001	Event Trojan: Asynchronous Event-based	001
002	Backdoor Attacks	002
003	Anonymous ECCV 2024 Submission	003
004	Paper ID $\#1155$	004
005	7 Overview	005
006	This supplementary document provides more event representation strategies,	006
007	backdoor attack training details, experimental results, and visualization exam-	007
800	ples that accompany the paper:	800

009	_	Sec. 8 illustrates the process of training a backdoor model on the event data.	009
010	_	Sec. 9 presents more details about the event data and popular event repre-	010
011		sentation strategies.	011
012	_	Sec. 10 provides detailed experimental results of 22 classifiers on N-Caltech101	012
013		and N-Cars datasets.	013
014	_	Sec. 11 shows more visualization results of triggered samples poisoned by	014
015		three types of triggers: representation trigger, immutable trigger, and mu-	015
016		table trigger, respectively. Additionally, the point sets of the poisoned event	016
017		data generated by immutable and mutable triggers are depicted in Fig. S4.	017

018 8 Backdoor attack on event vision models

In Fig. 2 of the main paper, we show the details of training event vision model 019 019 and the design of our proposed two triggers. For training a backdoored event 020 vision model, we need first to generate some poisoned samples by Fig. 2 (c) or 021 021 (d) of the main paper. Then, we can follow the pipeline shown in Fig. S1 to train 022 022 a victim model and evaluate the attacking performance. The backdoored model 023 023 can correctly classify benign event streams, such as the motorbike and airplane 024 024 shown in the first row of Fig. S1. However, once the attacker injects the specific 025 trigger into event samples, this model will output the predetermined label. For 026 026 instance, the poisoned motorbike and ferry (in the second row of Fig. S1) are all 027 027 misclassified as accordions. This kind of potential risk could severely impact the 028 performance of autonomous driving systems. 029 029

030 9 Event data

As depicted in Sec. 3.2 of the main paper, event data consists of a series of independent and discrete events (x_k, y_k, t_k, p_k) , a kind of sparse sequence data. In contrast to conventional images, event data is recorded by the event camera 133

018



Fig. S1: The framework of the backdoor attack on event vision models. Each "puzzle piece" represents an event data stream. \mathcal{R}_{ω} denotes the module of event representation with parameters ω (E. representation), and f_{θ} represents the victim model with parameters θ .



Fig. S2: Compared with conventional cameras, an event camera obtains the data (*e.g.*, an event) asynchronously. The event data consists of all discrete events within a certain time period.

with asynchronous sensors that respond to brightness changes in a scene asyn-chronously and independently for each pixel, as shown in Fig. S2. Hence, the event data is a variable data-rate sequence of digital "events", *i.e.*, $\mathcal{E} = \{e_k\}_{k=1}^N$, where N depends on the number of brightness changes in the scene. The faster the brightness changes, the more events per second are generated. The event data reacts rapidly to visual stimuli because the events are timestamped at mi-crosecond resolution and transmitted with less than a millisecond latency.

To accommodate the input requirements of deep neural networks, the event stream needs to be transformed into the corresponding representations, also known as event representation¹. Injecting triggers into the original event data ensures that the effectiveness of the proposed *Event Trojan* is not compromised by various event representation methods. As Table 6 illustrates, our method maintains attacking effectiveness across different event representations. The rep-resentation schemes we consider in our work are listed as follows:

¹ https://github.com/LarryDong/event_representation

- 048- Event Frame (EF). EF is a simple representation strategy that considers048049the polarity (positive / none / negative) within the event data to set the pixel049050value (+1 / 0 / -1) in the images [1,2]. Furthermore, some variant versions [3]050051convert events by counting events or accumulating polarity pixel-wise into051052an image compatible with image-based vision models.052
- 053- Time Surface (TS). A TS representation [4] is also a 2D image where each053054pixel stores a single time value, e.g., the time stamp of the last event at the054055selected pixel address. Thus, the event stream is converted into an image055056where only the most recent recorded timestamps at each pixel position are056057taken into account. It can be formulated as:057

058
$$TS(x,y) = p \times \exp^{-(t_{max}-t)/\tau},$$
 (1) 058

where τ is a time constant.

065

067

075

060- Voxel Grid (VG). VG [5] is a space-time (3D) histogram of events, where060061each voxel represents a particular pixel and time interval. This representation061062preserves better the temporal information of the events by avoiding collaps-062063ing them on a 2D representation. The VG representation can be generated063064by:064

$$V(x, y, t) = \Sigma p_k \phi(x - x_k) \phi(y - y_k) \phi(t - t_k^*), \qquad 065$$

$$\phi(a) = max(0, 1 - |a|),$$
066

$$t_k^* = (B-1)(t_k - t_1)/(t_N - t_1), \qquad (2) \qquad 067$$

where B bins are used to discretize the time dimension and N denotes the length of a set of input events.

070- Tencode. Tencode [6] considers both polarities and timestamps of the event070071stream to conduct the event representation. A temporal resolution Δt is071072defined to discretize the normalized time stamps in order to produce a three-072073channel frame I by:073

074
$$I[x, y, :] = (255, \frac{255 * (t_{max} - t)}{\Delta t}, 0) \leftarrow (x, y, t, +1),$$

$$I[x, y, :] = (0, \frac{255 * (t_{max} - t)}{\Delta t}, 255) \leftarrow (x, y, t, +1),$$
(3) 075

where
$$t_{max}$$
 represents the timestamp of the latest event in the temporal resolution Δt 077

10 Detailed experimental results

079Table S1 shows the detailed quantitative results of each classifier shown in Fig.0790804 and Fig. 5 of the main paper, respectively. It's clear that the mutable trigger080081achieves better attacking performance than the immutable trigger in almost all081082cases on two public datasets. On the other hand, these victim models achieve082083better performance on the N-Cars dataset [10] than that on the N-Caltech101083

3

059

078

dataset [9], primarily because N-Cars [10] has a larger number of training sam-084 084 ples and fewer categories. On Transformer-based models, ViTs [11] perform worst 085 085 because the extracted sequence features may not adequately satisfy the down-086 086 stream tasks especially when the event data contains much background activity 087 087 noise and fewer training samples (N-Caltech101 [9]). Due to the fact that poi-088 088 soning the event representations to initiate backdoor attacks is impossible in 089 089 real-world application scenarios, we haven't conducted more explorations about 090 090 representation triggers in the following experiments. Only the classical back-091 091 door method: BadNets [14] and the latest work: FIBA [15] are chosen in our 092 092 experiments (Table 1 in the main paper). 093 003

		N-Caltech	n101 [9]			N-Cars	s [10]	
	Immutab	le Trigger	Mutable	e Trigger	Immutabl	le Trigger	Mutable	e Trigger
	CDA	ASR	CDA	ASR	CDA	ASR	CDA	ASR
ResNet-18 [16]	0.8561	0.9673	0.8621	0.9971	0.9223	0.9967	0.9272	1.0000
ResNet-34 [16]	0.8698	0.7557	0.8598	0.9741	0.9190	0.9974	0.9279	0.9934
ResNet-50 [16]	0.8572	0.8194	0.8443	0.9747	0.9281	0.9933	0.9176	0.9981
ResNet-101 [16]	0.8578	0.9954	0.8534	0.9954	0.9159	0.9985	0.9132	0.9725
ResNet-152 [16]	0.8446	0.9839	0.8218	0.9885	0.9144	0.9859	0.9374	0.9949
VGG-16 [17]	0.7064	0.1812	0.8526	0.9765	0.9211	0.9970	0.9293	1.0000
VGG-19 [17]	0.6766	0.1692	0.8521	0.9719	0.4893	1.0000	0.9280	1.0000
EfficientNet-B0 [18]	0.8589	0.8630	0.8448	0.9443	0.9453	0.9964	0.9395	0.9933
EfficientNet-B1 [18]	0.8704	0.9386	0.8586	0.9644	0.9391	0.9988	0.9402	0.9972
EfficientNet-B2 [18]	0.8607	0.8291	0.8529	0.9718	0.9347	0.9988	0.9219	0.9908
EfficientNet-B3 [18]	0.8876	0.9828	0.8747	0.9868	0.9538	0.9993	0.9434	0.9904
EfficientNet-B4 [18]	0.8704	0.9025	0.8761	0.9718	0.9391	0.9926	0.9454	0.9955
Inception-v3 [19]	0.6852	0.6451	0.8477	0.6891	0.9495	0.9972	0.9327	0.9909
ViT-S [20]	0.5086	0.1474	0.4731	0.8773	0.8453	0.9729	0.8717	1.0000
ViT-B [20]	0.4851	0.0401	0.5189	0.9943	0.8113	0.9584	0.8815	1.0000
ViT-L [20]	0.4908	0.0860	0.4736	0.9874	0.8542	0.9807	0.8809	0.9987
Swin-S ^[21]	0.7494	0.2161	0.8899	0.9994	0.7974	0.5091	0.9476	1.0000
Swin-B [21]	0.7655	0.1799	0.9203	0.9977	0.7357	0.4105	0.9457	1.0000
Swin-L [21]	0.6247	0.3115	0.9203	0.9983	0.7981	0.5186	0.9536	1.0000
DeiT-S 22	0.4868	0.0860	0.6640	1.0000	0.8532	0.9845	0.9030	0.9991
DeiT-B [22]	0.4300	0.1067	0.8165	1.0000	0.8280	0.9721	0.8865	0.9997
DeiT-L 22	0.7425	0.1508	0.8773	1.0000	0.8641	0.9829	0.8978	0.9995

Table S1: Quantitative results of the immutable trigger and mutable trigger imposed on 22 classifiers on the N-Caltech101 [9] and N-Cars [10] datasets, respectively.

094 11 Visualization of triggers

Fig. S3 and Fig. S4 show more visualization examples of the benign event, poisoned event, and corresponding triggers. In Fig. S3, it's clear that the stealthiness of the poisoned event compromised by representation trigger (R. trigger) is lower than our *Event Trojan*. Notably, in BadNets [14], the noticeable white

patch in the top-left is easily detectable by users. FIBA [15] embeds a random image into the frequency domain of the event representations, vielding better performance than BadNets. However, it is still quite noticeable when compared to benign events. Our *Event Trojan* is designed to inject triggers directly into the event data, thereby avoiding abnormal anomalies in the corresponding event representations. Fig. S4 presents some point sets of the poisoned event data compromised by our two types of triggers. The mutable trigger exhibits a more stealthy pattern than the immutable trigger.



Fig. S3: From left to right, we show benign events, poisoned events with representation trigger (R. trigger), poisoned events with immutable trigger, and poisoned events with mutable trigger, respectively. Trigger details are zoomed in on the red square for better visibility. For the representation trigger, we show two types of triggers generated by BadNets [14] (first 3 rows) and FIBA [15] (last 3 rows), respectively.



Fig. S4: Point sets of triggered samples poisoned by our immutable and mutable triggers. For better visualization, we normalize these event data in the time dimension. Details are zoomed in on the green circle \bigcirc . Blue means the polarity p = 1.0 while red denotes the p = -1.0.

107 References

108	1.	Liu M, Delbruck T. Adaptive time-slice block-matching optical flow algorithm for	108
109		dynamic vision sensors. In BMVC, 2018. 3	109
110	2.	Rebecq H, Horstschaefer T, Scaramuzza D. Real-time visual-inertial odometry for	110
111		event cameras using keyframe-based nonlinear optimization. 2017. 3	111
112	3.	Maqueda A I, Loquercio A, Gallego G, et al. Event-based vision meets deep learning	112
113		on steering prediction for self-driving cars. In Proceedings of the IEEE conference	113
114		on computer vision and pattern recognition. 2018: 5419-5427. 3	114
115	4.	Lagorce X, Orchard G, Galluppi F, et al. Hots: a hierarchy of event-based time-	115
116		surfaces for pattern recognition. IEEE transactions on pattern analysis and ma-	116
117		chine intelligence, 2016, 39(7): 1346-1359. 3	117
118	5.	Zhu A Z, Yuan L, Chaney K, et al. Unsupervised event-based learning of opti-	118
119		cal flow, depth, and egomotion. In Proceedings of the IEEE/CVF Conference on	119
120		Computer Vision and Pattern Recognition. 2019: 989-997. 3	120
121	6.	Huang Z, Sun L, Zhao C, et al. EventPoint: Self-Supervised Interest Point Detec-	121
122		tion and Description for Event-based Camera. In Proceedings of the IEEE/CVF	122
123		Winter Conference on Applications of Computer Vision. 2023: 5396-5405. 3	123
124	7.	Schaefer S, Gehrig D, Scaramuzza D. Aegnn: Asynchronous event-based graph	124
125		neural networks. In Proceedings of the IEEE/CVF conference on computer vision	125
126		and pattern recognition. 2022: 12371-12381.	126
127	8.	Gehrig D, Loquercio A, Derpanis K G, et al. End-to-end learning of representations	127
128		for asynchronous event-based data. In Proceedings of the IEEE/CVF International	128
129		Conference on Computer Vision. 2019: 5633-5643.	129
130	9.	Orchard G, Jayawant A, Cohen G K, et al. Converting static image datasets to	130
131		spiking neuromorphic datasets using saccades. Frontiers in neuroscience, 2015, 9:	131
132		437. 4	132
133	10.	Sironi A, Brambilla M, Bourdis N, et al. HATS: Histograms of averaged time	133
134		surfaces for robust event-based object classification. In <i>Proceedings of the IEEE</i>	134
135		conference on computer vision and pattern recognition. 2018: 1731-1740. 3, 4	135
136	11.	Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16x16 words:	136
137		Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.	137
138	10	4 Lin 7 Lin V Ore V et al Coris transformers Illinoushied sizion transformers	138
139	12.	Liu Z, Lin Y, Cao Y, et al. Swin transformer: Hierarchical vision transformer	139
140		using shifted windows. In Proceedings of the IEEE/CVF international conference	140
141	19	on computer vision. 2021: 10012-10022.	141
142	15.	fourron H, Cord M, Douze M, et al. Training data-encient image transformers	142
143		a distination through attention. In <i>International conjetence on machine learning</i> .	143
144	14	PMLR, 2021: 1054/-1055/.	144
145	14.	Gu I, Doran-Gavitt D, Garg S. Dadnets. Identifying vulnerabilities in the machine	145
140	15	Fong V Ma B. Zhang L et al. Fiba: Frequency injection based backdoor attack in	140
147	10.	modical image analysis[C] In Proc. CVPR 2022: 20876 20885 4 5	147
140	16	He K Zhang X Ren S et al Deep residual learning for image recognition. In	140
150	10.	Proceedings of the IEEE conference on computer vision and nattern recognition	150
151		2016: 770-778 4	151
152	17	Simonyan K. Zisserman A. Very deep convolutional networks for large-scale image	152
153	±1.	recognition, arXiv preprint arXiv:1409.1556 2014 4	153
154	18	Tan M. Le Q. Efficientnet: Rethinking model scaling for convolutional neural net-	154
155	-0.	works. In International conference on machine learning. PMLR. 2019: 6105-6114	155
156		4	156

157	19.	Szegedy C, Vanhoucke V, Ioffe S, et al. Rethinking the inception architecture for	157
158		computer vision. In Proceedings of the IEEE conference on computer vision and	158
159		pattern recognition. 2016: 2818-2826. 4	159
160	20.	Dosovitskiy A, Beyer L, Kolesnikov A, et al. An image is worth 16×16 words:	160
161		Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929, 2020.	161
162		4	162
163	21.	Liu Z, Lin Y, Cao Y, et al. Swin transformer: Hierarchical vision transformer	163
164		using shifted windows. In Proceedings of the IEEE/CVF international conference	164
165		on computer vision. 2021: 10012-10022. 4	165
166	22.	Touvron H, Cord M, Douze M, et al. Training data-efficient image transformers	166
167		& distillation through attention. In International conference on machine learning.	167

168 PMLR, 2021: 10347-10357. 4