGrid	EST
CIFAR10-DVS	Baseline	83.20	78.71	78.85	77.47	78.81
	VPT-STs	84.40	79.58	79.12	79.62	79.37
N-Caltech101	Baseline	78.98	73.08	73.66	77.08	78.41
	VPT-STs	81.05	76.96	76.38	79.13	78.88
N-CARS	Baseline	95.40	94.44	94.76	93.86	94.97
	VPT-STs	95.85	94.60	94.81	94.30	94.99

the x and t dimensions can be formulated as:

$$\begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & \cos\theta & \tau \sin\theta & 0 \\ 0 & -\frac{1}{\tau} \sin\theta & \cos\theta & 0 \\ 0 & -x_c \cos\theta + \frac{1}{\tau} t_c \sin\theta + x_c & -\tau x_c \sin\theta - t_c \cos\theta + t_c & 1 \end{pmatrix} \quad (8)$$

Therefore, the viewpoint-transformed matrix M_{VPT}^{YT} and M_{VPT}^{XT} can be formulated as:

$$\left. \begin{aligned} M_{VPT}^{YT} &= M_\varepsilon M_{br}^{YT} \\ M_{VPT}^{XT} &= M_\varepsilon M_{br}^{XT} \end{aligned} \right\} \quad (9)$$

Furthermore, since events beyond the resolution will be discarded during the viewpoint transformation, we introduce spatiotemporal stretching (STS) to avoid potential information loss. STS stretches the temporal mapping in the VPT by a coefficient $\frac{1}{\cos\theta}$ while maintaining the spatial coordinates unchanged. Therefore, by setting $t_c = 0$, we get the transformed $(t)_{STS}^{YT}$ and $(t)_{STS}^{XT}$ from Eq. 7 and Eq. 8:

$$\left. \begin{aligned} (t_k)_{VPT}^{YT} &= (t_k) - \tau \tan\theta \cdot ((y_k) - y_c) \\ (t_k)_{VPT}^{XT} &= (t_k) - \tau \tan\theta \cdot ((x_k) - x_c) \end{aligned} \right\} \quad (10)$$

The time of STS is advanced or delayed according to the distance from the center $|x - x_c|$ ($|y - y_c|$), causing event stream to be stretched long the time axis according to the spatial coordinates.

3. EXPERIMENTS

3.1. Implementation

Extensive experiments are performed to demonstrate the superiority of the VPT-STs method on prevailing neuromorphic datasets, including CIFAR10-DVS(CIF-DVS) [10], N-Caltech101(N-Cal) [11], N-CARS [12] datasets. N-Caltech101 and CIFAR10-DVS datasets are generated by neuromorphic vision sensors on the basis of traditional datasets, while N-CARS is collected in the real world. For the convenience of comparison, the model without VPT-STs with the same parameters is used as the baseline. STBP [13] methods are used to train SNN-VGG9 network, other parameters

Table 2. Performance of VPT-STs and previous SOTAs on CIFAR10-DVS and N-CARS datasets.

Methods	References	Accuracy (%)	
		CIF-DVS	N-CARS
HATS[12]	CVPR 2018	52.40	81.0
Dart[19]	TPAMI 2020	65.80	-
Dspike [20]	NeurIPS 2021	75.40	-
STBP [13]	AAAI 2021	67.80	-
AutoSNN [21]	ICML 2022	72.50	-
RecDis [22]	CVPR 2022	72.42	-
DSR [23]	CVPR 2022	77.27	-
NDA [9]	ECCV 2022	81.70	90.1
VPT-STs	-	84.40	95.85

mainly refer to NDA [14]. For example, the Adam optimizer is used with an initial learning rate of $1e - 3$. The neuron threshold and leakage coefficient are 1 and 0.5, respectively. In addition, we also evaluate the performance of VPT-STs on various event representations with the Resnet9 network, including EST [15], VoxelGrid [16], EventFrame [17] and EventCount [18] representations.

3.2. Performance on various representations

Extensive experiments are conducted to evaluate the performance of VPT-STs method on different event representations, covering SNNs and CNNs. As shown in Tab. 1, SNNs with VPT-STs methods achieve significant improvements on three prevailing datasets. And VPT-STs also performs well on four representations commonly used by CNNs. It is worth noting that EST maintains the most spatiotemporal information from neuromorphic data and thus performs best overall. Furthermore, since the samples of N-CARS are collected in the real world, its initial viewpoint diversity is already enriched compared to the other two datasets. Considering the high baseline on N-CARS, VPT-STs still further improves the robustness of SNNs.

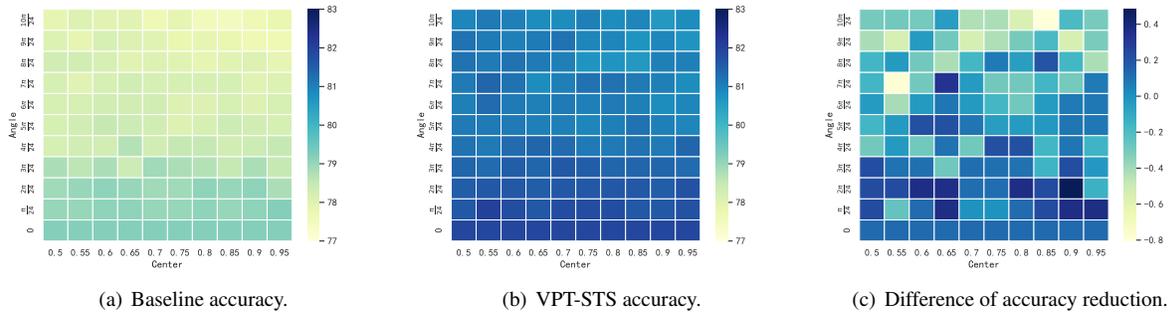


Fig. 2. Performance of VPT-STs and Baseline under different perturbations.

3.3. Compared with SOTAs

As shown in Tab. 2, we compare VPT-STs with recent state-of-the-art results on neuromorphic datasets. The results show that VPT-STs achieves substantial improvements over previous SOTAs. It is worth noting that VPT-STs significantly outperforms NDA, which is an ensemble of six geometric transformations. The experimental results demonstrate the superiority of combining spatiotemporal information for data augmentation. Since VPT-STs is orthogonal to most training algorithms, it can provide a better baseline and improve the performance of existing models.

3.4. Ablation Studies on VPT-STs

As shown in Fig. 3, the performance of VPT-STs with different rotation angles is evaluated on the N-Caltech101 dataset. It turns out that a suitable rotation angle is important for the performance of data augmentation, which can increase data diversity without losing features.

3.5. Analysis of VPT-STs

To gain further insight into the workings of VPT-STs, we add different strategies on the baseline to analyze the effective components of VPT-STs. As shown in Table 3, spatial rotation (Rotation) is performed as a comparative experiment for VPT-STs. It turns out that both VPT and STs including spatiotemporal transformations are significantly better than pure spatial geometric transformations on all three datasets, which illustrate the importance of spatiotemporal transformations.

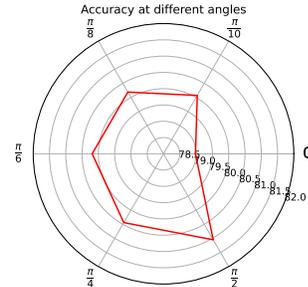


Fig. 3. Performance of VPT-STs at different angles.

While VPT and STs are implemented with operations similar to rotation, it actually improves the robustness of SNNs to different viewpoints. Furthermore, we evaluate the robustness of SNNs to viewpoint fluctuations by adding different degrees of spatiotemporal rotation to the test data. Figures 2(a) and 2(b) show the performance of the baseline model and the model trained by VPT-STs under different disturbances, respectively. The results show that the general trend of the accuracy change is to decrease with the increase of the perturbation amplitude. In addition, Fig. 2(c) shows the difference in the accuracy reduction of the VPT-STs compared to baseline. As the perturbation amplitude increases, the difference in the accuracy reduction of the two models is less than zero, and the absolute value grows, which illustrate that the accuracy reduction of baseline is larger than that of VPT-STs. Experimental results show that the model trained with VPT-STs generalize better and improves the robustness of SNNs against spatial location variances.

Table 3. Comparison of Different Strategies.

Methods	Accuracy (%)		
	CIF-DVS	N-Cal	N-CARS
Baseline	83.20	78.98	95.40
Rotation	83.90	80.19	95.46
VPT	84.40	81.05	95.56
STs	84.30	80.56	95.85

4. CONCLUSION

We propose a novel data augmentation method suitable for events, viewpoint transformation and spatiotemporal stretching (VPT-STs). Extensive experiments on prevailing neuromorphic datasets show that VPT-STs is broadly effective on multiple event representations and significantly outperforms pure spatial geometric transformations. It achieves substan-

tial improvements over previous SOTAs by improving the robustness of SNNs to different viewpoints.

5. REFERENCES

- [1] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbrück, “A 128 x 128 120db 30mw asynchronous vision sensor that responds to relative intensity change,” in *ISSCC*. IEEE, 2006, pp. 2060–2069.
- [2] Siwei Dong, Tiejun Huang, and Yonghong Tian, “Spike camera and its coding methods,” *arXiv preprint arXiv:2104.04669*, 2021.
- [3] Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Tabbara, Andrea Censi, Stefan Leutenegger, Andrew J. Davison, Jörg Conradt, Kostas Daniilidis, and Davide Scaramuzza, “Event-based vision: A survey,” *TPAMI*, pp. 154–180, 2022.
- [4] Kaushik Roy, Akhilesh Jaiswal, and Priyadarshini Panda, “Towards spike-based machine intelligence with neuromorphic computing,” *Nature*, vol. 575, pp. 607–617, 2019.
- [5] Terrance Devries and Graham W. Taylor, “Improved regularization of convolutional neural networks with cutout,” *CoRR*, vol. abs/1708.04552, 2017.
- [6] Zhun Zhong, Liang Zheng, Guoliang Kang, Shaozi Li, and Yi Yang, “Random erasing data augmentation,” in *AAAI*, 2020.
- [7] Hongyi Zhang, Moustapha Cissé, Yann N. Dauphin, and David Lopez-Paz, “mixup: Beyond empirical risk minimization,” in *ICLR*. 2018, OpenReview.net.
- [8] Fuqiang Gu, Weicong Sng, Xuke Hu, and Fangwen Yu, “Eventdrop: Data augmentation for event-based learning,” in *IJCAI*, Zhi-Hua Zhou, Ed. 2021, pp. 700–707, ijcai.org.
- [9] Yuhang Li, Youngeun Kim, Hyungseob Park, Tamar Geller, and Priyadarshini Panda, “Neuromorphic data augmentation for training spiking neural networks,” *arXiv preprint arXiv:2203.06145*, 2022.
- [10] Hongmin Li, Hanchao Liu, Xiangyang Ji, Guoqi Li, and Luping Shi, “Cifar10-dvs: an event-stream dataset for object classification,” *Frontiers in neuroscience*, vol. 11, pp. 309, 2017.
- [11] Garrick Orchard, Ajinkya Jayawant, Gregory K Cohen, and Nitish Thakor, “Converting static image datasets to spiking neuromorphic datasets using saccades,” *Frontiers in neuroscience*, vol. 9, pp. 437, 2015.
- [12] Amos Sironi, Manuele Brambilla, Nicolas Bourdis, Xavier Lagorce, and Ryad Benosman, “Hats: Histograms of averaged time surfaces for robust event-based object classification,” in *Proceedings of the IEEE*

Conference on Computer Vision and Pattern Recognition, 2018, pp. 1731–1740.

- [13] Hanle Zheng, Yujie Wu, Lei Deng, Yifan Hu, and Guoqi Li, “Going deeper with directly-trained larger spiking neural networks,” in *AAAI*, 2021, pp. 11062–11070.
- [14] Yuhang Li, Youngeun Kim, Hyoungeob Park, Tamar Geller, and Priyadarshini Panda, “Neuromorphic data augmentation for training spiking neural networks,” *arXiv preprint arXiv:2203.06145*, 2022.
- [15] Daniel Gehrig, Antonio Loquercio, Konstantinos G. Derpanis, and Davide Scaramuzza, “End-to-end learning of representations for asynchronous event-based data,” in *ICCV*. 2019, pp. 5632–5642, IEEE.
- [16] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis, “Unsupervised event-based learning of optical flow, depth, and egomotion,” in *CVPR*. 2019, pp. 989–997, Computer Vision Foundation / IEEE.
- [17] Henri Rebecq, Timo Horstschaefer, and Davide Scaramuzza, “Real-time visual-inertial odometry for event cameras using keyframe-based nonlinear optimization,” in *BMVC*. 2017, BMVA Press.
- [18] Ana I. Maqueda, Antonio Loquercio, Guillermo Gallego, Narciso García, and Davide Scaramuzza, “Event-based vision meets deep learning on steering prediction for self-driving cars,” in *CVPR*. 2018, pp. 5419–5427, Computer Vision Foundation / IEEE Computer Society.
- [19] Bharath Ramesh, Hong Yang, Garrick Orchard, and et al., “Dart: Distribution aware retinal transform for event-based cameras,” *TPAMI*, vol. 42, no. 11, pp. 2767–2780, 2020.
- [20] Yuhang Li, Yufei Guo, Shanghang Zhang, Shikuang Deng, Yongqing Hai, and Shi Gu, “Differentiable spike: Rethinking gradient-descent for training spiking neural networks,” *NeurIPS*, vol. 34, pp. 23426–23439, 2021.
- [21] Byunggook Na, Jisoo Mok, Seongsik Park, Dongjin Lee, Hyeokjun Choe, and Sungroh Yoon, “Autosnn: Towards energy-efficient spiking neural networks,” *arXiv preprint arXiv:2201.12738*, 2022.
- [22] Yufei Guo, Xinyi Tong, Yuanpei Chen, Liwen Zhang, Xiaode Liu, Zhe Ma, and Xuhui Huang, “Recdis-snn: Rectifying membrane potential distribution for directly training spiking neural networks,” in *CVPR*, 2022, pp. 326–335.
- [23] Qingyan Meng, Mingqing Xiao, Shen Yan, Yisen Wang, Zhouchen Lin, and Zhi-Quan Luo, “Training high-performance low-latency spiking neural networks by differentiation on spike representation,” in *CVPR*, 2022, pp. 12444–12453.