

Event Transformer⁺. A multi-purpose solution for efficient event data processing

Alberto Sabater, Luis Montesano, and Ana C. Murillo



arXiv:2211.12222v2 [cs.CV] 3 Sep 2023

Abstract—Event cameras record sparse illumination changes with high temporal resolution and high dynamic range. Thanks to their sparse recording and low consumption, they are increasingly used in applications such as AR/VR and autonomous driving. Current top-performing methods often ignore specific event-data properties, leading to the development of generic but computationally expensive algorithms, while event-aware methods do not perform as well. We propose *Event Transformer⁺*, that improves our seminal work *EvT* with a refined patch-based event representation and a more robust backbone to achieve more accurate results, while still benefiting from event-data sparsity to increase its efficiency. Additionally, we show how our system can work with different data modalities and propose specific output heads, for event-stream classification (i.e. action recognition) and per-pixel predictions (dense depth estimation). Evaluation results show better performance to the state-of-the-art while requiring minimal computation resources, both on GPU and CPU.¹

1 INTRODUCTION

Event cameras register changes in intensity at each pixel of the sensor array providing, with minimal power consumption, asynchronous sparse information with an increased High Dynamic Range and a high temporal resolution (in the order of microseconds). Many applications such as AR/VR or autonomous driving can take advantage of this type of cameras, especially when computational power is limited or when dealing with challenging motion and lighting conditions. Although this type of sensors are relatively recent, they have already shown good results in action recognition [4], [15], tracking [6], [28], depth estimation [10], [41] or odometry [18]. The most common way to process event streams converts them into frame representations and use state of the art algorithms based on Convolutional Neural Networks [1], [3], [7], [15] and/or Recurrent Layers [7], [15]. These frame-like representations ignore the natural sparsity of event cameras. There are also methods that try to exploit this sparsity and, consequently, are more efficient,

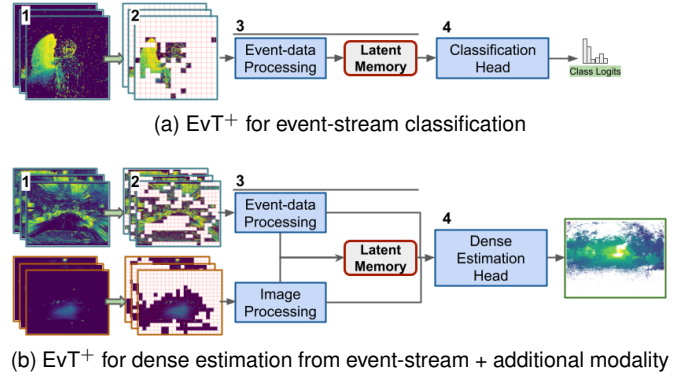


Fig. 1. **Framework overview.** Areas (*activated patches* (2)) from the input data (event frames and images (1)) with sufficient information are extracted and processed by the EvT⁺ backbone (3) to update a set of *latent memory* vectors. Different output heads (4) are used for: a) event-stream classification by processing the latent memory, and b) multi-modal dense estimation by updating and further processing the input information with the latent memory vectors.

e.g. PointNet-like Neural Networks [40], Graph Neural Networks [4], [8] or Spike Neural Networks [17], [32]. However, they usually obtain lower accuracy.

This paper proposes Event Transformer⁺, a novel solution (overview in Fig. 1) for efficient event data processing without sacrificing performance. It extends our previous work on Event Transformer [29] in three ways. First, using a finer patch-based event data representation with richer spatio-temporal information, while still benefiting from its sparsity. Second, improving the Event Transformer backbone with a more robust data processing, adapted to jointly use information from different data modalities (e.g., event data and grayscale images). In addition to this, Event Transformer⁺ can be combined with different output heads to perform either event-stream classification, i.e., action recognition, or per-pixel predictions, i.e. dense estimation.

We evaluate the new EvT⁺ in real event data benchmarks of two different tasks. First, event-stream classification (i.e. gesture recognition), where it improves prior work performance, including our previous EvT solution. Second, dense per-pixel estimation tasks (i.e. dense depth estimation), including also multi-modal inputs, events and grayscale images. In all cases, our validation demonstrates that EvT⁺ obtains better results than the state-of-the-art for the different tasks, while performing very efficiently.

- A. Sabater, L. Montesano, and A.C. Murillo are with DIIS-ISA, Universidad de Zaragoza, Spain
- L. Montesano is also with Bitbrain Technologies, Spain.
- This work was supported by DGA project T45_23R and project PID2021-125514NB-I00, funded by MCIN/AEI/10.13039/501100011033, by ERDF A way of making Europe and by the European Union NextGenerationEU/PRTR

¹. Code and trained models are available in <https://github.com/AlbertoSabater/EventTransformerPlus>

2 RELATED WORK

This section summarizes the most common approaches for event data representation as well as event-based Neural Network architectures that process them. It also includes a brief summary of available event-based datasets.

2.1 Event data representation

Event data representations encode the event information related to a time-interval or temporal-window extracted from an event-stream. These representations can be divided in two categories: **event-level representations** usually treat the event data as graphs [4], [4], [8], [40] or point-clouds [30], [37] with minimal pre-processing and keeping the event data sparsity; differently, **frame-based representations** group incoming events into dense frame-like arrays, ignoring the event data sparsity but easing a later learning process. Our work is built on the top of frame-based representations, where we find plenty of variations in the literature. The *time-surfaces* [19] build frames encoding the last generated event for each pixel. SP-LSTM [24] builds frames where each pixel contains a value related to the existence of an event in a time-window and its polarity. The *Surfaces of Active Events* [22] builds frames where each pixel contains a measurement of the time between the last observed event and the beginning of the accumulation time. *Motion-compensated* [27], [38] generate frames by aligning events according to the camera ego-motion. [11] binarizes frame representations in the temporal dimension, achieving a better time-resolution. TBR [15] aggregates binarized frame representations into single-bins frames. M-LSTM [7] uses a grid of LSTMs that processes incoming events at each pixel to create a final 2D representation. TORE [3] uses FIFOs to retain the last events for each pixel. EvT [29] build histogram-like representations for each pixel and divide the frame representation into patches, ignoring the ones with not enough event information.

The present work builds patch representations as EvT, but from frames constructed using FIFOs to retain events distributed sparsely on time. The proposed solution benefits from the sparsity of the event data, but also benefits from the robustness of using frame-based representations.

2.2 Neural Network architectures for event data

Event-stream classification has been addressed in different ways in the literature. First, we find some efficient architectures that process sparse event representations such as Spike Neural Networks [17], [32], [43], PointNet-style Networks [40] or Graph neural Networks [4], [8]. Most common, other rely on CNNs to process event-frame representations [1], [3], [7], [15]. For long event-stream processing, they are split into shorter time-windows that are frequently processed independently, and then aggregated with Recurrent Networks [15], [42], CNNs [1], [15], temporal buffers [3], [8], or voting between the intermediate results [15].

Depending on the aggregation strategy, we consider that an event-processing algorithm is able to perform **online inference** if it can evaluate the information within each time-window incrementally, as it is generated, and then perform the final visual recognition with minimal latency,

as opposed to the processing of all the captured information in a large batch. Our approach performs online inference by updating incrementally a set of latent memory vectors with simple addition operations, and processing the resulting vectors with a simple classifier.

When it comes to **event-stream dense estimation**, dense event representations and CNNs to process them are the most common scenario. LMDDE [12] uses fully convolutional networks to process the event-data and ConvLSTMs to handle their temporality. ULODE [47] trains a CNN to deblur event representations and predict optical flow, egomotion, and depth. ECN [44] uses an Evenly-Cascaded Convolutional Network to predict optical flow, egomotion and depth. DTL [39] use CNNs to translate events to images for semantic segmentation and depth estimation. RAM-Net [10] uses CNNs to encode both grayscale and event frames and ConvGRUs to update a hidden state with the temporal information, used later to perform multi-modal depth estimation. LMDDE [12] and RAM-Net [10] propose synthetic datasets to be used as pre-training.

Differently, we complement our sparse tokens (already processed by our backbone) with dummy tokens to create a dense representation that is then updated with the information from our latent memory vectors. Then, similar to RGB solutions, [26] we use skip connections between the self-attention blocks in the encoder and the dense output head to generate the final dense output.

2.3 Event dataset recordings

Large-scale public datasets recorded with event cameras in real scenarios are scarce. This lack has motivated many works to propose different approaches to translate RGB datasets to their event-based counterpart. Earlier approaches [4], [13], [20], [25], [31] display RGB data in a LCD monitor and then record the display with an event-camera. More recent approaches introduce the use of learning-based emulators [9], [14], [23] to generate event data. Unfortunately, these translated datasets cannot fully mimic the event-data nature and introduce certain artifacts, specially on their sparsity and latency. In order to have a more reliable evaluation setup, we focus our experimentation on datasets recorded with event-cameras on real scenarios. More specifically, we train and evaluate EvT⁺ for event-stream classification [1], [34], and multi-modal dense estimation [46].

3 EVENT TRANSFORMER FRAMEWORK

Different to traditional RGB cameras, event cameras log the captured visual information in a sparse and asynchronous manner. Each time the event camera detects an intensity change, it triggers an event $e = \{x, y, t, p\}$ defined by its location (x, y) within the sensor grid $(H \times W)$, its timestamp t (in the order of μs) and its polarity p (either positive or negative change). In the following, we detail the EvT⁺ contributions in terms of event-data representation and processing for classification and dense estimation tasks.

3.1 Patch-based event data representation

Similarly to previous work [1], [3], [4], [15], [29], [40], we create a frame representation for each time-window Δt that

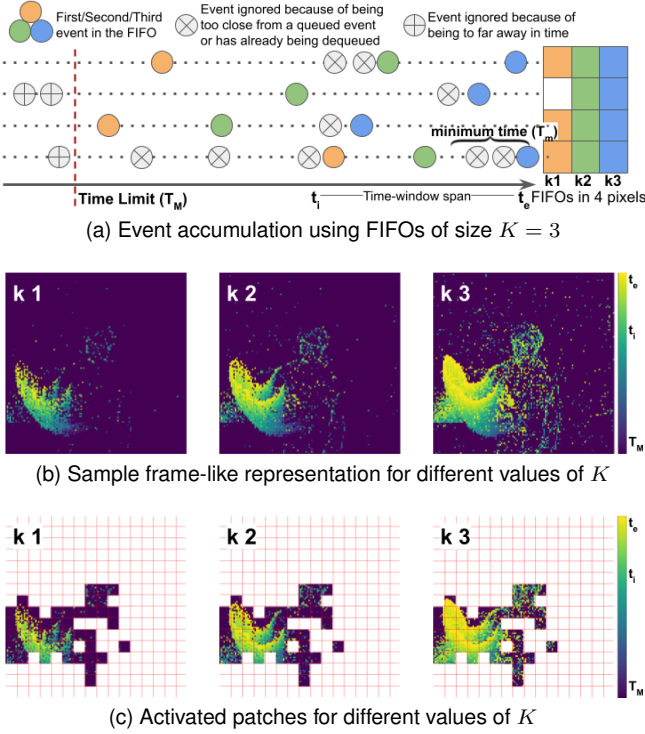


Fig. 2. Patch-based event data representation. (a) For each pixel, we retain the last K events with sufficient sparsity in time. (b) Frame representations are built with the time-stamps of the queued events. (c) Frames are split into patches, keeping only the activated patches, i.e., with enough event information generated during the time-window span.

covers a time-span from t_i to t_e . Like TORE [3], we model this event information with queues $FIFO(x, y, p, k)$ (see Fig. 2a) that retain $k \in K$ events for each pixel ($y \in H, x \in W$) within the sensor array and polarity ($p \in \{0, 1\}$). But differently, for each pixel we do not retain the last K events but the last K events that are separated by at least a minimum time of $T_m = \frac{\Delta t}{K}$. This threshold is intended to avoid the over-representation of the information provided by the events that happen consecutively in time. And additionally, when no events are registered in this time-window span for a certain pixel, we account for the ones triggered up to a maximum time T_M ($T_M \gg \Delta t$ and $T_M \ll t_i$).

Once the events are queued for a given time-window, we build an intermediate frame representation $F^{H \times W \times K \times 2}$ with their time-stamps (see Fig. 2b). We normalize the pixels to have a value in the range from 0 to T_M and then to scale their values to a 0 – 1 range (Eq. 1):

$$F = F - (t_e - T_M), F = F/T_M \quad (1)$$

Therefore, the events queued at the end of the time-window will have values close to 1 and the ones close to T_M will have a value close to 0.

Then, similar to EvT, we split the generated frame-representations into non-overlapping patches of size $P \times P$ (see Fig. 2c), and we set each patch as *activated* if it contains a minimum m percent of pixels that have information of events triggered between t_i and t_e . Note that events triggered between T_M and t_i are not involved in the patch activation decision, since they have been considered in previous time-windows, but they complement the patch information to ease later their processing. Activated patches are finally

flattened to create tokens T of size $(P^2 \cdot K \cdot 2)$, input of the transformer backbone detailed in the next subsection.

3.2 Event Transformer

Transformers [36] are a natural way to process the patch-based representation we propose. Different to other architectures, they are able to ingest lists of tokens of variable length and process them with attention mechanisms. The later, different to convolutions, focus on the whole input data (structuring it as a Query (Q), Key (K) and Value (V)) to capture both local and long-range token dependencies:

$$Attention(Q, K, V) = softmax\left(\frac{QK^T}{\sqrt{d_k}}\right)V. \quad (2)$$

The processing core of our work, Event Transformer⁺ (EvT⁺), is motivated by these ideas. As the patches T (whose length varies on the event sparsity, as described in Section 3.1) from new time-windows are being generated, they are processed by this backbone with attention mechanisms. This process ends with the update of a set of M latent vectors. These vectors act as a memory that logs the key information seen so far and its final processing allows to perform tasks such as event-stream classification or dense per-pixel estimations. This whole process (detailed in Figure 3) is divided in the following steps:

1. Patch pre-processing. In this step, the set patch tokens related to single time-windows are processed by analyzing their spatial affinity. For this purpose, each one of the T input activated patches is mapped to a vector of dimensionality D , constant along all the network. This transformation ($FF1$) consists of an initial Feed Forward layer (FF), the concatenation of 2D-aware positional embeddings, and a last FF layer. The use of positional embeddings to augment the patch information is required since Transformers, unlike CNNs, cannot implicitly know the locality of the input data. An initial set of N_1 Self-Attention blocks is then used to analyze long and short-range spatial dependencies between tokens. In the case of *multi-modal data processing*, a different patch pre-processing branch is used for each data modality (as summarized in Fig. 1b).

2. Backbone processing. In this step, the incoming spatial information from the pre-processed patch tokens is fused with the one from the latent memory vectors. This is done with a single Cross-Attention Module that processes the latent memory vectors M (as Q) with token information (as $K - V$). The resulting M' vectors are then refined with N_2 Self-Attention layers.

3. Memory update. The latent memory vectors M are updated given the new generated vectors M' with a simple sum operation and normalization:

$$M = \|M + M'\| \quad (3)$$

This augmented version of the latent vectors encodes longer spatio-temporal information and is used to perform the final downstream task, for which we have implemented the following two options.

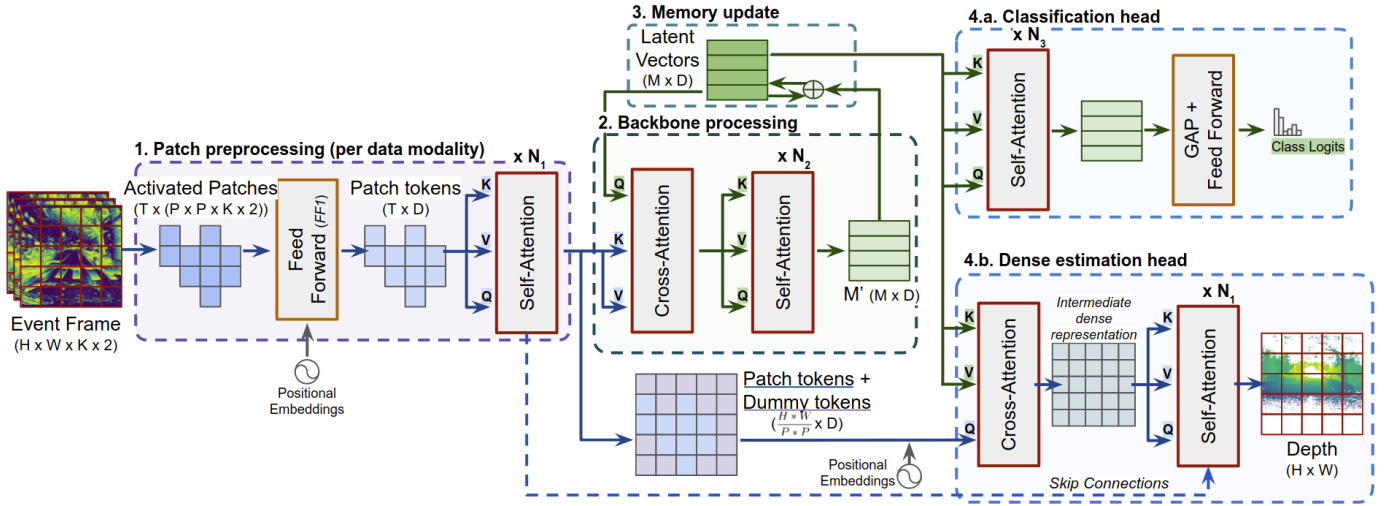


Fig. 3. Event Transformer⁺ overview. The input is a set of time-window representations (e.g., event frames or images) that are processed sequentially. Each time-window representation generates a set of patch tokens T that is processed (1., 2.) to update a set of latent memory vectors (3.), which encodes the information seen so far. For event-stream classification (4.a.), the latent vectors are directly processed with a simple classifier. For dense estimation (4.b.), we convert the input sparse representation to a dense one by adding *dummy tokens* and positional information, and we process it along with the information encoded in the latent memory vectors, generating the final dense prediction.

4.a. Classification head. Event-stream classification is performed by processing the refined latent memory vectors, that contain the key spatio-temporal information of the event-stream seen so far. This processing consists of N_3 Self-Attention modules and then, similar to EvT, processing the resulting vectors with two Feed Forward layers and Global Average Pooling (GAP).

4.b. Dense estimation head. Given the sparse set of tokens processed at step 1 and their initial location in the frame representation, we convert them back to a dense representation by adding *dummy tokens* (initialized with zeroes) that take the place of the patches filtered out due to their event-data sparsity. We then add positional information to each patch of this dense representation and update it with the information contained in the latent memory vectors (used as $K - V$) with a Cross-Attention layer. The resulting set of tokens is then refined with N_1 Self-Attention layers, with skip-connections from the N_1 Self-Attention layers of the patch pre-processing step. In this process, dense information at non-activated patches is inferred by jointly processing their positional information, their surrounding activated and non-activated (*dummy*) tokens, and the latent memory vectors, which encode the information processed in previous time-windows.

In the case of *multi-modal data processing*, the skip connections propagate the information jointly for each data modality, whose tokens are merged with a simple addition and normalization operation.

Attention modules. All the Cross and Self-Attention modules from EvT⁺ share the same architecture, similar to previous transformer related works [2], [16], [29], [36], composed of a Multi-Head Attention layer [36], normalization layers, skip connections and Feed Forward layers.

3.3 Optimization

EvT⁺ is optimized differently for different downstream tasks. In the case of **event-stream classification**, EvT⁺ is

optimized with the Negative Log-Likelihood loss:

$$\mathcal{L}_{NLL} = - \sum_i^n (Y_i \log \hat{Y}_i + (1 - Y_i) \log (1 - \hat{Y}_i)) \quad (4)$$

between the predicted labels \hat{Y}_i and the groundtruth labels Y_i , and label smoothing [45] for regularization.

In the case of **monocular dense estimation**, we train EvT⁺ in the sparse depth labels measured by a LiDAR sensor. These ground-truth depth maps, similar to other methods [10], [12], are clipped to a range $[D_m - D_M]$ captured by the sensor ($[2 - 80]$ and $[224 - 1881]$ in our case of the MVSEC and EV-IMO2 datasets) and we train EvT⁺ to predict its normalized log depth representation $\hat{Y} \in [0 - 1]$:

$$\bar{Y} = \frac{\log(Y) - \log(D_m)}{\log(D_M) - \log(D_m)}. \quad (5)$$

We optimize EvT⁺ as in [10] with a scale-invariant loss

$$\mathcal{L}_{si} = \frac{1}{n} \sum_i^n (R_i)^2 - \frac{1}{n^2} \left(\sum_i^n R_i \right)^2, \quad (6)$$

and a multi-scale invariant loss

$$\mathcal{L}_{msi} = \frac{1}{n} \sum_k^4 \sum_i^n (|\nabla_x R_i^k| + |\nabla_y R_i^k|), \quad (7)$$

where n are the valid depth ground-truth points, R_i is the log-depth difference map $\|\bar{Y}_i - \hat{Y}_i\|$ at the point i , R_i^k is the log-depth difference map at the scale $k \in [0 - 4]$. Both losses are combined as $\mathcal{L} = \mathcal{L}_{si} + \lambda \mathcal{L}_{msi}$, with $\lambda = 0.25$.

4 EXPERIMENTS

This section includes the implementation and training details of the proposed Event Transformer⁺ (EvT⁺) along with its experimental validation. We evaluate EvT⁺ in two tasks (event-stream classification and monocular dense estimation), analyze its efficiency, and main design choices.

4.1 Implementation and training details

Patch-based event representation: We set a patch size of 10×10 for event-stream classification and 12×12 for the depth estimation task, that has larger frame representations. In all cases the number of events in the FIFO, K , is set to 3, M_T is set to 256 ms, and the threshold m for the patch activation is set as in EvT (7.5%). Dataset-specific hyperparameters are discussed in the following Section 4.2.

Event Transformer: The latent vectors and the vector dimensionality D is set to 160. The latent memory is composed of 32 latent vectors. The positional encodings are initialized with 6 bands of 2D Fourier Features [33] ($\frac{H}{P} \times \frac{W}{P} \times 24$, being H and W the specific sensor height and width from each dataset). Both the latent vectors and positional encodings are learned as the rest of the parameters of the network during training. The amount of Attention layers N_1 , N_2 and N_3 are set to 1, but in the case of depth estimation, N_1 is set to 2. All the Multi-Head Attention layers use 8 heads, but in the case of depth estimation that has bigger input size, we use only 4 in the pre-processing and decoding steps to increase their efficiency.

Training details: The whole framework is optimized with the AdamW optimizer [21] in a single NVIDIA Tesla V100, with the learning rate set to $1e - 3$ and using gradient clipping. The batch-size is 128 for event-stream classification and 24 for depth estimation. Data augmentation used consists of spatial and temporal random cropping, dropout, drop token, and repetition of each sample within the training batch twice with different augmentations.

4.2 Evaluation

The proposed Event Transformer⁺ (EvT⁺) is evaluated in two scenarios of real event-camera recordings that represent different use cases. First, we evaluate EvT⁺ to classify event-streams of human actions and gestures. Second, we demonstrate that our solution is also suitable for dense inference from a sparse input, in particular, using multi-modal (grayscale image and event data) dense depth estimation. Finally, we provide a detailed analysis on how EvT⁺ takes advantage from the event-data sparsity to increase its efficiency, performing inferences with minimal latency.

4.2.1 Event-stream classification

EvT⁺ is evaluated in two benchmarks for event-stream classification. The **DVS128 Gesture Dataset** [1] is composed of 1342 event-streams capturing 10 different human gestures (plus an optional extra category for random movements) and recorded with 29 different subjects under three different illumination conditions. The **SL-Animals-DVS Dataset** [34] is composed of 1121 event-streams capturing 19 different sign language gestures, executed by 58 different subjects, under different illumination conditions. Recordings from these two datasets last between 1 and 6 seconds and contain continuous repetitions of shorter human gestures. As detailed in the supplementary material, these recordings are cropped to up to 1298 and 1792 ms and split into time-windows Δt of 24 and 48 ms for the DVS128 and SL-Animals-DVS datasets respectively.

Table 1 shows the accuracy of top-performing models in the DVS128 Dataset, with and without including the

Model	10 Classes	11 Classes	Online
RG-CNN [4]	N/A	97.2	x
3D-CNN + Voting [15]	99.58	99.62	x
CNN [1]	96.49	94.59	✓
Space-time clouds [40]	97.08	95.32	✓
CNN + LSTM [15]	97.5	97.53	✓
TORE [3]	N/A	96.2	✓
EvT	98.46	96.20	✓
EvT⁺ (Ours)	99.24	97.57	✓

TABLE 1

Classification Accuracy in DVS128 Gesture Dataset. N/A = Not Available at the source reference

Model	3 Sets	4 Sets
SLAYER [35]	78.03	60.09
STBP [35]	71.45	56.20
DECOLLE [17]	77.6	70.6
TORE [3]	N/A	85.1
EvT	87.45	88.12
EvT⁺ (Ours)	92.34	94.39

TABLE 2

Classification Accuracy in SL-Animals-DVS. N/A = Not Available at the source reference

extra additional distractor class of random movements (11 and 10 classes respectively). The column *Online* highlights the ability of each model to perform online inference, i.e., incremental processing of the event data and classification with low latency. Similarly, Table 2 shows the accuracy of top-performing methods evaluated on the SL-Animals-DVS Dataset, a more demanding benchmark with lower state-of-the-art accuracy. *3 Sets* results exclude the samples recorded indoor with artificial lighting from a neon light source, since they include noise related to the reflection of clothing and the flickering of the fluorescent lamps. *4 Sets* evaluates all the samples within the dataset.

Results from Table 1 show how our approach obtains better results than prior works, improving EvT in both data set-ups. Only [15] is more accurate than EvT⁺ but it uses offline inference and 3D-CNNs, which are computationally expensive but have a good inductive bias, useful when training with small datasets like DVS128 and with random movements (as it is the case of 11 Classes). As for the more challenging SL-Animals Dataset, EvT⁺ achieves a new state-of-the-art, outperforming prior methods by a large margin. Interestingly, our solution presents higher robustness to the different lightning conditions of the *4 Sets*, being able to take advantage of larger training set to achieve better accuracy.

4.2.2 Monocular multi-modal dense depth estimation

To evaluate EvT⁺ for dense depth estimation we use the **MVSEC Dataset** [46]. This dataset includes stereo automovilistic recordings with event-data, grayscale images and depth maps captured by a LiDAR, captured by day (2 recordings) and at night (3 recordings). Similar to previous work, we use the *outdoor_day_2* sequence for training and the remaining 4 sequences for testing. Due to the corruption of different sequences, we limit the training and validation to the data recorded from the left sensors. Depth maps are always recorded at 20 Hz, but grayscale images are recorded at 45 Hz by day and 10 Hz by night, therefore, they are not synced with the depth maps.

Model	Events	Images	outdoor day 1			outdoor night 1			outdoor night 2			outdoor night 3		
			10	20	30	10	20	30	10	20	30	10	20	30
ULODE [47]	✓	✗	2.72	3.84	4.40	3.13	4.02	4.89	2.19	3.15	3.92	2.86	4.46	5.05
LMDDE [12]	✓	✗	2.70	3.46	3.84	5.36	5.32	5.40	2.80	3.28	3.74	2.39	2.88	3.39
LMDDE [12]*	✓	✗	1.85	2.64	3.13	3.38	3.82	4.46	1.67	2.63	3.58	1.42	2.33	3.18
EvT ⁺ (Ours)	✓	✗	1.31	1.92	2.32	1.54	2.31	2.96	1.47	2.22	2.92	1.36	2.13	2.80
RAM Net [10]**	✓	✓	1.39	2.17	2.76	2.50	3.19	3.82	1.21	2.31	3.28	1.01	2.34	3.43
EvT ⁺ (Ours)	✓	✓	1.24	1.91	2.36	1.45	2.10	2.88	1.48	2.13	2.90	1.38	2.03	2.77

* Pre-training on DENSE [12] dataset

** Pre-training on EventScape [10] dataset

TABLE 3

Evaluation on the MVSEC Dataset. Average absolute depth error in meters (lower is better) at different cut-off depth distances in meters (10, 20, 30). First block shows models trained just on event data. Second block shows models trained jointly with event and image (grayscale) data.

Since the recorded sequences are very long (262 to 653 seconds), due to computational restrictions during training, we only consider the previous 512 ms (as detailed in the supplementary material) of information before the timestamp of a depth map for its inference. The event information from these sequences is split in time-windows Δt of 50 ms that are synced with the depth maps, and is complemented with a grayscale image generated at least $\Delta t/2$ ms away from the end of each time-window. Therefore, all time-steps contain event-data but might not contain grayscale image information. This issue is more frequent in the night sequences, where the grayscale image frequency is lower. When there is no grayscale information, only the event tokens update the memory and are used to the later dense depth estimation.

Table 3 shows the average depth error of different models at different cut-off depths, i.e. pixels whose groundtruth depth information is under the specified threshold (10, 20, or 30 meters). As observed, EvT⁺ is able to largely outperform previous methods in most of the set-ups, with no specific pre-training. More importantly, when including image data EvT⁺ improves its accuracy to achieve higher robustness even in the most challenging scenarios.

Additionally, we have evaluated EvT⁺ in the newly introduced **EV-IMO2 dataset** [5], which presents a different use case of indoor depth estimation and a larger sensor size. Although there is no prior work tested in this dataset, EvT⁺ is able to perform this indoor depth estimation with a mean absolute error of 176, 121, 176, and 181 mm for the different *imo*, *imo_ll*, *sfm*, and *sfm_ll* dataset splits respectively. For this experiment, we have reduced the sensor size by half by filtering out one out of two consecutive pixels from the sensor array.

4.2.3 Event sparsity and model efficiency analysis.

We now provide a deeper analysis of the theoretical computational cost of EvT⁺ and how it benefits from the data sparsity. In the case of event-stream classification, different from EvT where the computational cost ($O(|T| \times M)$) depends on the cross-attention layer, the computational cost of EvT⁺ depends on the initial self-attention pre-processing ($O(|T|^2)$), where $|T|$ stands for the amount of activated patches and M for the amount of latent memory vectors. This means that similar to EvT, the cost lowers as the input data is more sparse (less activated tokens $|T|$ are generated), but it is not bounded by M . Although this cost is theoretically higher, as observed in Table 4, in practice different implementation improvements such as bigger patch sizes P and less latent

vectors M , make EvT⁺ even more efficient (FLOPs and latency) than EvT for event-stream classification.

In the case of depth estimation, the existence of a dense output head that does not work with sparse information increases the computational cost to $O(\frac{H*W}{P*P}^2)$. In the following, we use $|P|$ to refer to this cost, defined as the number of the total generated patches. This higher cost and the use of a bigger sensor size that generates more activated patches, make EvT⁺ to demand more computational resources than for event-stream classification. Still, the rest of the Neural Network benefits from the data sparsity by suppressing non-activated patches (details in Table 5). This is also true for the ones generated from grayscale images, when pixels are black, especially in the night sequences.

However, since our network is very shallow, in all cases the final latency of EvT⁺ is minimal, being able to perform inference in a time-span significantly shorter than the time-window Δt processed, both in GPU and CPU.

As observed in Table 4, the practical computational cost (FLOPs) depends mainly on the number of tokens generated, which depends on the sensor and patch sizes, but also on the event data sparsity. It is interesting to notice that when not filtering the non-activated patches, i.e. using the whole set of generated patches $|P|$, the FLOPs that EvT⁺ demands to process a single time-window scale up between 1.35 and 2.3 times. Figure 4 also shows, for different datasets, how the demanded FLOPs depend on the event data sparsity and the patch size P . In particular, smaller patch sizes, bigger recording sensors (i.e. MVSEC dataset), and denser event information generate more activated patches, increasing the computational cost.

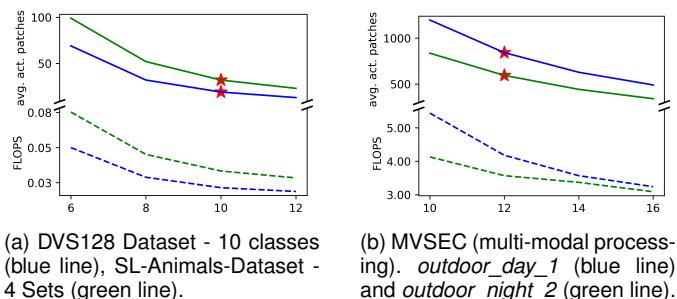


Fig. 4. Avg. number of activated patches (vertical axis) generated at each time window on different datasets with different patch sizes (horizontal axis). Stars: the selected hyperparameter value.

The supplementary material extends this efficiency analysis including efficiency statistics and comparisons with

Model	Sensor size	Dataset	P	$ T / P $	Δt ms	Latency (GPU/CPU)	FLOPs ($ T / P $)	#Params
EvT	128×128	SL-Animals	8	80 / 256	48	3 / 5 ms	0.09 / - G	0.48 M
EvT	128×128	DVS128	8	45 / 256	24	2 / 4 ms	0.08 / - G	0.48 M
EvT ⁺	128×128	SL-Animals	10	32 / 169	48	2 / 4 ms	0.04 / 0.08 G ($\times 2$)	0.66 M
EvT ⁺	128×128	DVS128	10	18 / 169	24	3 / 3 ms	0.03 / 0.07 G ($\times 2.3$)	0.66 M
EvT ⁺	346×260	MVSEC	12	318 / 638	50	10 / 25 ms	2.94 / 3.96 G ($\times 1.35$)	1.98 M
EvT ⁺ *	$346 \times 260 \times 2$	MVSEC	12	318 + 319 / 638 + 638	50	15 / 37 ms	3.68 / 4.98 G ($\times 1.35$)	2.50 M
EvT ⁺	320×240	EV-IMO2	12	237 / 540	50	10 / 42 ms	2.27 / 3.11 G ($\times 1.37$)	1.99 M

* uses multi-modal processing (events + grayscale images)

P : size of the generated patches in pixels

$|T|$: amount of activated patches

$|P|$: total amount of generated patches (no filtering)

Δt : time-window length

GPU: NVIDIA GeForce RTX 2080 Ti

CPU: Intel Core i7-9700K

TABLE 4

EvT⁺ efficiency analysis: execution time and FLOPs per Δt . Average results for all validation samples in each dataset.

EvT ⁺ Step	FLOPs
Sparse token pre-processing	0.61 G
Backbone processing and Latent Vectors update	0.06 G
Dense output head	2.27 G

TABLE 5

Average FLOPs required to process a single time-window Δt and generate a dense output for depth estimation.

standard visual backbones (ResNets), and our seminal EvT work [29] examines the benefits of our proposed representation against other non-frame event representations and non-transformer processing models.

4.3 Framework design study

The main hyperparameters that are key for achieving high efficiency while maintaining a high performance are the patch size P and the number of self-attention layers N_1 .

The **patch size** P influences the number of activated patches that are generated. Larger patch sizes generate fewer activated patches, benefiting efficiency. In our experiments, we set the patch size P as 10 for the classification task and a value of 12 for the dense estimation task, whose datasets present a larger recording sensor and more dense event information. These values allow to achieve high performance while being highly efficient.

The number of **self-attention layers** N_1 determine the quadratic computational cost of EvT⁺, so a large amount will significantly decrease the framework efficiency. In our experiments, we set the number of MHSA layers to 1 in the case of event-stream classification and a number of 2 MHSA layers in the case of dense estimation, which requires more detailed event data processing.

A thorough analysis of the key EvT⁺ hyperparameters can be found in the supplementary material.

5 CONCLUSIONS

This work presents the Event Transformer⁺ (EvT⁺) framework for event data processing, improving the seminal version of EvT. The proposed refined patch-based event representation and backbone compared to EvT are shown to provide more accurate results and increase further its efficiency for event-stream classification. This framework is demonstrated with a more complete validation, including the usage of data from different modalities and dense per

pixel estimation tasks (in particular dense depth estimation) in addition to event-stream classification tasks. Evaluation results show better or comparable accuracy to the state-of-the-art while requiring minimal computational resources, which makes EvT⁺ able to work with minimal latency both on GPU and CPU. This work shows how patch-based representations and transformers are a promising line of research for efficient event-data processing and opens opportunities for further contributions with different kinds of sparse data such as LiDAR data.

REFERENCES

- [1] Arnon Amir et al. A low power, fully event-based gesture recognition system. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017.
- [2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. *arXiv preprint arXiv:2103.15691*, 2021.
- [3] Raymond Baldwin and Ruixu et al. Liu. Time-ordered recent event (tore) volumes for event cameras. *IEEE Trans. on Pattern Analysis and Machine Intelligence*, 2022.
- [4] Yin Bi, Aaron Chadha, Alhabib Abbas, Eirina Bourtsoulatze, and Yiannis Andreopoulos. Graph-based spatio-temporal feature learning for neuromorphic vision sensing. *IEEE Trans. on Image Processing*, 2020.
- [5] Levi Burner, Anton Mitrokhin, Cornelia Fermüller, and Yiannis Aloimonos. Evimo2: An event camera dataset for motion segmentation, optical flow, structure from motion, and visual inertial odometry in indoor scenes with monocular or stereo algorithms. *arXiv preprint arXiv:2205.03467*, 2022.
- [6] Enrico Calabrese and Gemma et al. Taverni. Dhp19: Dynamic vision sensor 3d human pose dataset. In *Proc. of the IEEE/CVF Conf. on computer vision and pattern recognition workshops*, 2019.
- [7] Marco Cannici, Marco Ciccone, Andrea Romanoni, and Matteo Matteucci. A differentiable recurrent surface for asynchronous event-based data. In *European Conf. on Computer Vision*, 2020.
- [8] Yongjian Deng, Hao Chen, Huiying Chen, and Youfu Li. Ev-gcn: A voxel graph cnn for event-based object classification. *arXiv preprint arXiv:2106.00216*, 2021.
- [9] Daniel Gehrig, Mathias Gehrig, Javier Hidalgo-Carrió, and Davide Scaramuzza. Video to events: Recycling video datasets for event cameras. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2020.
- [10] Daniel Gehrig, Michelle Rüegg, Mathias Gehrig, Javier Hidalgo-Carrió, and Davide Scaramuzza. Combining events and frames using recurrent asynchronous multimodal networks for monocular depth prediction. *IEEE Robotics and Automation Letters*, 2021.
- [11] Rohan Ghosh, Anupam Gupta, Andrei Nakagawa, Alcimar Soares, and Nitish Thakor. Spatiotemporal filtering for event-based action recognition. *arXiv preprint arXiv:1903.07067*, 2019.
- [12] Javier Hidalgo-Carrió, Daniel Gehrig, and Davide Scaramuzza. Learning monocular dense depth from events. In *2020 International Conference on 3D Vision (3DV)*. IEEE, 2020.

- [13] Yuhuang Hu, Hongjie Liu, Michael Pfeiffer, and Tobi Delbruck. Dvs benchmark datasets for object tracking, action recognition, and object recognition. *Frontiers in neuroscience*, 2016.
- [14] Yuhuang Hu, Shih-Chii Liu, and Tobi Delbruck. V2e: From video frames to realistic dvs events. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2021.
- [15] Simone et al. Innocenti. Temporal binary representation for event-based action recognition. In *International Conference on Pattern Recognition (ICPR)*. IEEE, 2021.
- [16] Andrew Jaegle, Felix Gimeno, Andrew Brock, Andrew Zisserman, Oriol Vinyals, and Joao Carreira. Perceiver: General perception with iterative attention. *arXiv preprint arXiv:2103.03206*, 2021.
- [17] Jacques Kaiser, Hesham Mostafa, and Emre Neftci. Synaptic plasticity dynamics for deep continuous local learning (decolle). *Frontiers in Neuroscience*, 2020.
- [18] Simon Klenk, Jason Chui, Nikolaus Demmel, and Daniel Cremers. Tum-vie: The tum stereo visual-inertial event dataset. *arXiv preprint arXiv:2108.07329*, 2021.
- [19] Xavier et al. Lagorce. Hots: a hierarchy of event-based time-surfaces for pattern recognition. *IEEE Trans. on pattern analysis and machine intelligence*, 2016.
- [20] Hongmin Li, Hanchao Liu, Xiangyang Ji, Guoqi Li, and Luping Shi. Cifar10-dvs: an event-stream dataset for object classification. *Frontiers in neuroscience*, 2017.
- [21] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [22] Elias Mueggler, Christian Forster, Nathan Baumli, Guillermo Gallego, and Davide Scaramuzza. Lifetime estimation of events from dynamic vision sensors. In *2015 IEEE international conference on Robotics and Automation (ICRA)*. IEEE, 2015.
- [23] Jalees et al. Nehvi. Differentiable event stream simulator for non-rigid 3d tracking. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2021.
- [24] Anh Nguyen, Thanh-Toan Do, Darwin G Caldwell, and Nikos G Tsagarakis. Real-time 6dof pose relocation for event cameras with stacked spatial lstm networks. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition Workshops*, 2019.
- [25] Garrick Orchard, Ajinkya Jayawant, Gregory K Cohen, and Nitish Thakor. Converting static image datasets to spiking neuromorphic datasets using saccades. *Frontiers in neuroscience*, 2015.
- [26] René Ranftl, Alexey Bochkovskiy, and Vladlen Koltun. Vision transformers for dense prediction. In *Proc. of the IEEE/CVF Int. Conf. on Computer Vision*, 2021.
- [27] Henri Rebecq, Timo Horstschaefer, and Davide Scaramuzza. Real-time visual-inertial odometry for event cameras using keyframe-based nonlinear optimization. In *British Machine Vision Conference (BMVC)*, 2017.
- [28] Viktor et al. Rudnev. Eventhands: Real-time neural 3d hand pose estimation from an event stream. In *Proc. of the IEEE/CVF Int. Conf. on Computer Vision*, 2021.
- [29] Alberto Sabater, Luis Montesano, and Ana C. Murillo. Event transformer: a sparse-aware solution for efficient event data processing. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition Workshops*, June 2022.
- [30] Yusuke Sekikawa, Kosuke Hara, and Hideo Saito. Eventnet: Asynchronous recursive event processing. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2019.
- [31] Teresa Serrano-Gotarredona and Bernabé Linares-Barranco. A $128 \times 1281.5\%$ contrast sensitivity 0.9% fpn $3 \mu\text{s}$ latency 4 mw asynchronous frame-free dynamic vision sensor using transimpedance preamplifiers. *IEEE Journal of Solid-State Circuits*, 2013.
- [32] Sumit Bam Shrestha and Garrick Orchard. Slayer: Spike layer error reassignment in time. In *NeurIPS*, 2018.
- [33] Matthew Tancik et al. Fourier features let networks learn high frequency functions in low dimensional domains. *arXiv preprint arXiv:2006.10739*, 2020.
- [34] Ajay Vasudevan, Pablo Negri, Camila Di Ielsi, Bernabe Linares-Barranco, and Teresa Serrano-Gotarredona. Sl-animals-dvs: event-driven sign language animals dataset. *Pattern Analysis and Applications*, 2021.
- [35] Ajay Vasudevan, Pablo Negri, Bernabe Linares-Barranco, and Teresa Serrano-Gotarredona. Introduction and analysis of an event-based sign language dataset. In *Int. Conf. on Automatic Face and Gesture Recognition*. IEEE, 2020.
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, 2017.
- [37] Sai Vemprala, Sami Mian, and Ashish Kapoor. Representation learning for event-based visuomotor policies. *ArXiv*, abs/2103.00806, 2021.
- [38] Antoni Rosinol Vidal, Henri Rebecq, Timo Horstschaefer, and Davide Scaramuzza. Ultimate slam? combining events, images, and imu for robust visual slam in hdr and high-speed scenarios. *IEEE Robotics and Automation Letters*, 2018.
- [39] Lin et al. Wang. Dual transfer learning for event-based end-task prediction via pluggable event to image translation. In *Proc. of the IEEE/CVF Int. Conf. on Computer Vision*, 2021.
- [40] Qinyi et al. Wang. Space-time event clouds for gesture recognition: From rgb cameras to event cameras. In *IEEE Winter Conf. on Applications of Computer Vision (WACV)*, 2019.
- [41] Ziwei Wang, Liyuan Pan, Yonhon Ng, Zheyu Zhuang, and Robert Mahony. Stereo hybrid event-frame (shef) cameras for 3d perception. In *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems*, 2021.
- [42] Wenming Weng, Yueyi Zhang, and Zhiwei Xiong. Event-based video reconstruction using transformer. In *Proc. of the IEEE/CVF Int. Conf. on Computer Vision*, 2021.
- [43] Yujie Wu, Lei Deng, Guoqi Li, Jun Zhu, and Luping Shi. Spatio-temporal backpropagation for training high-performance spiking neural networks. *Frontiers in neuroscience*, 2018.
- [44] Chengxi Ye and Anton et al. Mitrokhin. Unsupervised learning of dense optical flow, depth and egomotion with event-based sensors. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2020.
- [45] Chang-Bin Zhang, Peng-Tao Jiang, Qibin Hou, Yunchao Wei, Qi Han, Zhen Li, and Ming-Ming Cheng. Delving deep into label smoothing. *IEEE Trans. on Image Processing*.
- [46] Alex Zihao Zhu, Dinesh Thakur, Tolga Özaslan, Bernd Pfrommer, Vijay Kumar, and Kostas Daniilidis. The multivehicle stereo event camera dataset: An event camera dataset for 3d perception. *IEEE Robotics and Automation Letters*, 2018.
- [47] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised event-based learning of optical flow, depth, and egomotion. In *Proc. of the IEEE/CVF Conf. on Computer Vision and Pattern Recognition*, 2019.



Alberto Sabater received his Ph.D. (2023) and B.S. in Computer Engineering (2016) from the University of Zaragoza, and his M.S. in Decision Making Engineering (2017) from the University Rey Juan Carlos. His research is focused on efficient computer vision algorithms for video-based scene understanding.



Luis Montesano received his B.S. and Ph.D from the University of Zaragoza, Spain. He is a Full Professor in Computer Science at the University of Zaragoza and CTO at Bitbrain Technologies. His work focuses on EEG-based neurotechnology and Artificial Intelligence applied to neuroscience.



Ana C. Murillo received her B.S. and Ph.D. from the University of Zaragoza. She is an Associate Professor in Computer Science at the University of Zaragoza, and her research is focused on computer vision and robotics.

Event Transformer⁺. A multi-purpose solution for efficient event data processing

Supplementary material

1 ABLATION STUDY

This section analyzes the influence of key components in our framework, both for the event-stream classification and dense estimation tasks.

1.1 Event-stream classification

In the following, we analyze the hyperparameters that are key to achieving better results for event-stream classification. This evaluation has been performed with both the DVS128 Gesture (10 classes) and the SL-Animals-DVS (4 sets) datasets.

Event representation hyperparameters: The key components of our event representation are the patch size (P) used to split the event frame representations, the number of events K queued for each FIFO and the maximum sequence length used to represent each event-stream sequence. As observed in Fig. 1b, a higher patch size reduces the floating point operations required to process a single time-window, but from some point it also reduces the model accuracy. Differently, a higher value of the parameter K (Fig. 1c), used to split the event information for a certain time-window, also increases the required FLOPs without a significant increment of the model accuracy. In both cases, we select, as observed in the Figures, the value that presents a good trade-off between computational cost and accuracy. Finally, we observe in Fig. 1a that, in general terms, the accuracy increases with the sequence length used to describe an event-stream. It is important to notice that this length converges to a maximum accuracy since the sequences found in the datasets are made out of the repetition of shorter action movements. Also, larger event-stream sequence lengths could not be tested because of GPU memory restrictions.

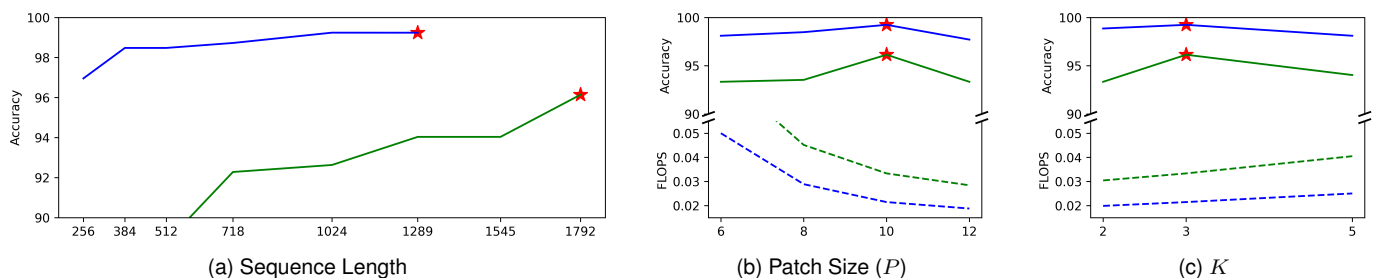


Fig. 1. EvT⁺ accuracy with different data hyperparameters. DVS128 Dataset with 10 classes (blue). SL-Animals-Dataset - 4 Sets (green). Dashed lines: average FLOPs required to process a single time-window. Stars: the selected hyperparameter value.

Backbone hyperparameters: The most relevant backbone hyperparameter for the computational cost is the amount of self-attention layers applied over the input patch tokens. As observed in Fig. 2a, its use increases the model accuracy, but the application of many of these layers increases significantly the model complexity, without a benefit over its accuracy. This is probably due to an increase of the model instability during training. As observed in Fig. 2b, the use of self-attention after the cross-attention layers does not benefit EvT⁺, both in accuracy and model complexity. Differently, using self-attention in the classification decoder does help the classification accuracy with a minimal overhead. Finally, as observed in Fig. 2d, we note that different from EvT, using more latent vectors does not help the training accuracy.

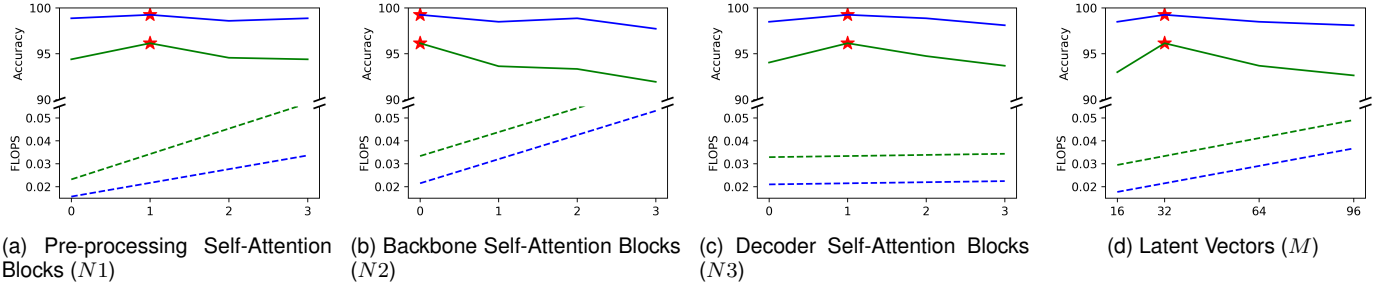


Fig. 2. EvT⁺ accuracy with different model hyperparameters for the DVS128 Dataset - 10 classes (blue line) and for the SL-Animals-Dataset - 4 Sets (green line). Dashed lines: average FLOPs operations required. Stars: the selected hyperparameter value.

1.2 Dense estimation

In the following, we analyze the hyperparameters that are key to achieving better results for monocular depth estimation by analyzing performance results over two different evaluation sequences, *outdoor_day_1* and *outdoor_night_2*, from the MVSEC dataset. For this study, we use the multi-modal version of EvT⁺.

Event representation hyperparameters: In the case of dense estimation, as observed in Fig. 3b, a higher patch size reduces the average floating point operations required to process individual single time-windows, but also leads to worse performance. For our experiments, we select a patch size P of 12, as minimal values lead to impracticable model complexity. Differently but similar to the classification task, higher values of the parameter K (Fig. 3c) used to model the time-window representations, slightly increase the required FLOPs but, at some point, it does not benefit the final performance. In both cases, we select the value that presents a good trade-off between computational cost and accuracy as observed in the Figure. Finally, we observe in Fig. 3a that, in general terms, the accuracy increases with the sequence length used to describe an event-stream. It is essential to notice that the performance starts to converge to a minimum when using 512 ms of information, depending on the evaluated sequence, which seems enough to predict the depth of a given time-stamp. Also, larger event-stream sequence lengths could not be tested because of GPU memory restrictions.

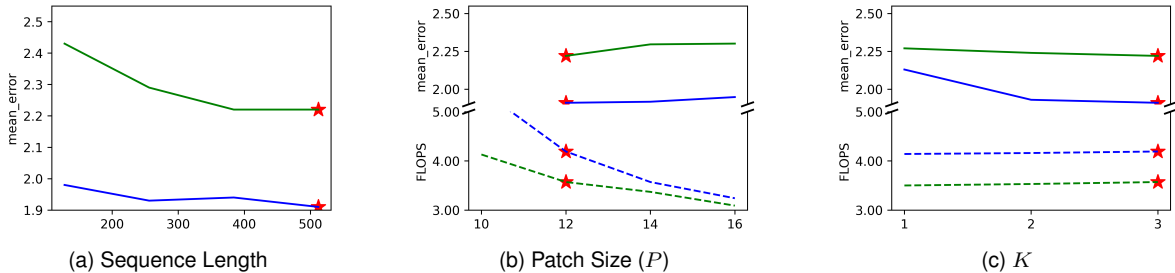


Fig. 3. EvT⁺ average depth error at 20m crop rate and FLOPs on the *outdoor_day_1* (blue) and *outdoor_night_2* (green) sequences from MVSEC. Lower is better. Dashed lines: average amount FLOPs required to process a single time-window. Stars: selected hyperparameter value.

Backbone hyperparameters: As in the case of event-stream classification, the most relevant backbone hyperparameter for the computational cost is the number of self-attention layers $N1$. This hyperparameter acquires even more importance in this task since it not only defines the number of MHSA layers for the data pre-processing, but also the number of layers in the dense output head, which has a higher computational cost, and its skip-connections. As observed in Fig. 4a, its use increases the model accuracy, but the application of many of these layers increases significantly the model complexity, without a benefit over its accuracy. As in the case of event-stream classification, this is probably due to an increase in the model instability during training. In the case of the number of MSHA layers in the backbone in Fig. 4b, we note how extra layers barely affect the final performance and efficiency for most of the evaluated dataset sequences. Finally, as observed in Fig. 4b, we note that the use of more latent memory vectors, although it does not significantly increase the model complexity, they do not benefit the final model performance.

2 EXTENDED EFFICIENCY ANALYSIS

Table 1 shows the FLOPs required by standard frame-based backbones (ResNet models) to process frames with the same sizes as the recording sensors from the evaluated datasets. As the Table shows, standard backbones require heavier computations than EvT⁺, not only because they process the whole frame even if no meaningful information is present in certain areas of the frame, but also because they are more complex and deeper neural networks. Note that these FLOPs are only calculated by processing individual frames and, different from the EvT⁺ statistics, do not include the FLOPs required

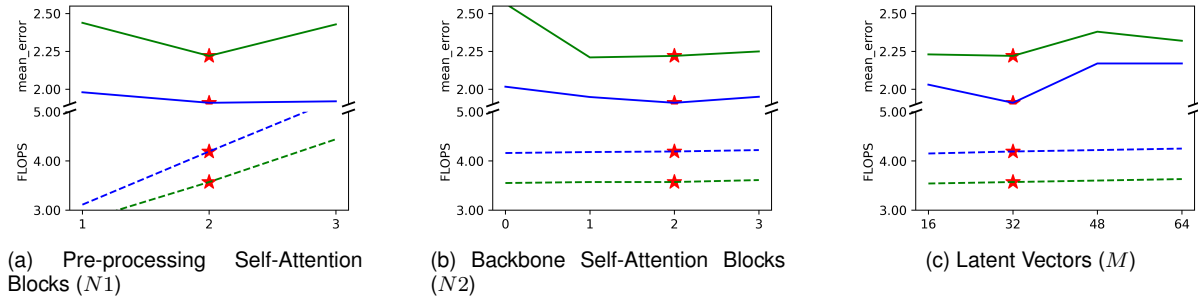


Fig. 4. EvT⁺ average depth error at 20m crop rate and FLOPs on the *outdoor_day_1* (blue) and *outdoor_night_2* (green) sequences from MVSEC. Lower is better. Dashed lines: average amount FLOPs required to process a single time-window. Stars: selected hyperparameter value.

to handle the temporal dimension or to generate the final predictions, which would significantly increase the demanded computational resources depending on the architecture used for that purpose.

Frame Size	ResNet-18	ResNet-34	ResNet-50	EvT ⁺
(128 x 128) (as in SL-Animals [2], DVS128 [1])	1.14 G	2.35 G	2.64 G	0.08 G
(346 x 260) (as in MVSEC [3])	8.55 G	17.68 G	19.86 G	3.96 G

TABLE 1

ResNet FLOPs calculated when processing different frame sizes.

EvT⁺ FLOPs also include the cost related to the decoding part, classification or dense estimation in each case.

Additionally, Table 2 shows the number of parameters of each of these standard backbones (not including specific architectures to handle the temporal dimension or decoder heads). Similarly, these common processing architectures require significantly more parameters than our proposed EvT⁺.

	Params.
ResNet-18	11.17 M
ResNet-34	21.28 M
ResNet-50	23.50 M
EvT ⁺	0.66 - 1.98 M

TABLE 2

Amount of parameters of different ResNet backbones.

EvT⁺ params. also include the ones related to the decoding part of the model.

REFERENCES

- [1] Arnon Amir et al. A low power, fully event-based gesture recognition system. In *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017.
- [2] Ajay Vasudevan, Pablo Negri, Camila Di Ielsi, Bernabe Linares-Barranco, and Teresa Serrano-Gotarredona. Sl-animals-dvs: event-driven sign language animals dataset. *Pattern Analysis and Applications*, 2021.
- [3] Alex Zihao Zhu, Dinesh Thakur, Tolga Ozaslan, Bernd Pfrommer, Vijay Kumar, and Kostas Daniilidis. The multivehicle stereo event camera dataset: An event camera dataset for 3d perception. *IEEE Robotics and Automation Letters*, 2018.