

---

# A NOVEL APPROACH FOR NEUROMORPHIC VISION DATA COMPRESSION BASED ON DEEP BELIEF NETWORK

---

**Sally Khaidem**

Department of Electrical Engineering  
Indian Institute of Technology, Madras  
Chennai, 600036, India  
ee20d041@smail.iitm.ac.in

**Mansi Sharma**

Department of Electrical Engineering  
Indian Institute of Technology, Madras  
Chennai, 600036, India  
mansisharma@ee.iitm.ac.in

**Abhipraay Nevatia**

Department of Mechanical Engineering  
Indian Institute of Technology, Madras  
Chennai, 600036, India  
me20b007@smail.iitm.ac.in

October 28, 2022

## ABSTRACT

A neuromorphic camera is an image sensor that emulates the human eyes capturing only changes in local brightness levels. They are widely known as event cameras, silicon retinas or dynamic vision sensors (DVS). DVS records asynchronous per-pixel brightness changes, resulting in a stream of events that encode the brightness change's time, location, and polarity. DVS consumes little power and can capture a wider dynamic range with no motion blur and higher temporal resolution than conventional frame-based cameras. Although this method of event capture results in a lower bit rate than traditional video capture, it is further compressible. This paper proposes a novel deep learning-based compression scheme for event data. Using a deep belief network (DBN), the high dimensional event data is reduced into a latent representation and later encoded using an entropy-based coding technique. The proposed scheme is among the first to incorporate deep learning for event compression. It achieves a high compression ratio while maintaining good reconstruction quality outperforming state-of-the-art event data coders and other lossless benchmark techniques.

**Keywords** Event computing · entropy coding · dynamic vision sensor · deep belief network

## 1 Introduction

Sight, along with the brain, is the dominant sense in humans for perceiving the world and learning new things. “Silicon Retina” [1] mimics the neural architecture of human eyes and reveals a new, powerful way of computations, sparking the emerging field of neuromorphic engineering. Bio-inspired novel sensors such as Dynamic Vision Sensors (DVS) [2] measure intensity changes asynchronously rather than capturing intensity images at a fixed rate. As a result, it generates a stream of events that encodes the time, location, and polarity of brightness changes, where the data rate depends on scene complexity and camera speed. When compared to traditional cameras, DVS have superior properties. They have a very high dynamic range (140 dB versus 60 dB), no motion blur, and measurements with latency on the order of microseconds. DVS devices, such as Dynamic and Active-pixel Vision Sensor (DAVIS) [3] and Asynchronous Time-based Image Sensor (ATIS) [4] are a viable alternative in challenging conditions for standard cameras, such as high-speed high-dynamic-range motion photography, robotic automation, and intelligent surveillance [5, 6, 7, 8].

The neuromorphic silicon technology uses Address-Event-Representation (AER) [9], a communication protocol for transferring spikes events between bio-inspired chips. A tuple  $(X, Y, p, t)$  represents each event, where  $X$  and  $Y$  denote

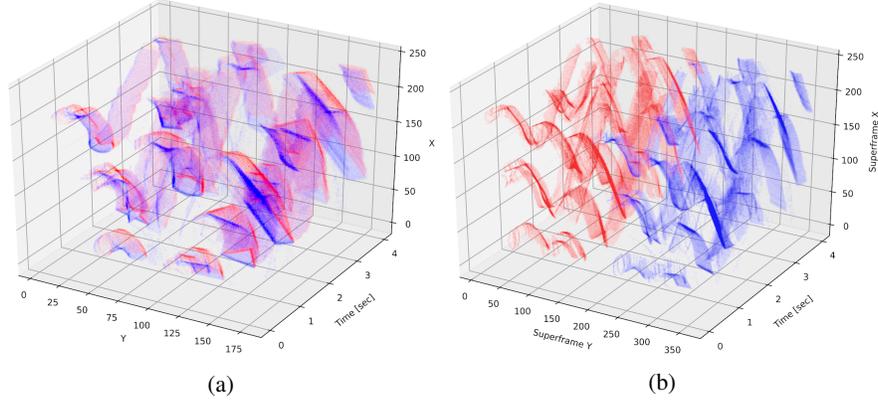


Figure 1: Visualization of *Box* event output in space-time. Red and Blue represents events with polarity ‘0’ and ‘1’ respectively. (a) without polarity separation (b) with polarity separation and creation of super-frame.

the location of the event at a particular timestamp  $t$  with polarity  $p$  indicating an increase or decrease in event brightness. Each tuple is represented by 64 bits, with the timestamp being 32 bits and the remaining three fields being 32 bits. The goal is to gather helpful information from event data and utilize it for processing.

DVS acquire information asynchronously and sparsely, with high temporal resolution and low latency. Hence, the temporal aspect, particularly latency, is critical in the event data processing. The output stream cannot use traditional vision algorithms since it is a series of asynchronous events rather than actual intensity images. Therefore, development of new algorithms that take advantage of the sensor’s high temporal resolution and asynchronous nature is necessary. There are two types of algorithms based on the processing number of events at the same time. The first approach operates on an event-by-event basis, in which the system’s state changes upon the occurrence of a single event, resulting in minimal latency. The second approach involves latency because it operates on groups or packets of events. It can still provide a system state update upon the occurrence of each event if the window moves by one event at a time. The data storage and transmission bandwidth limitation for onboard DVS processing is an open challenge and requires immediate solutions. Spike coding [10] is a dedicated lossless compression strategy that exploits event data’s time-series and asynchronous nature. It follows a cube-based coding framework where the spike sequence is divided into multiple macro-cubes and encoded accordingly. Entropy-based coding strategies like Huffman and Arithmetic can effectively encode DVS data by treating each spike event field as an input symbol. Existing lossless coding schemes such as dictionary-based [11, 12, 13] and fast-integer [14, 15] encoders can also compress the DVS data after converting the spike events into a multivariate stream of integers.

The applications of DVS range from self-driving cars [16] to robotics [17] and drones [18]. Applications such as coordinating multiple intelligent vehicles (IoV) (cars, drones, etc.) having onboard processing constraints require real-time data sharing and feedback. In comparison to traditional sensing techniques, neuromorphic sensing provides an intrinsic compression. Further compression of event data is advantageous for transmission in the Internet of Things (IoT) and the Internet of IoV. This paper presents a novel approach suitable for DVS data compression based on a deep learning algorithm, Deep Belief Network (DBN). Figure 2 depicts the complete workflow of event compression. The entire stream of events is converted into a dimensionally reduced latent representation by multiple code layer blocks using the DBN. The compact latent code blocks contain recurring information suitable for lossless symbol-based encoders. Hence, we compress the latent code using an entropy-based Huffman coding technique. The primary contributions of the proposed scheme are as follows:

- The proposed framework is among the first to incorporate deep learning techniques for event data processing. High-dimensional event data is transformed into low-dimensional latent code using a multilayer neural network called a deep belief network. We perform lossless encoding of low-dimensional latent features using entropy-based encoders to achieve a further compressed representation.
- We formulated a unique events arrangement deemed more suitable for processing by the proposed framework. The events are time-aggregated by accumulating spike events over time as super-frame sequences, as explained in Section 2.1. Super-frames result in high spatial and temporal correlation among the event data.
- We conducted extensive comparisons with lossless benchmark strategies on a diverse standard dataset with varying scene complexity and camera movement. As a result of the learning-based framework, we obtain a

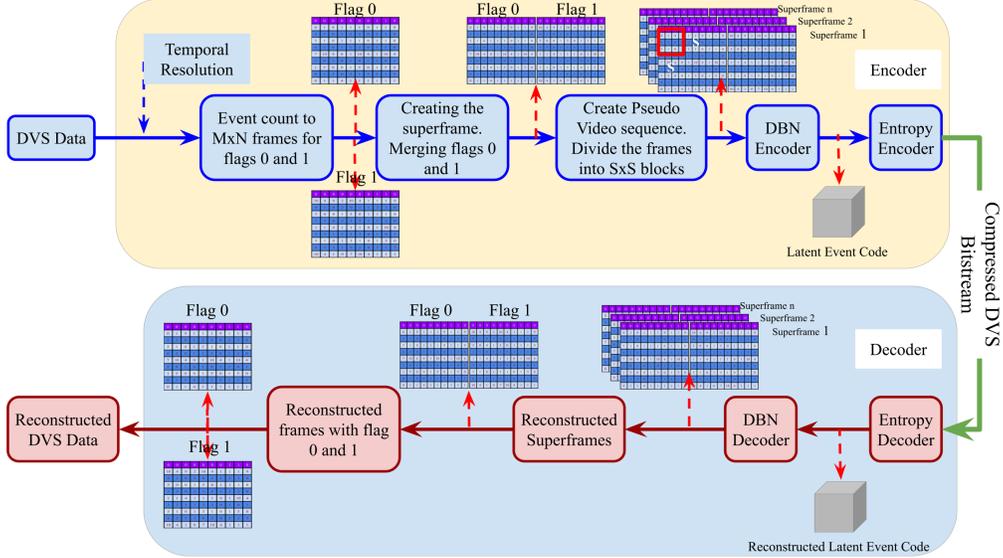


Figure 2: Complete workflow of proposed coding scheme.

concise latent code of high dimensional event data, increasing the compression gain while maintaining good reconstruction quality. Hence, the proposed framework outperforms the benchmark strategies.

## 2 Proposed Architecture

DVS generates a high-dimensional multivariate data stream indicating each event’s occurring timestamp, location and polarity. The proposed methodology’s primary goal is to obtain an efficient representation of the high-dimensional input event data. Algorithms like auto-encoder encodes the input data into a much lower dimensional representation and store latent information about the input data distribution. Autoencoders typically find poor local minima with large initial weights; with small initial weights, the gradients in the early layers are tiny, making training autoencoders with many hidden layers impossible. In DBN, multiple RBMs obtain the initial weights for the autoencoder network. DBN serves as the foundation for the coding framework and reduces the high-dimensional event data into a compressed form. The latent version of events is a string of repetitive integer values that are further compressible using an entropy-based encoder such as the Huffman encoder.

### 2.1 Event Arrangement

Even though the coding framework can encode DVS data on an event-by-event basis, we devise a unique arrangement of event data by applying time aggregation at particular time intervals. The primary benefit of event aggregation is the generation of temporal frames which reduces data size by projecting the DVS spike event stream into a series of frames recording the location histogram count, i.e., the number of event counts at each pixel. We segregate the spike sequence into two separate frames, one for increasing luminance (polarity 1) and one for decreasing luminance (polarity 0). On a particular pixel, if the previous event’s polarity was zero (or one), the next event has a high probability of having zero (or one) polarity due to smooth change in luminance [10]. The polarities of the accumulated spike events have a strong correlation so we record the event count separately for each polarity with each frame having full pixel array resolution. This increases the temporal correlation between frames of the same polarity [19]. We combine the frames with the same timestamp from each polarity into a single super-frame as shown in Figure 2. It comprises a ‘0’ polarity frame on the left and ‘1’ polarity frame on the right resulting in high inter-frame correlation. Hence, super-frames have inherent polarity information and reduce the required frame rate for processing, saving space to store the data.

### 2.2 Super-frames to Latent Representation

Exploiting spatial, temporal, and statistical correlations among the event frames associated with the DVS sequence, we compress the accumulated spike events into latent code representation using a DBN. DBN is a stack of two-layered stochastic network with visible and hidden layers where the hidden layer of one RBM is the visible layer of next RBM.

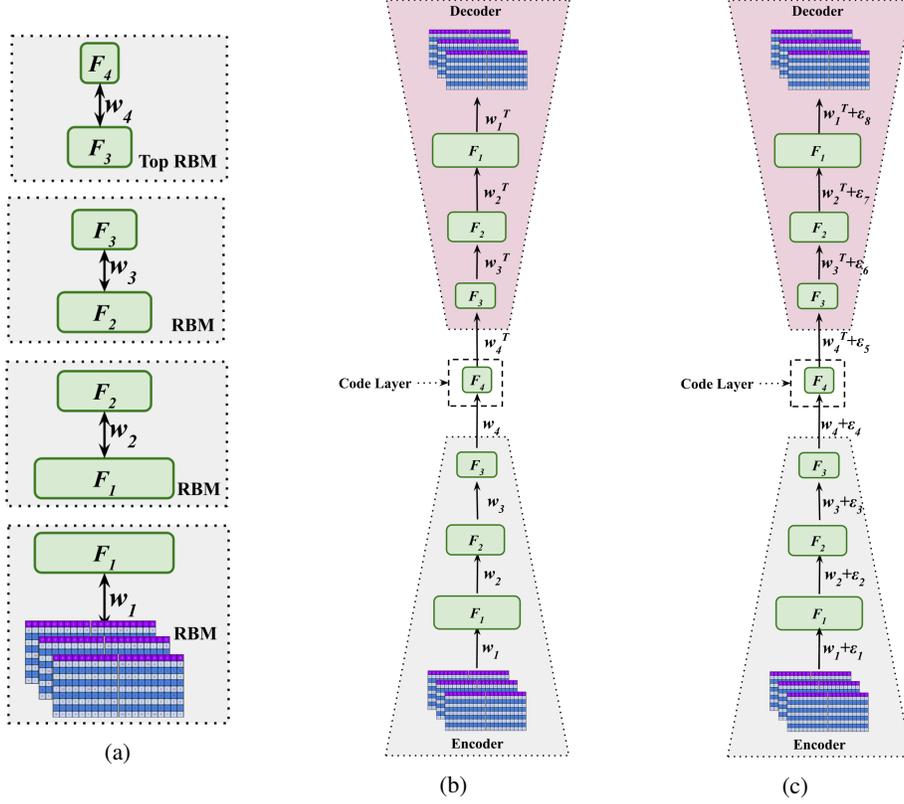


Figure 3: Workflow of deep auto-encoder with RBM pretraining. (a) Pre-training of RBM weights (b) Unrolling of stacked RBMs (c) Fine-tuning of the unrolled stacked RBMs network.

In binary RBMs, the random variables take the value  $(v, h) \in \{0, 1\}^{m+n}$  and the Gibbs distribution describes model's joint probability distribution  $p(v, h) = \frac{1}{Z} e^{-E(v, h)}$  with the energy function  $E(v, h)$ .

$$E(v, h) = - \sum_{i=1}^m \sum_{j=1}^n w_{ij} h_i v_j - \sum_{j=1}^n b_j v_j - \sum_{i=1}^m c_i h_i \quad (1)$$

where,  $w_{ij}$  is the weight associated with the nodes  $v_i$  and  $h_j$ . The bias terms for  $j^{th}$  visible and  $i^{th}$  hidden node are  $b_j$  and  $c_i$  respectively. The change in weight is given as

$$\Delta w_{ij} = \epsilon (\langle v_i h_j \rangle_{input} - \langle v_i h_j \rangle_{recon}) \quad (2)$$

where,  $\epsilon$  is the learning rate and  $\langle v_i h_j \rangle_{input}$ ,  $\langle v_i h_j \rangle_{recon}$  denote the fraction of time event  $i$  and feature detector  $j$  are 'ON' at the same time while driven by input and reconstructed output respectively.

The layer-by-layer learning algorithm is a very effective way to pre-train the weights of a deep auto-encoder. Each feature layer captures high-order correlations between the events in the layer below. This gradually reveals the low-dimensional, nonlinear structures in a wide range of data sets. The workflow of DBN starts with pretraining process of feature detector layers with initial weights update (2), as shown in Figure 3a. After pretraining is completed, the model is unfolded (Figure 3b) to form a deep auto-encoder network initialised with the pre-trained weights. The pre-trained weights of deep auto-encoder are fine-tuned by replacing stochastic activities with deterministic, real-valued probabilities using the back-propagation method (Figure 3c). For fine-tuning, we utilised the conjugate gradient method. The number of features extractors  $F_1, F_2, F_3$  and  $F_4$  in each layers are varied according to the requirement of input data such that  $F_2 > F_3 > F_4$  and  $F_2 > F_1$ . The code layer with  $F_4$  features is the latent space representation of the event super-frames. The latent code of accumulated spike events super-frame is a string of integer values. We used Huffman coder, an entropy-based encoder, to compress the latent code further. Huffman coding aims to reduce the expected code length by assigning shorter codes to frequently used characters and longer ones to rarely used characters. Adding Huffman coding to the DBN framework improves compression performance as discussed in section 3.2, even though Huffman coding alone produces low compression gains. The encoder and decoder network share a symmetrical

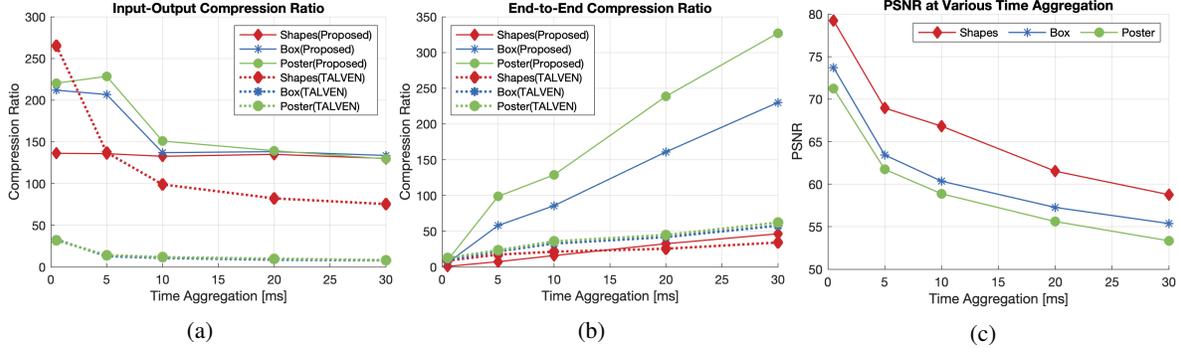


Figure 4: Comparison of compression performance of Proposed Scheme with respect to TALVEN. (a) Input-Output Compression Ratio (b) End-to-end Compression Ratio (c) PSNR at various Time Aggregation.

structure and number of features in each layer. The implementation details and experimental results are analysed in the later sections.

### 3 Experiments

In this work, we first evaluate the proposed compression scheme’s performance based on end-to-end and input-output compression ratios. The end-to-end compression ratio is between uncompressed events and compressed event bitstream size. In contrast, the input-output compression ratio is the input and output frame size ratio. Subsequently, we provide an extensive performance comparison with benchmark strategies regarding compression gains and peak-signal-to-noise ratio (PSNR).

#### 3.1 Implementation Details

To test our compression scheme’s performance, we used standard event camera datasets [20]. We choose three different datasets *Boxes*, *Poster* and *Shapes* to evaluate the competence of our proposed scheme at varying degree of scene complexity and camera movement speed. The evaluation was done on a workstation with Intel Xeon 2.30GHz CPU, Nvidia K80 12GB GPU and 12GB RAM specifications. All the experiments were conducted using the Pytorch framework.

The auto-encoder framework comprises of an encoder with layers size  $(30 \times 30)$ -1000-500-250-20 and a symmetric decoder. The units in the code layer are linear, while other units are logistics. The model’s accuracy increases with lesser units in the code layer, but it has the drawback of longer training time and an increase in low-dimensional code blocks. Hence, we settled with a 20-dimensional code layer maintaining a balance between training time, accuracy and the number of code blocks. The 20-dimensional code blocks are compressed using an entropy-based Huffman coding technique generating a compressed DVS bitstream. The network was trained on 96,000 event blocks derived from the first 10 seconds of the event sequence with 10ms time-aggregation for 20 epochs and 10 batch size. The average training time for each sample was around 30 minutes, with the learning rate set at 0.1. The network was evaluated with test samples obtained with various degree of time-aggregation (0.5ms, 5ms, 10ms, 20ms and 30ms).

#### 3.2 Experimental Results

As shown in Table 1, we compared our performance with existing benchmark strategies regarding end-to-end compression ratio on standard datasets with diverse scene complexity and camera movement speed. We outperform all the benchmark coders at almost all time-aggregation values. We also compared our performance with a state-of-the-art event data coder proposed by Khan et al. [21] in terms of end-to-end compression and input-output compression ratio. In both cases of compression metrics, our proposed compression scheme outperforms the TALVEN method [21] at all time-aggregation values as depicted in Figure 4b and Figure 4a respectively. In addition, higher event rate sequences produce better compression performance than low event rate sequences. The *Poster* sequence, for example, has a very high event rate (4.01 Mega-events/s), while the *Shapes* sequence has the lowest event rate (0.245 Mega-events/s). The *Shapes* sequence has low scene complexity and camera speed. Naturally, it should have high compression gains while, in reality, the *Poster* sequence yields high compression gain. In sequences with high event rate, the delta-coded timestamp need lesser bits to be encoded when compared to low event rate sequences. Additionally, we computed

Table 1: Comparison of End-to-End Compression Ratio of Proposed Scheme with Benchmark Coders.

	<b>Time Aggregation</b>	<b>Shapes</b>	<b>Boxes</b>	<b>Poster</b>
<b>Proposed</b>	<b>0.5ms</b>	0.837	5.885	9.516
	<b>5ms</b>	7.406	57.983	98.697
	<b>10ms</b>	15.935	85.715	128.834
	<b>20ms</b>	32.554	160.974	238.774
	<b>30ms</b>	46.381	230.125	327.093
<b>Benchmark</b>	<b>SPIKE Coding</b>	3.78	4.95	4.88
	<b>Huffman</b>	1.79	1.96	1.79
	<b>Zstd</b>	2.67	4.13	4.02
	<b>LZMA</b>	3.04	4.92	4.77
	<b>LZ4</b>	2.19	3.03	2.97
	<b>ZLib</b>	2.8	4.21	4.12
	<b>Brotli</b>	2.78	4.38	4.26
	<b>Snappy</b>	2.21	2.98	2.92

the reconstruction quality of our proposed model with the help of the peak-signal-to-noise ratio (PSNR) metric. The proposed model can reconstruct the dynamic vision sensor event data from the compressed bitstream with a high PSNR value, as shown in Figure 4c.

## 4 Conclusion

Dynamic vision sensor captures the change in per pixel intensity value producing an asynchronous stream of event data. We proposed a novel framework for event data arrangement and its compression based on DBN. It is among the first few to utilise a deep learning framework for event data processing. Time aggregation of spike events is advantageous because it enhances efficiency by bringing a specific group of events into context with one another. The super-frame arrangement yields high spatial and temporal correlations among events which the DBN utilises to obtain latent feature code. The proposed coding framework outperforms the benchmark compression schemes as supported by the extensive performance comparison study presented in the paper.

## References

- [1] Misha Mahowald. The silicon retina. In *An Analog VLSI System for Stereoscopic Vision*, pages 4–65. Springer, 1994.
- [2] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbruck. A  $128 \times 128$  120 db  $15\mu\text{s}$  latency asynchronous temporal contrast vision sensor. *IEEE journal of solid-state circuits*, 43(2):566–576, 2008.
- [3] Christian Brandli, Raphael Berner, Minhao Yang, Shih-Chii Liu, and Tobi Delbruck. A  $240 \times 180$  130 db  $3\mu\text{s}$  latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits*, 49(10):2333–2341, 2014.
- [4] Christoph Posch, Daniel Matolin, Rainer Wohlgenannt, Michael Hofstätter, Peter Schön, Martin Litzenberger, Daniel Bauer, and Heinrich Garn. Live demonstration: Atis camera with full-custom ae processor. In *Proceedings of IEEE ISCAS*, pages 1392–1392. IEEE, 2010.
- [5] Hanme Kim, Ankur Handa, Ryad Benosman, Sio-Hoi Ieng, and Andrew J Davison. Simultaneous mosaicing and tracking with an event camera. *J. Solid State Circ*, 43:566–576, 2008.
- [6] 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*, 3(2):994–1001, 2018.
- [7] Henri Rebecq, René Ranftl, Vladlen Koltun, and Davide Scaramuzza. Events-to-video: Bringing modern computer vision to event cameras. In *Proceedings of the IEEE/CVF CVPR*, pages 3857–3866, 2019.
- [8] Alex Zihao Zhu, Liangzhe Yuan, Kenneth Chaney, and Kostas Daniilidis. Unsupervised event-based learning of optical flow, depth, and egomotion. In *Proceedings of the IEEE/CVF CVPR*, pages 989–997, 2019.
- [9] Vincent Chan, Shih-Chii Liu, and Andr van Schaik. Aer ear: A matched silicon cochlea pair with address event representation interface. *IEEE TCAS-I*, 54(1):48–59, 2007.

- [10] Zhichao Bi, Siwei Dong, Yonghong Tian, and Tiejun Huang. Spike coding for dynamic vision sensors. In *2018 Data Compression Conference*, pages 117–126. IEEE, 2018.
- [11] Yann Collet and Murray Kucherawy. Zstandard compression and the application/zstd media type. Technical report, 2018.
- [12] Peter Deutsch and Jean-Loup Gailly. Zlib compressed data format specification version 3.3. Technical report, 1996.
- [13] Jyrki Alakuijala and Zoltan Szabadka. Brotli compressed data format. Technical report, 2016.
- [14] Davis Blalock, Samuel Madden, and John Guttag. Sprintz: Time series compression for the internet of things. *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, 2(3):1–23, 2018.
- [15] Daniel Lemire and Leonid Boytsov. Decoding billions of integers per second through vectorization. *Software: Practice and Experience*, 45(1):1–29, 2015.
- [16] 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 *Proceedings of the IEEE/CVF CVPR*, pages 5419–5427, 2018.
- [17] Amin Rigi, Fariborz Baghaei Naeini, Dimitrios Makris, and Yahya Zweiri. A novel event-based incipient slip detection using dynamic active-pixel vision sensor (davis). *Sensors*, 18(2):333, 2018.
- [18] Elias Mueggler, Basil Huber, and Davide Scaramuzza. Event-based, 6-dof pose tracking for high-speed maneuvers. In *2014 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2761–2768. IEEE, 2014.
- [19] Xavier Lagorce, Garrick Orchard, Francesco Galluppi, Bertram E Shi, and Ryad B Benosman. Hots: a hierarchy of event-based time-surfaces for pattern recognition. *IEEE TPAMI*, 39(7):1346–1359, 2016.
- [20] Elias Mueggler, Henri Rebecq, Guillermo Gallego, Tobi Delbruck, and Davide Scaramuzza. The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and slam. *The International Journal of Robotics Research*, 36(2):142–149, 2017.
- [21] Nabeel Khan, Khurram Iqbal, and Maria G Martini. Time-aggregation-based lossless video encoding for neuromorphic vision sensor data. *IEEE Internet of Things Journal*, 8(1):596–609, 2020.