

Asynchronous Corner Tracking Algorithm based on Lifetime of Events for DAVIS Cameras*

Sherif A.S. Mohamed¹, Jawad N. Yasin¹, Mohammad-Hashem Haghbayan¹, Antonio Miele³, Jukka Heikkonen¹, Hannu Tenhunen², and Juha Plosila¹

¹ University of Turku, 20500 Turku, Finland

² KTH Royal Institute of Technology, 11428 Stockholm, Sweden

³ Politecnico di Milano, 20133 Milano, Italy

Abstract. Event cameras, i.e., the Dynamic and Active-pixel Vision Sensor (DAVIS) ones, capture the intensity changes in the scene and generates a stream of events in an asynchronous fashion. The output rate of such cameras can reach up to 10 million events per second in high dynamic environments. DAVIS cameras use novel vision sensors that mimic human eyes. Their attractive attributes, such as high output rate, High Dynamic Range (HDR), and high pixel bandwidth, make them an ideal solution for applications that require high-frequency tracking. Moreover, applications that operate in challenging lighting scenarios can exploit from the high HDR of event cameras, i.e., 140 dB compared to 60 dB of traditional cameras. In this paper, a novel asynchronous corner tracking method is proposed that uses both events and intensity images captured by a DAVIS camera. The Harris algorithm is used to extract features, i.e., frame-corners from keyframes, i.e., intensity images. Afterward, a matching algorithm is used to extract event-corners from the stream of events. Events are solely used to perform asynchronous tracking until the next keyframe is captured. Neighboring events, within a window size of 5x5 pixels around the event-corner, are used to calculate the velocity and direction of extracted event-corners by fitting the 2D planar using a randomized Hough transform algorithm. Experimental evaluation showed that our approach is able to update the location of the extracted corners up to 100 times during the blind time of traditional cameras, i.e., between two consecutive intensity images.

Keywords: Corner · asynchronous tracking · Hough transform · event cameras · lifetime.

1 Introduction

Simultaneous Localization And Mapping (SLAM) approaches use onboard sensors, such as Lidar, camera, and radar to observe the environment and estimate the position and the orientation of the robot [9] [18] [19]. For example, visual SLAM methods use single or multiple cameras to obtain the pose by extracting

* Supported by the Academy of Finland under the project (314048)

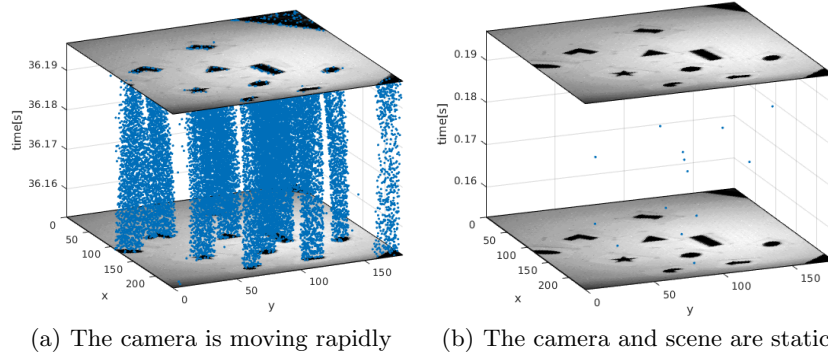


Fig. 1. Events (blue) triggered between two consecutive frame-images in two different scenarios.

and tracking a sufficient number of features from consecutive images. Over the past years, there are many methods have been presented to detect frame-corners (e.g., *Harris* [6]) and edges (e.g., *Canny* [4]).

Most of SLAM approaches in the literature obtain the pose of robots by progressing a sequence of intensity images captured by CMOS sensors, i.e., frame-based cameras. However, such cameras mostly suffer from some major hardware limitations. They generate grayscale or RGB images at a fixed rate, typically 60Hz. Moreover, they produce unnecessary information when the camera is not moving and the scene is static, which increases the computational cost dramatically. On the other hand, when the camera and objects in the scene are highly dynamic they are affected by the motion blur phenomena, i.e., they generate insufficient information. In addition, a huge amount of information might be missed during the blind time, i.e., the time between two successive images, which might result in inaccurate pose estimates. In conclusion, the limitations of conventional cameras would degrade the tracking performance, resulting in inaccurate localization.

Newly emerged event cameras, such as DVS [7] and ATIS [13], overcome these problems through transmitting a stream of events by capturing the intensity changes instead of the “absolute” intensity. Events are triggered when the intensity of any pixel changes by a certain threshold. Thus, the number of triggered events depends on the texture of the environment and the motion of the camera and objects in the scene. rightness increases “1” or decreases “0”. Figure 1 illustrates the output of the event camera in two different.

When the camera and the scene are static or the camera is moving slowly in a texture-less scene, event cameras only generate a few numbers of events and therefore only a small amount of resources would be enough to process those events (see Figure 1(b)). Contrarily, when the scene is highly dynamic or the camera is moving rapidly, a large number of events (in order of millions) will be generated (see Figure 1(a)). These powerful attributes make event cameras

a great solution for applications that operates in high dynamic scenes and need to get their pose updated frequently, i.e., drones.

Given to the discussed necessity of robust and high-frequency corner detection algorithms, in this paper we present a novel algorithm that exploits the attractive attributes of event cameras and operates in an asynchronous fashion. Our algorithm uses both events and grayscale images captured by the DAVIS [3] cameras to extract and track corners. The proposed algorithm is composed of three main phases: observation, detection, and asynchronous tracking. A sufficient number of salient frame-corners are detected using Harris [6] from grayscale images. Afterward, a matching filter is applied to obtain event-corners by selecting the first event that occurs on the same pixel location of extracted frame-corners. A window size of 5x5 pixels around each event-corner is used to compute the lifetime of each event-corner. We use a randomized Hough transform to fit a local plane to the points in the 5x5 pixels matrix [17] and compute the lifetime and the direction of each event-corner. Lifetime indicates the time a event-corner will take to move from its current pixel to one of the eight neighboring pixels. Our algorithm is able to update the extracted corners up to 2.4×10^3 times per second compared to 24 using only images.

The rest of the paper is organized as follows. Section 2 reviews related works. In Section 3, the system architecture of the proposed work is presented. An experimental evaluation of the proposed approach is discussed in Section 4. Last section draws conclusion and presents future work.

2 Related Work

To unlock events' camera potential, novel methods are required to process the output, i.e., a stream of events. Methods in the literature can be categorized as indirect and direct.

Indirect methods pre-process the captured events to generate artificially synthesized frames, i.e., event-frames and then apply one of the state-of-the-art frame-based methods to extract and track frame-corners. In [2], the authors generate event-frames by accumulating events during a fixed time interval. In [14], the authors proposed another approach to generate event-frames by accumulating a fixed number of events, (e.g., 2000 events) to form a frame. In [8], the authors presented a dynamic slicing method to generate event-frames based on the velocity of the event camera and the number of edges in the environment, i.e., *entropy*. The main disadvantage of using such techniques is that they omit the innate asynchronous nature of the event cameras.

Direct approaches, on the other hand, process the asynchronous stream of events directly. In [16], Vasco et al. proposed an algorithm to detect asynchronous corners using an adaptation of the original Harris score [6] on a spatio-temporal window, i.e., surface Active Events (SAE). The algorithm demands a lot of computational resources to compute the gradients of each incoming event. In [10], the authors present an algorithm inspired by FAST [15] to detect event-corners, they call it eFAST. Corners are extracted on SAE by comparing the

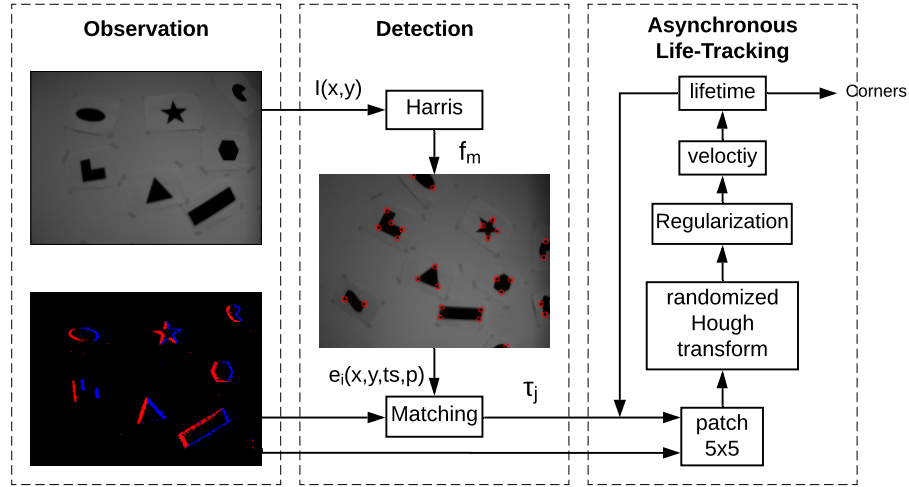


Fig. 2. Overview of the overall asynchronous corner detection and tracking algorithm.

timestamp of the latest events, i.e., pixels on two circles. The radius of the inner and the outer circle is 3 and 4 pixels respectively. In [1], the authors used a similar technique (Arc) with an event filter to detect corner corners 4.5x faster than the eFAST. There are two main disadvantages of using such techniques mainly related to computational cost and accuracy. Since event cameras can generate millions of events per second and processing each incoming event to extract corners is computationally expensive. Thus, they are not feasible to run in real-time performance and especially on resource contained systems. On the other hand, some approaches rely on simple operations to perform real-time corner extraction. However, these methods produce inaccurate corners (50 ~ 60%) compared to frame-based methods, such as Harris.

In conclusion, indirect methods omit one of the main characteristics of event-cameras, i.e., asynchronous. Direct methods, on the other hand, mistakenly detect more false event-corners and subsequently require a lot of computational resources to process all the incoming events.

3 Proposed approach

The proposed approach, illustrated in Figure 2, is composed of three main phases: observation, detection, and asynchronous tracking. Details of the various phases are discussed in the following subsections.

3.1 Observation

The sensor of DAVIS [3] cameras has an Active Pixel Sensor (APS) [5] and a Dynamic Vision Sensor (DVS) for each pixel, making them able to generate

synchronized images and events at the same time. Generated images contain the absolute brightness of the scene and have a resolution of 240x180 pixels. Grayscale images are normally captured at a constant rate (equal to 24Hz). In our approach, we use images as a *keyframe* to extract strong and robust corners. The other output of the DAVIS camera is events. Events represent the change of brightness of each pixel independently. Events are triggered asynchronously at a high rate (up to 1MHz), which makes them ideal for tracking in a highly dynamic environment. Events contain simple information: the position of the pixel, the timestamp, and the polarity of the brightness changes [1, -1].

3.2 Detection

In the detection phase, we extract 2-dimensional interest points, i.e., corners from the scene, using Harris [6] algorithm. As illustrated in Figure 3(a), a point in an image is considered as a corner when its local pixels have a significant change in all eight directions. Edges have significant changes in six direction but not along the edge direction Figure 3(b). On the other hand, flat regions have no change in all directions Figure 3(c). Thus, corners are important for location and object detection algorithm based in vision sensors, since they are invariant to motion changes, such as rotational and translation and illumination changes. Harris is considered one of the most robust corner detection algorithm in the literature.

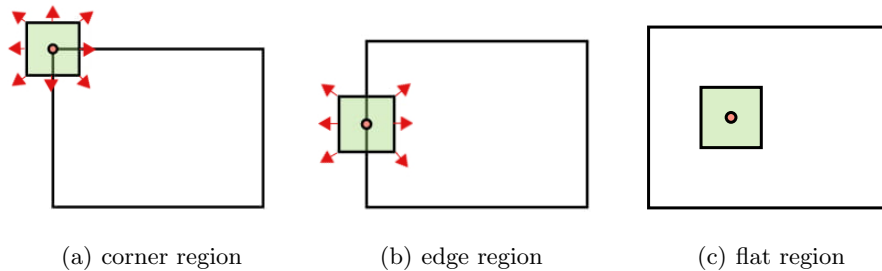


Fig. 3. A general presentation of three different regions in images: corners, edges, and flat regions.

Harris algorithm uses a score technique to extract strong corners from the capture images. A point is considered as a corner if its score is larger than a certain threshold. A window size of 3x3 pixels is used to calculate the S score of each pixel in the image, as follows

$$S = \lambda_1 \lambda_2 - k(\lambda_1 + \lambda_2)^2 \quad (1)$$

where λ_1 and λ_2 are the eigenvalues of M , i.e.,

$$M = \sum w(x, y) \begin{bmatrix} I_x^2 & I_x I_y \\ I_x I_y & I_y^2 \end{bmatrix} \quad (2)$$

where $w(x, y)$ denotes the local window of each pixel. The horizontal and vertical gradients is denoted by I_x and I_y respectively. Point P is considered as a corner if its score S is large, in other words its λ_1 and λ_2 are large and $\lambda_1 \sim \lambda_2$. After we detect corners from grayscale images, we use a matching algorithm to extract event-corners.

Algorithm 1 Matching Unit

Input: frame-corners, events

Output: event-corners

- 1: Initialization
 - 2: **if** e.position = corner.position **then**
 - 3: Select L_SAE
 - 4: Determine binary L_SAE
 - 5: Compute I_x and I_y
 - 6: Compute Score S
 - 7: **if** $S >$ threshold **then**
 - 8: Process event
 - 9: **else**
 - 10: Discard event
 - 11: **end if**
 - 12: **else**
 - 13: Discard event
 - 14: **end if**
-

The matching algorithm, as summarized in Algorithm 1, filters the incoming event and only processes the first corner-candidate that occurs on the same location of the extracted corners. Processed events, i.e., event-corners, are then fed to the next phase to track the event-corners in asynchronous fashion using neighbouring events. The inputs of the algorithm are the triggered events and the extracted frame-corners from Harris. If the incoming event is on the same location of a corner, a local surface active events (L_SAE) is extracted around the event. The L_SAE has a size of 7x7 pixels and contains the timestamp on neighboring events. As shown in Figure 4, only the most recent N neighbours ($N = 12$) are included and labeled as “1” in binary L_SAE to compute the vertical and horizontal gradients. Those gradient are then used to compute the score of the incoming event. If the score is larger than the threshold, the incoming event is considered as event-corner and it’s only updated when a new image, i.e., keyframe is captured.

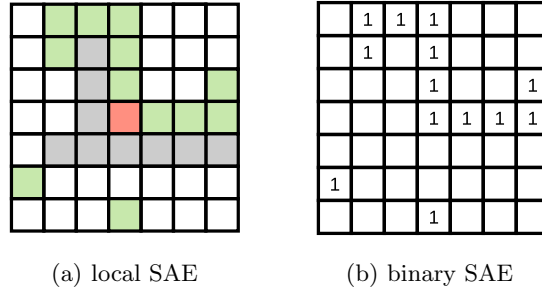


Fig. 4. A presentation of local and binary SAE. Incoming event, i.e., corner-candidate, is in the center of the SAE in red. The most recent N neighbours are in green and are labeled as “1” in the binary SAE. Old neighbouring events are in gray and they are not included in the binary SAE.

3.3 Asynchronous Life-Tracking

In this section, we explain our asynchronous tracking algorithm which is based on the lifetime of events. The concept of lifetime refers to the time an event will take to move from the current pixel to one of the eight neighboring pixels. Since event cameras have a high temporal rate up to 8 million events per second, lifetime is feasible for event-based tracking. The summary of the asynchronous tracking algorithm is illustrated in Algorithm 2, in which we first extract a local L_SAE of size 5×5 pixels around each event-corner. The randomized Hough transform algorithm [17] is used to fit a local plane using the neighboring events in the L_SAE. The vector $(a, b, \text{ and } c)$ which is orthogonal onto the obtained plane is calculated. Then we compute the velocity of the event-corner in x - and y -directions v_x and v_y respectively, as follows:

$$v_x = c * (-a) / (a^2 + b^2) \quad (3)$$

$$v_y = c * (-b) / (a^2 + b^2) \quad (4)$$

The plane fitting algorithm is inspired by [11], in which we robustly compute candidate planes using three points, i.e., the event-corner and two additional neighboring events from L_SAE. The plane fitting algorithm is illustrated in Algorithm 3, where two neighboring are chosen randomly plus the event-corner to fit a local plane (Line 1). The algorithm calculated the Hough space parameters $H(\theta, \phi, \rho)$ and store in different spaces in cells C (Lines 3-4). For each vote, a cell counter C increments by one $C + 1$ (Lines 5-6) and the most voted cell is considered as the fitting plane. The Hough space is computed using two randomly picked points p_2, p_3 from the point set in L_SAE and the event-corner p_1 , i.e.,

$$\rho = v \cdot p_1 = ((p_3 - p_1) \times (p_1 - p_2)) \cdot p_1 \quad (5)$$

$$p_x \cdot \cos \theta \cdot \sin \phi + p_y \cdot \sin \phi \cdot \sin \theta + p_z \cdot \cos \phi = \rho \quad (6)$$

Algorithm 2 Asynchronous Life-Tracking

Input: $e_i = (x_i, y_i, t_i, pol_i)$ and $\tau = \{e_1, e_2, \dots, e_j\}$ **Output:** new location of events in τ

- 1: Extract SAE of size 5x5 for each event in τ
 - 2: Fit a local plane for each SAE
 - 3: Calculate a_j, b_j , and c_j
 - 4: Calculate velocity $v_j(v_x, v_y)$
 - 5: Calculate lifetime
-

where the distance to the centroid is denoted by ρ , ϕ is the angle between the vector v and the XY-plane in the z-direction and θ denotes the angle of v on the XY-plane.

Algorithm 3 Randomized Hough transform

Input: $\mathbf{P} = \{p_1, p_2, \dots, p_n\}$ **Output:** $H(\theta, \phi, \rho)$

- 1: **while** points in $\mathbf{P} > 2$ **do**
 - 2: Select event-corner p_1 and two additional points p_2, p_3
 - 3: Calculate Hough space
 - 4: Store planes in cells
 - 5: Points vote for cells
 - 6: For each vote cell's Counter $C = C + 1$
 - 7: **if** $C = \text{threshold}$ **then**
 - 8: Parameterize the detected plane
 - 9: Delete p_2 and p_3 from \mathbf{P}
 - 10: $C = 0$
 - 11: **else**
 - 12: Continue
 - 13: **end if**
 - 14: **end while**
-

4 Experimental Evaluation

We evaluated the proposed approach on the publicly available datasets [12]. Subsets used in the experiment consists of slow and fast camera motions and low- and high-textured scenes to ensure comprehensive evaluation scheme. In general, DAVIS camera can generate up to 10 million events per second in textured and highly dynamic scenes. On other hand, in slow camera motion and texture-less scene, event camera triggers few events. As mentioned previously, DAVIS camera produce a sequence of intensity images along side a stream of asynchronous events. The subsets used in the experiment are recorded by a DAVIS240 camera which generates images with resolution of 240x180. We implemented our algorithm in C++, OpenCV and Eigen libraries. All experiments were carried out on an embedded system with an ARM-based processor, hosted on Jetson TX2 board and running at 2 GHz.

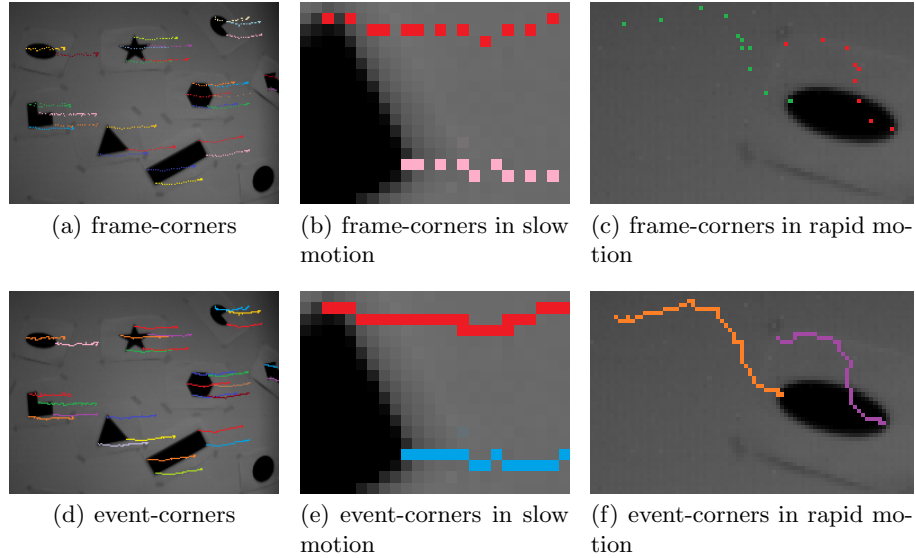


Fig. 5. The qualitative result of proposed method compared with Harris method on intensity images.

The experimental results of the proposed asynchronous corner tracking algorithm are reported in Figure 5. For this evaluation we used *shapes_6dof* dataset. The dataset includes various 2D geometric shapes mounted on a wall and recorded freely moving DAVIS-240 cameras with different speeds. We compared the proposed algorithm with normal Harris corner detection and at different speeds of the camera. In slow camera motions, frame-corner detection method, i.e., Harris performances considerably well as only break of 1 pixel appears in the detection, as shown in Figure 5(b). However, when the speed of the camera, or some objects in the environment, is high, feature extraction and tracking based on intensity images cannot provide sufficient information about the environment. This can be shown in Figure 5(c) where the increase in the speed of the camera imposes a break in detecting up to 4 pixels between two consecutive detections. In such an environment, an event-based camera and the proposed feature tracking technique can provide a seamless feature tracking for the algorithm, see Figure 5(f).

Table 1 summarizes the execution time of different units in the proposed method. The intensity images have a resolution of 240x180, which leads to fast execution time for Harris algorithm, i.e., the average execution time is 1.8 milliseconds. Execution time for Harris does not change significantly since the algorithm processes all pixels in the image and all captured images have the same size. On the other hand, event-based units, such as matching unit and lifetime calculation depends on the number of generated events that is based on the speed of the camera and the amount of information in the scene [8]. The execution time of both units is very low and especially the matching unit which has an average

execution time of 3.6 microseconds. One of the main advantages of the proposed solution w.r.t. the state-of-the-art event-based approaches is the fact that it is able to provide high performance in all the scenarios; in particular, in case the camera moves very fast, up to 10 Millions events can be generated. This tremendous load causes an execution time of the classical matching unit equal to 36 seconds and of the lifetime calculation equal to 4900 seconds, thus not affordable in any real-time scenario. At the opposite, the proposed filtering solutions has an average execution time of 5.62 seconds. We extract a maximum of 50 corners per image using Harris algorithm from images. The DAVIS-240 camera captures a total number of 24 images per second. For each corner an average number of 10 events are checked to detect event-corners using the matching unit, which totals a number of 12×10^3 events are checked per second resulting a maximum execution time of 43.2 milliseconds.

Table 1. The average execution time of units of the proposed method

Unit	Time (ms)
Harris per image	1.8
matching per event	3.6×10^{-3}
lifetime per event	0.49

5 Conclusion

Performing high-frequency corner detection and tracking based on traditional cameras is not feasible due to the limitations of such cameras, such as motions blur and blind time. In this paper, we exploit event-based cameras since they can cope well with rapid camera movements and highly dynamic scenes. We proposed an asynchronous tracking algorithm based on the lifetime of extracted event-corners. We first extract strong frame-corners from intensity images using Harris. Afterward, we match those corners with first generated events with the same pixel location. The tracking algorithm fits a local plane to the neighboring events of each extracted event-corner to compute the lifetime and velocity of each event-corner. A randomize Hough transform is used to fit a plane. Finally, we evaluated our proposed method by performing experiments on datasets with different speeds and compared the results with the extracted corner using only intensity images. In addition, the execution time of each unit of the proposed method is reported. The results show that our method has an average execution time of 5.62 seconds.

References

1. Alzugaray, I., Chli, M.: Asynchronous corner detection and tracking for event cameras in real time. *IEEE Robotics and Automation Letters* **3**(4), 3177–3184 (Oct 2018). <https://doi.org/10.1109/LRA.2018.2849882>

2. Alzugaray, I., Chli, M.: Asynchronous corner detection and tracking for event cameras in real time. *IEEE Robotics and Automation Letters* **3**(4), 3177–3184 (2018)
3. Brandli, C., Berner, R., Yang, M., Liu, S.C., Delbruck, T.: A 240x180 130db 3 μ s latency global shutter spatiotemporal vision sensor. *IEEE Journal of Solid-State Circuits* **49**(10), 2333–2341 (2014)
4. Canny, J.: A computational approach to edge detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **PAMI-8**(6), 679–698 (Nov 1986)
5. Fossum, E.R.: Cmos image sensors: electronic camera-on-a-chip. *IEEE Transactions on Electron Devices* **44**(10), 1689–1698 (1997)
6. Harris, C., Stephens, M.: A combined corner and edge detector. In: *Proceedings of the 4th Alvey Vision Conference*. pp. 147–151 (1988)
7. Lichtsteiner, P., Posch, C., Delbrück, T.: A 128 \times 128 120 db 15 μ s latency asynchronous temporal contrast vision sensor. *J. Solid-State Circuits* **43**(2), 566–576 (2008)
8. Mohamed, S.A.S., Haghbayan, M., Heikkonen, J., Tenhunen, H., Plosila, J.: Towards dynamic monocular visual odometry based on an event camera and imu sensor. In: *Intelligent Transport Systems, From Research and Development to the Market Uptake (INTSYS2019)*. Springer International Publishing (2020)
9. Mohamed, S.A.S., Haghbayan, M., Westerlund, T., Heikkonen, J., Tenhunen, H., Plosila, J.: A survey on odometry for autonomous navigation systems. *IEEE Access* pp. 97466–97486 (2019)
10. Mueggler, E., Bartolozzi, C., Scaramuzza, D.: Fast event-based corner detection. In: *British Machine Vision Conference (BMVC)* (2017)
11. Mueggler, E., Forster, C., Baumli, N., Gallego, G., Scaramuzza, D.: Lifetime estimation of events from dynamic vision sensors. *Proceedings - IEEE International Conference on Robotics and Automation* **2015** (05 2015)
12. Mueggler, E., Rebecq, H., Gallego, G., Delbrück, T., Scaramuzza, D.: The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and SLAM. *I. J. Robotics Res.* **36**(2), 142–149 (2017)
13. Posch, C., Serrano-Gotarredona, T., Linares-Barranco, B., Delbruck, T.: Retinomorphc event-based vision sensors: Bioinspired cameras with spiking output. *Proceedings of the IEEE* **102**(10), 1470–1484 (2014)
14. Rebecq, H., Horstschaefer, T., Gallego, G., Scaramuzza, D.: EVO: A geometric approach to event-based 6-dof parallel tracking and mapping in real time. *IEEE Robotics and Automation Letters* **2**(2), 593–600 (2017)
15. Rosten, E., Drummond, T.: Machine learning for high-speed corner detection. In: *Proceedings of the 9th European Conference on Computer Vision - Volume Part I*. pp. 430–443. *ECCV'06*, Springer-Verlag, Berlin, Heidelberg (2006)
16. Vasco, V., Glover, A., Bartolozzi, C.: Fast event-based harris corner detection exploiting the advantages of event-driven cameras. *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)* pp. 4144–4149 (2016)
17. Xu, L., Oja, E.: Randomized hough transform (rht): Basic mechanisms, algorithms, and computational complexities. *CVGIP: Image Understanding* **57**(2), 131 – 154 (1993)
18. Yasin, J.N., Mohamed, S.A.S., Haghbayan, M., Heikkonen, J., Tenhunen, H., Plosila, J.: Unmanned aerial vehicles (uavs): Collision avoidance systems and approaches. *IEEE Access* **8**, 105139–105155 (2020)
19. Yasin, J.N., Mohamed, S.A.S., Haghbayan, M.H., Heikkonen, J., Tenhunen, H., Yasin, M.M., Plosila, J.: Energy-efficient formation morphing for collision avoidance in a swarm of drones. *IEEE Access* **8**, 170681–170695 (2020)