

000  
001  
002054  
055  
056

# Real-time 6-DoF Pose Estimation by an Event-based Camera using Active LED Markers

003  
004  
005  
006  
007057  
058  
059  
060  
061

Anonymous WACV Applications Track submission

008  
009  
010  
011062  
063

Paper ID 282

012  
013  
014064  
065

## Abstract

015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044066  
067  
068  
069  
070  
071  
072  
073  
074  
075  
076  
077  
078  
079  
080  
081  
082  
083  
084  
085  
086  
087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

Real-time applications for autonomous operations depend largely on fast and robust vision-based localization systems. Since image processing tasks require processing large amounts of data, the computational resources often limit the performance of other processes. To overcome this limitation, traditional marker-based localization systems are widely used since they are easy to integrate and achieve reliable accuracy. However, classical marker-based localization systems significantly depend on standard cameras with low frame rates, which often lack accuracy due to motion blur. In contrast, event-based cameras provide high temporal resolution and a high dynamic range, which can be utilized for fast localization tasks, even under challenging visual conditions. This paper proposes a simple but effective event-based pose estimation system using active LED markers (ALM) for fast and accurate pose estimation. The proposed algorithm is able to operate in real time with a latency below 0.5 ms while maintaining output rates of 3 kHz. Experimental results in static and dynamic scenarios are presented to demonstrate the performance of the proposed approach in terms of computational speed and absolute accuracy, using the OptiTrack system as the basis for measurement. Moreover, we demonstrate the feasibility of the proposed approach by deploying the hardware, i.e., the event-based camera and ALM, and the software in a real quadcopter application. Our project page is available at: [almpose.github.io](https://almpose.github.io)

045

086

## 1. Introduction

046  
047  
048  
049  
050  
051  
052  
053087  
088  
089  
090  
091  
092  
093  
094  
095  
096  
097  
098  
099  
100  
101  
102  
103  
104  
105  
106  
107

Fast and reliable spatial localization is essential in a wide range of robotic applications. For example, in collaborative scenarios, the ability to accurately and rapidly estimate the pose of the end effector is a key component for achieving a safe, reliable and robust execution of corresponding tasks. Vision-based methods [3, 14, 37] are the most common ap-

\*These authors contributed equally to this work.

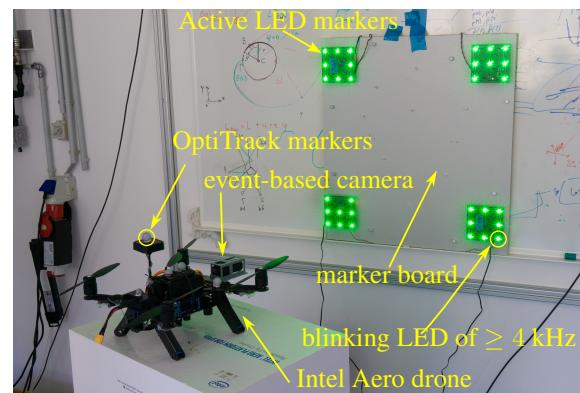


Figure 1. Overview of the experimental setup. Active LED markers (ALM) are attached to a marker board. An event-based camera mounted on a drone is used to estimate the pose of the marker board.

proaches for obtaining the relative localization of objects within the line of sight. These methods achieve significantly better accuracy compared to other non-contact localization methods, e.g. radio-based localization approaches [11, 33]. Vision-based approaches are, however, computationally expensive and typically require more than one sensor, e.g. infrared-based systems [26]. To reduce the computational overhead, classical markers [15, 36] serving as easy-to-detect anchors are often integrated into vision-based systems. Since conventional RGB-D cameras are often used in these systems, the latency of detection cannot be reduced beyond the limit determined by the frame rate of the utilized cameras.

Event-based vision is an emerging field that has attracted much attention in recent years [10, 24]. An event-based camera consists of an array of independent pixels measuring changes in luminosity  $L = \log(I)$ , based on the photocurrent  $I$  [18]. A change in the continuous luminosity signal

$$\Delta L(\mathbf{u}_k, t_k) = L(\mathbf{u}_k, t_k) - L(\mathbf{u}_k, t_k - \Delta t_k) > p_k C \quad (1)$$

triggers an event  $\mathbf{e}_k = (\mathbf{u}_k, t_k, p_k)$  at pixel location  $\mathbf{u}_k =$

108      $(u_k, v_k)$  due to a temporal contrast threshold  $\pm C$ ,  $p_k \in$   
 109      $\{+1, -1\}$  being its polarity and  $\Delta t_k$  the time since the last  
 110     event (at  $\mathbf{u}_k$ ) occurred at  $t_k$  [10]. The state-of-the-art event-  
 111     based sensors can produce up to 1.2 Giga events per second  
 112     (Geps) [10] with a microsecond range timestamp accuracy.

113     Compared to frame-based cameras that deliver periodically  
 114     dense (*i.e.* full-frame) information, the event stream is  
 115     sparse and contains information that relates only to changes  
 116     in the scene. Additionally, event-based cameras have a high  
 117     temporal resolution and a large dynamic range. These fea-  
 118     tures make them ideal for applications requiring fast and  
 119     accurate detection. Starting from the early years of event-  
 120     based vision development [9], advantages given by those  
 121     sensors for robotic applications are noticeable. Recent ad-  
 122     vancements in event-based sensor development [10] have  
 123     enabled them to compete with the precision of other lo-  
 124     calization methods [7] due to increased resolution and re-  
 125     duced noise. Avoiding accumulated event representations  
 126     (*i.e.* frames), markers can be tracked online utilizing the  
 127     event-based camera's high temporal resolution of up to 1  $\mu\text{s}$ .

128     An active LED marker (ALM) is a fixed geometric arrange-  
 129     ment of individual LEDs, each unambiguously identifiable  
 130     by its unique blinking frequency. By identifying the indi-  
 131     vidual LEDs in the event stream and knowing the geo-  
 132     metric arrangement of the LEDs, the pose of the ALM can be  
 133     retrieved.

134  
 135     In this paper, we propose a fast and simple method  
 136     employing an event-based camera together with ALM for  
 137     simultaneous detection and tracking of the 6 degrees-of-  
 138     freedom (DoF) pose of a rigid object in the 3D space. An  
 139     overview of our proposed approach is depicted in Fig. 1  
 140     with four ALMs attached to a marker board. The event-  
 141     based camera is mounted on the drone, which is utilized  
 142     to estimate the pose of the marker board with respect to  
 143     the camera's base frame. To estimate the pose, the blink-  
 144     ings of the LEDs are logged with an event-based camera  
 145     to identify the corresponding frequencies of each LED in  
 146     the ALM. These blinking frequencies are utilized to iden-  
 147     tify each individual LED and match it with the known geo-  
 148     metry of the ALM. With this mapping of the individual  
 149     points on the camera's sensor plane and the known geo-  
 150     metry of the ALM, the pose of the ALM can be computed  
 151     by utilizing a Perspective-n-Point (PnP) algorithm. In the  
 152     presented approach, by tuning the biases, *i.e.*, parameters  
 153     for tuning the analog front-end of the event-based camera,  
 154     and using a priori knowledge about timing, the complexity  
 155     of the ALM tracking can be simplified. This aids in reduc-  
 156     ing the tracking latency. During the tracking, the initial de-  
 157     tection is continuously refined, resulting in subpixel resolu-  
 158     tion. Such an approach can still precisely estimate the pose  
 159     even under fast rotational and linear motion. The proposed  
 160     approach was tested and verified extensively using an ex-  
 161     ternal infrared-based positioning system. Our contributions

162  
 163  
 164  
 165  
 166  
 167  
 168  
 169  
 170  
 171  
 172  
 173  
 174  
 175  
 176  
 177  
 178  
 179  
 180  
 181  
 182  
 183  
 184  
 185  
 186  
 187  
 188  
 189  
 190  
 191  
 192  
 193  
 194  
 195  
 196  
 197  
 198  
 199  
 200  
 201  
 202  
 203  
 204  
 205  
 206  
 207  
 208  
 209  
 210  
 211  
 212  
 213  
 214  
 215

are listed in the following.

- We propose a fast event-based pose estimation system using ALM achieving a latency below 0.5 ms while maintaining an output rate of 3 kHz.
- We analyze the proposed system in static and dynamic scenarios for several in-depth aspects, *e.g.*, absolute accuracy, static noise, and latency. Translational errors of  $34.5 \text{ mm} \pm 16 \text{ mm}$  and  $0.74^\circ \pm 0.15^\circ$  orientation errors at distances of 2.1 m to 4.8 m between the camera and the marker were achieved. Together with the fast computing speed, this proves that the proposed algorithm is promising for real-time applications.
- We integrate the proposed system into a quadcopter application for the 6-DoF pose estimation task. For indoor experiments, the proposed system outperforms the ORB-SLAM algorithm. Furthermore, in outdoor experiments, the proposed system can simultaneously detect and track the ALM in very aggressive flights at velocities of up to  $10 \text{ m s}^{-1}$  and up to 10 m away from marker.

The paper is organized as follows: Section 2 presents the related work in the field of pose estimation with event-based cameras and active markers. Section 3 describes the proposed method for marker detection and tracking. In Section 4, we present the experimental setup and results. Finally, we conclude the paper in Section 5 with a summary of our contributions and suggestions for future work.

## 2. Related Work

Visual localization systems show improved accuracy compared to systems based on other physical principles [27], [7]. Fiducial marker-based systems [13] constitute the most common choice for robotic applications. Due to the limited range and the dependence on the lighting conditions, some studies proposed LED-based solutions based on standard RGB cameras [35], infrared [34], or ultraviolet [31] spectrum. However, the latency cannot be reduced beyond the camera's frame rate.

One of the first works in the direction of localization based on event-based sensors was the 2D localization method [32]. The known shape (contours) was tracked, and the relative localization was determined by event-based vision. The high temporal resolution of the event-based sensors was used in [21] to localize an Unmanned Aerial Vehicle (UAV) during high-speed maneuvers. The pose information was retrieved using a black square as a known shape. In [16], a visual odometry method was proposed based on the feature tracking algorithm. In this direction, multiple methods were developed [16], [20], which show a significant improvement compared to the RGB-based approach for high-speed applications.

216 The utilization of ALMs was proposed first in [22],  
 217 where the authors tracked the 2D position of the LED and  
 218 used it as a feedback signal for the robot homing and a pan  
 219 tilt system. Later, the first method for pose estimation using  
 220 ALMs was presented in [4]. Therein, ALMs were used to  
 221 detect and estimate the position of a flying quadrocopter.  
 222 LEDs were recognized and detected using event polarity  
 223 changes in the event stream. In [4], the authors used an  
 224 accumulated event representation to decode the frequency  
 225 and estimate the pose. In [6], a Gaussian mixture proba-  
 226 bility hypothesis density filter was proposed to localize the  
 227 camera with respect to the active marker. Therein, online  
 228 tracking was presented to increase the robustness and reli-  
 229 ability of the pose estimation. The achieved results indicate  
 230 a localization error lower than 3 cm in scenarios where the  
 231 camera was within 1 m relative to the active marker.

232 Most recent works using ALMs propose the additional  
 233 fusion of inertial measurements [28]. The error in the pre-  
 234 dicted relative position is in the subcentimeter range. How-  
 235 ever, utilizing only the vision-based approach increases the  
 236 error by the order of one magnitude. Compared to previous  
 237 methods, the marker size is significantly larger. The LEDs  
 238 are placed 1 m apart. Current work in active marker-based  
 239 solutions also focuses on the visual communication aspect  
 240 of modulated light.

241 Different from other approaches in the literature, our ap-  
 242 proach simplifies the complexity by tuning the biases and  
 243 using a priori knowledge about timing. This helps to reduce  
 244 the tracking latency. To the best of the authors' knowledge,  
 245 this work achieves the lowest latency compared to other  
 246 methods in the literature.

### 247 3. Active Marker Tracking and Pose Es- 248 timation

249 Using ALMs, periodic and dense signals can be gen-  
 250 erated as a projection of the LED on the camera's sensor  
 251 plane.

252 To reduce the computational complexity and the required  
 253 bandwidth, the biases of the sensor are tuned to generate a  
 254 single event per pixel on every LED blink while suppressing  
 255 all other background events to increase the signal-to-noise  
 256 ratio, as presented in Figure 2. While [4] uses events of  
 257 both polarities and relatively low frequencies (1-2kHz), the  
 258 amount of noise can be reduced by using higher frequencies  
 259 and disabling one polarity.

260 The ALM's structure is an arrangement of high-  
 261 frequency blinking LEDs, where a unique frequency of the  
 262 blinking pattern can individually recognize each LED (*e.g.*  
 263 different blinking frequencies). The arrangement of the  
 264 LEDs has to be fixed and determined in the 3D coordi-  
 265 nate space. However, it can also be arranged on a plane,  
 266 as utilized in this work. Based on the 2D projection of the  
 267 LEDs, knowing their 3D arrangement and camera intrin-  
 268

269 sics, the relative pose of the marker with respect to the sen-  
 270 sor can be reconstructed using a Point-n-Perspective (PnP)  
 271 algorithm [19]. In this work, the *IPPE* PnP algorithm is  
 272 used [8].

273 The proposed approach is divided into four parts, as il-  
 274 lustrated in Fig. 2. To reduce the noise in the signal, the bias  
 275 settings are tuned to produce a single event per pixel on ev-  
 276 ery blink of a LED. Next, events are accumulated over the  
 277 time  $\frac{2}{f_{\min}}$ , which is two times the period of the LED's mini-  
 278 mal frequency  $f_{\min}$  (typ.  $f_{\min} \approx 2 \text{ kHz}$ ). For those accumu-  
 279 lated event clusters, frequencies are recognized. These fre-  
 280 quencies are used to identify newly appearing ALMs. For  
 281 each of the ALM's LEDs, trackers are spawned that keep  
 282 track of the LEDs' center points based on single events. The  
 283 tracking of the LEDs is independent of the detection loop.  
 284 The pose of the ALM is estimated utilizing the trackers of  
 285 the ALM. The accuracy of the pose estimation can be ob-  
 286 tained using the reprojection error. When an ALM leaves  
 287 the field of view or the reprojection error exceeds the de-  
 288 fined maximal value, the corresponding trackers are deleted.  
 289 If the ALM enters the field of view again, the detection al-  
 290 gorithm respawns it. Such an approach reduces the latency  
 291 and increases the accuracy of the solution.

#### 292 3.1. Detection

293 For the detection of an ALM, the geometrical arrange-  
 294 ment and the blinking frequency of the ALM LEDs, have to  
 295 be provided in advance.

296 The range of possible frequencies for the LEDs is wide:  
 297 from tests conducted, frequencies higher than 4kHz and  
 298 lower than 40kHz work best. To detect lower frequencies,  
 299 biases have to be adapted to maximize the signal-to-noise  
 300 ratio. As the timestamp is quantized, it is advisable to use  
 301 LED frequencies with an integer microsecond period.

302 Due to the limited noise, detection can be simplified by  
 303 using only single types of events. In [28] and [4], the de-  
 304 tections rely on the transitions between event polarities. In-  
 305 stead, in the proposed approach, we use the timing informa-  
 306 tion between consecutive events generated by a single pixel.

307 For detection, an event frame generated over the period  
 308  $T_d$  is used, where  $T_d$  has to be larger than  $\frac{2}{f_{\min}}$  to ensure that  
 309 at least two blinks are visible for every LED. Candidates  
 310 for the blinking LEDs can be retrieved by selecting the con-  
 311 nected regions where more than  $T_d f_{\min}$  events per pixel are  
 312 generated. Each region with an area larger than a defined  
 313 minimal area is selected as a potential candidate. Due to  
 314 the short accumulation period, even under fast motion, the  
 315 LEDs' center points can be calculated by computation of  
 316 the center of mass on a 2D plane. The error introduced by  
 317 the relative movement of the LED is refined by the tracking  
 318 procedure.

319 For the frequency estimation of the LEDs, a histogram of  
 320 the time differences between events in a given area is used.

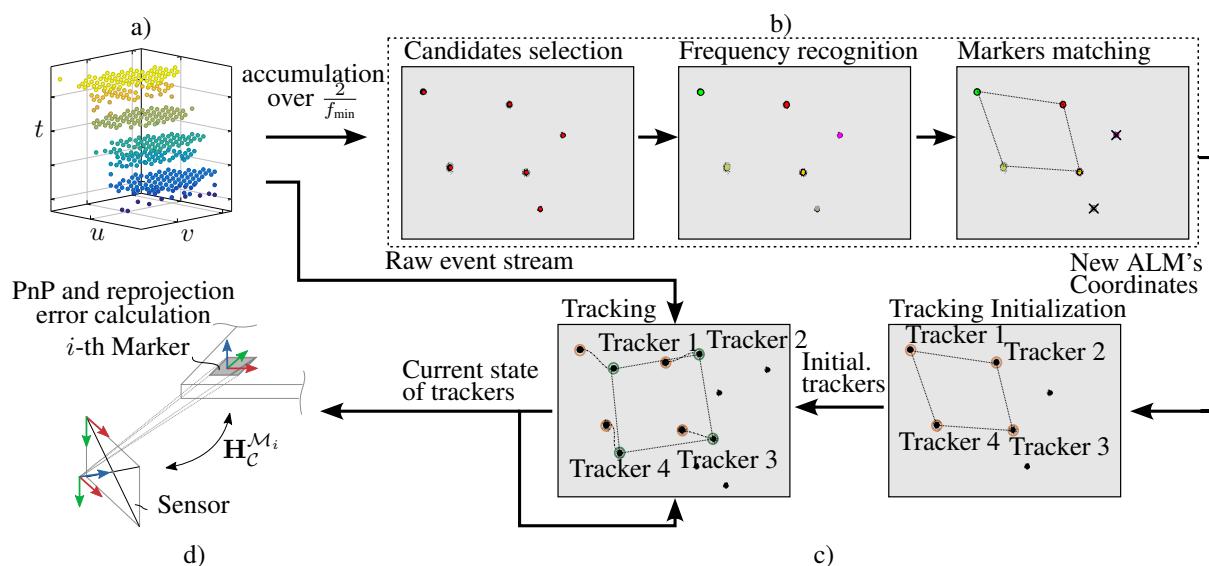


Figure 2. Overview of the proposed approach. Our pipeline consists of four asynchronous parts. *First*, illustrated in (a), to reduce the noise in the signal, biases are tuned to produce a single event per pixel on every blink. The proposed detection algorithm accumulates the events over a short period of time (two times the period of the LED's minimal frequency  $f_{\min}$  to ensure at least two blinks are visible for every LED). *Second*, for the accumulated blinks, frequencies are recognized and assigned to the specified markers. For every detected marker (b), trackers (c) are spawned for each individual LED. *Third*, using a simple tracking procedure, the LEDs are being tracked independently. *Fourth*, to quantify tracking quality during runtime, the resulting solution (d) is used to compute the reprojection errors.

In the case of frequencies with an integer microsecond period, the histogram has a pronounced peak, while for other frequencies, the histogram follows a wider Gaussian distribution. The frequency estimation follows the procedure proposed in [4].

### 3.2. Tracking

While detection relies on an accumulated representation of the events, the tracking can be performed online to reduce latency. Using the initial guess from the detection of the ALM's LED center points, trackers are spawned for every LED. The  $i$ -th tracker is characterized by its frequency  $f_i$ , center point  $\mathbf{c}_i = [x_i, y_i]$ , and radius  $r_i$ . In comparison to the assumptions of [17] and [25], the distribution of the generated events (within one blink) follows a spatially uniform distribution and hence, produces a dense event stream in this region. This allows us to simplify the tracking algorithm while maintaining precise tracking with sub-pixel accuracy.

For every LED's blink, the tracker's center of mass  $\bar{\mathbf{c}}_i$  is calculated using all events within its current radius  $r_i$ . The update term

$$\bar{\mathbf{c}}_i = \tilde{\beta} \mathbf{u}_k + (1 - \tilde{\beta}) \bar{\mathbf{c}}_i \quad (2)$$

introduces low-pass filtering, where every new event  $\mathbf{u}_k$  updates the current solution directly with an update factor  $\tilde{\beta}$  of typically 0.02. The radius  $r_i$  is updated every  $N$  events and set to twice the average distance of the events from the cen-

ter point of the tracker.

### 3.3. Pose Estimation

The 6-DoF pose is estimated asynchronously, using the current center points of the ALM's trackers. To increase the update rate of the algorithm, a PnP algorithm is started whenever the previous iteration is done. Due to the simplicity of the tracking, the PnP calculation is decisive in terms of latency and output rate.

To ensure stability and detect tracking failures, the reprojection error is computed and compared to the tracker's center points. When the reprojection error of one tracker exceeds the mean distance of the events from the center point, a tracking lost signal is generated, and tracking is stopped. It is reinitialized with the first new detection of a given marker.

## 4. Experiments

For the experimental setup, the EVK4 HD evaluation kit from Prophesee is used. It includes the event-based vision sensor IMX636ES providing HD resolution ( $1280 \times 720$  pixels) and the Soyo SFA0820-5M lens. The ALM consists of printed circuit boards with 8 LEDs arranged in a square of 9 cm side length. Each ALM has a base frequency (first LED), and the remaining frequencies are selected to match integer microsecond period times. For the experiments, four markers are arranged in a square on a marker board with a

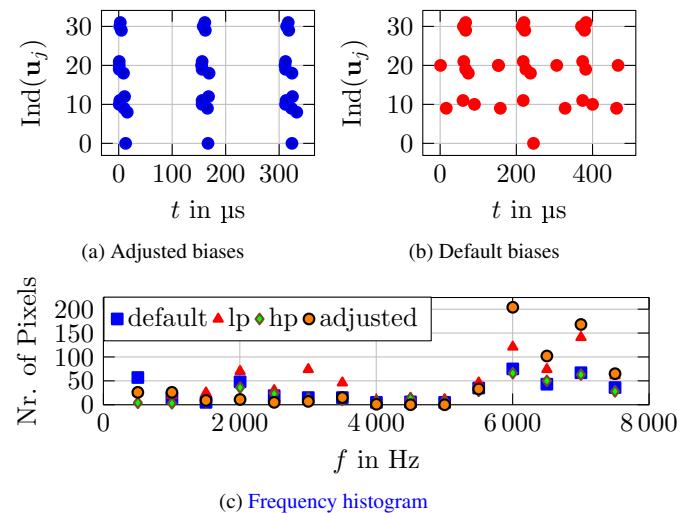


Figure 3. Visualization of bias adjustment. The vertical axis in (a) and (b) represents the flattened indices of the pixels in the region of interest (ROI) around an LED light. Plot (a) and (b) demonstrate the event distribution with optimal and default bias settings, respectively. In the frequency histogram (c), the effect of the low-pass (lp) and high-pass (hp) settings are illustrated beside the default and adjusted bias settings.

side length of 59 cm. The 8 outermost LEDs were chosen to create a single marker. The event stream is processed on a Desktop PC (Ubuntu 20.04, Intel i9-12900K, 32GB RAM) and on an Intel Aero Compute Board. As ground truth, the commercial infrared 3D tracking system OptiTrack is used. The OptiTrack recordings are triggered with the same trigger signal as the event-based camera via its external trigger input.

#### 4.1. Bias Adjustment

The bias adjustment of the event camera is essential for the proposed system's performance. The IMX636ES sensor biases [1] allow control over analog pixel gate thresholds to achieve the desired sensor response. The adjustment goal is to minimize the number of activated pixels between two LED blinks, as shown in Figure Fig. 3, thereby reducing processing complexity. The proposed method employs a single event polarity for simplicity.

By adjusting the refractory period setting, a pixel should be rendered insensitive to subsequent changes in LED brightness. An optimal value during adjustments should filter out all events between two consecutive LED-triggered events, as depicted in Fig. 3a. By utilizing high-pass and low-pass filter setups, the number of environment-generated events (excluding those by LEDs) can be limited to prevent sensor overflow and maintain manageable event blob density. The histogram of frequencies at which events occur over an accumulation time of 10 ms is depicted in Figure

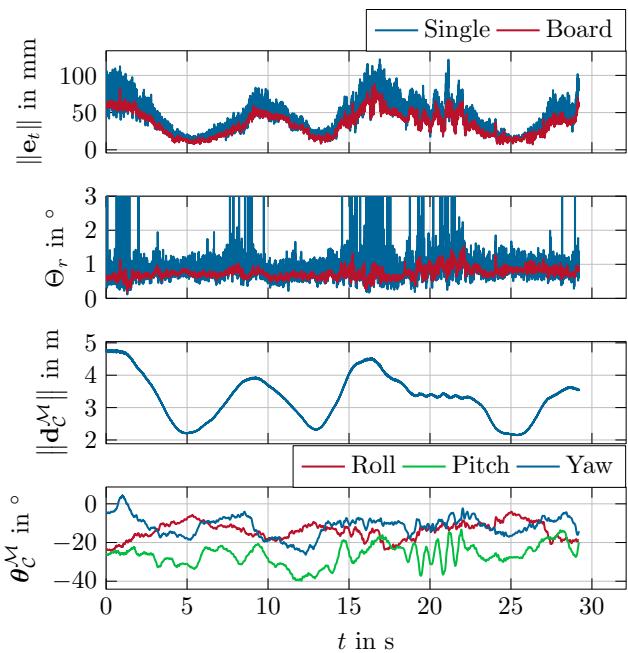


Figure 4. Analysis of the orientation and position error with reference distances and orientations.

3c. It illustrates the effect of the low-pass (lp) and high-pass (hp) bias settings acting as a band-pass filter for the LED frequencies of the ALM. Please note that the bias values may affect each other [1]. This causes the high event count, with the adjusted bias settings, at 6 kHz in Fig. 3c, where the value exceeds low-pass and high-pass settings.

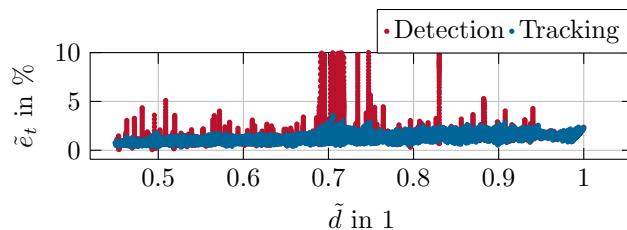
A detailed explanation of the sensor biases is provided in [1] and the description of the bias adjustment procedure is detailed in the supplementary document.

#### 4.2. Absolute Accuracy

To evaluate the absolute accuracy, the pose estimation of the ALMs and the marker board are compared with the synchronized measurements of the OptiTrack system (ground truth). For this experiment, the marker board is placed statically in the scene. The camera moves from close to far, covering the working distance of the setup, which is limited by the OptiTrack setup. The kinematic relations of the experimental setup are described in detail in the supplementary document.

The magnitude of the absolute position error  $\mathbf{e}_t$  and the orientation error  $\Theta_r$  with the distance between the marker and the camera  $\|\mathbf{d}_c^M\|$ , ranging from 2.1 m to 4.8 m, as well as the orientation  $\Theta_c^M$  of the marker with respect to the camera, is depicted in Fig. 4. Moreover, Fig. 4 illustrates the difference of using a single ALM with a side length of 9 cm for pose estimation compared to a larger marker board with a side length of 59 cm. From the plot of the position

432  
433  
434  
435  
436  
437  
438  
439  
440  
441  
442  
443  
444  
445  
446  
447  
448  
449  
450  
451  
452  
453  
454  
455  
456  
457  
458  
459  
460  
461  
462  
463  
464  
465  
466  
467  
468  
469  
470  
471  
472  
473  
474  
475  
476  
477  
478  
479  
480  
481  
482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493  
494  
495  
496  
497  
498  
499  
500  
501  
502  
503  
504  
505  
506  
507  
508  
509  
510  
511  
512  
513  
514  
515  
516  
517  
518  
519  
520  
521  
522  
523  
524  
525  
526  
527  
528  
529  
530  
531  
532  
533  
534  
535  
536  
537  
538  
539

Figure 5. Normalized relative position error  $\tilde{e}_t$ .

error  $\|\mathbf{e}_t\|$  in Fig. 4 it can be seen that the marker board has less noise but a comparable error magnitude. The second plot displaying the orientation error  $\Theta_r$  shows significant spikes for the single ALM. This indicates flips in the estimated pose, especially for medium to far distances. Hence, the usage of a marker board with increased side length is beneficial for accurate orientation estimations. The plot of the marker orientation  $\Theta_c^M$  in Fig. 4 shows fast orientation changes beginning at 20 s. This demonstrates the ability of the proposed method to accurately estimate pose information even in highly dynamic scenes.

In order to compare the performance between the detection-based pose estimation (Sec. 3.1) and the tracking-based pose estimation (Sec. 3.2), the relative position error  $\tilde{e}_t = \frac{\|\mathbf{e}_t\|}{\|\mathbf{d}_c^M\|}$  is displayed in Fig. 5 as a function of the normalized distance  $\tilde{d} = \frac{\|\mathbf{d}\|_c^M}{\max\|\mathbf{d}\|_c^M}$ . The results for the tracking-based approach indicate a better consistency, *i.e.* less noise, and altogether lower error numbers in the 1 % to 2 % range. The expected linear increase of the position error with the distance can also be inferred from Fig. 5. The statistical values of the data illustrated in Fig. 5 are summarized in Tab. 1. The maximum position error of 87.8 mm at a distance of 4.8 m and the maximum orientation error of 1.55° indicate the excellent performance of the tracking-based approach. The standard deviation of the position and orientation error of 16.2 mm and 0.146°, respectively, show the robustness of our method. In Table 4, we contextualize our results within different types of positioning system.

Table 1. Statistical values of the absolute accuracy measurements.

	Tracking		Detection	
	$\ \mathbf{e}_t\ $	$\Theta_r$	$\ \mathbf{e}_t\ $	$\Theta_r$
Mean	34.5 mm	0.738°	64.9 mm	1.55°
Std. Dev.	16.2 mm	0.146°	121 mm	5.12°
Maximum	87.8 mm	1.55°	1.233 m	71.9°

### 4.3. Static Noise

In Tab. 2, the noise floor of the proposed method is characterized for different distances between the camera and the marker board. The low standard deviation values indicate the stability of the pose estimation. This data can be uti-

lized to tune Bayesian filters (*e.g.* Kalman filters).

Table 2. Statistics of the noise in static scenes.

$\ \mathbf{d}_c^M\ $	Std. Dev.	Maximum
6 m	1.4 mm	5.5 mm
4 m	0.68 mm	2.95 mm
2 m	0.25 mm	2.17 mm

### 4.4. Latency Measurement and Output Rate

To determine the latency of the proposed system, the execution priority was elevated. Additionally, the visualization, as well as background tasks of the operating system, were disabled. This avoids unintentional interrupts and stalls during the execution of the pose estimation.

The latency and output rate values are listed in Tab. 3. The output rate of the tracking-based approach outperforms the detection-based method while achieving comparable latency results. As shown in Tab. 3, the proposed method is capable of running even on an embedded PC of a drone. While maintaining real-time performance, we can notice a reduced output rate (limited by the PnP computation time) as well as an increased average delay (limited by a number of concurrent threads) compared to a desktop PC. Latency is measured using a precise synchronization trigger signal and is equal to the time difference between the trigger and the time when pose estimation for this timestamp is available. Our proposed method achieves lower mean latency combined with low standard deviation compared to the state of the art. A large part of the resulting latency is due to communication overhead.

Table 3. Latency and output rates using a Desktop PC (PC) and Intel Aero Compute Board (Drone).

	Tracking		Detection	
	latency	rate	latency	rate
PC	354 µs	3.81 kHz	699 µs	670 Hz
	92 µs	64 Hz	35 µs	288 Hz
Drone	1232 µs	1.32 kHz	1953 µs	223 Hz
	194 µs	140 Hz	240 µs	94 Hz

### 4.5. Application to 6-DoF position estimation for a quadcopter

In this section, indoor and aggressive outdoor flights are considered for 6-DoF position estimation of a quadcopter.

#### 4.5.1 Indoor flight experiments

Compared to stationary robots, *e.g.* articulated manipulators [29, 30], which are equipped with high-precision encoders to monitor their state, flying robots mainly rely on IMUs, barometers, and vision-based systems to estimate their state. While the pose estimation module equipped with

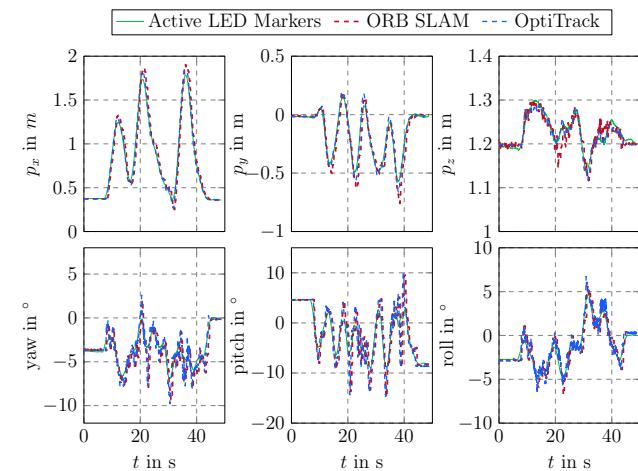
648 Table 4. Comparison of different marker based approaches and visual odometry/SLAM methods to contextualize the performance of the  
 649 system. All the results are taken from corresponding papers. \*Possible types E-Events, F-Frames, I-IMU. †Absolute Trajectory Error  
 650 (RMS) reported in [12] using "kitchen" sequence. ‡Sequences used for evaluation include only slow motion.  
 651

Method	Input*	Rate/FPS	Range	Markers	Positioning error	Resolution	Dynamic motion	Marker size
Ours	E	3.8 kHz	up to 10 m	Active	1.89 % 2.11 %	720 p	✓	59×59 cm 9×9 cm
Censi et al. 2013 [4]	E	250 Hz	-	Active	8.9 cm	128 × 128	✓	20×20 cm
Salah et al. 2022 [28]	E+I	200 Hz	up to 7 m	Active	0.074 %	480 p	✗‡	3×3×3 m
Chen et al. 2020 [6]	E	-	1 m (tests)	Active	3 %	346 × 240	✗	40×30 cm
STag [2]	F	56 FPS	up to 3 m	Passive	1.32 %	720 p	✗	15×15 cm
ORB-SLAM2 [23]	F	-	-	SLAM	13.0 cm†	480 p	✓	-
EDS [12]	E+F	-	-	VO	9.6 cm†	480 p	✓	-

663 only IMUs and barometers often suffers from the problem  
 664 of drift, vision-based systems ensure a more reliable measurement.  
 665 In this experiment, the Intel® Aero Ready to Fly (RTF) Drone, shown in Fig. 1, is employed as it offers  
 666 enough computational power for on-board processing.  
 667

668 For absolute drift-free pose information, the ORB-  
 669 SLAM2 [23] algorithm is utilized. The ORB-SLAM2 was  
 670 chosen for its impressive performance and open-source im-  
 671 plementation. Note that ORB-SLAM3 [3], as the suc-  
 672 cessor of ORB-SLAM2, is a more robust version compared to  
 673 ORB-SLAM2. However, the accuracy of these approaches  
 674 on stereo and RGB-D cameras is still comparable since the  
 675 key concepts of the estimation module and the relocalization  
 676 method, remain unchanged. We are aware that comparing  
 677 the proposed system with other SLAM algorithms, e.g.,  
 678 feature-based SLAM [23] and event-based SLAM [5], may  
 679 not be fair because the key concept is different. However,  
 680 this comparison could provide a qualitative guide for choos-  
 681 ing the right method in a given situation and contextualize  
 682 results as shown in Table 4. Similar to previous subsections,  
 683 the OptiTrack serves as the source of ground truth.  
 684

685 For the position errors, the metric for comparison is  
 686 the difference between a ground-truth position and the es-  
 687 timated position as  $\mathbf{p}_e = [p_{e,x} \ p_{e,y} \ p_{e,z}]^T = \mathbf{p}_g - \hat{\mathbf{p}}$ , where the hat ( $\hat{\cdot}$ ) indicates quantities estimated by  
 688 the ORB-SLAM2 algorithm [23] or with the event-based  
 689 marker, respectively, the subindex  $(\cdot)_g$  stands for the three-  
 690 dimensional ground-truth quantities, and the subindex  $(\cdot)_e$   
 691 for the resulting three-dimensional errors. The orientation  
 692 errors are represented as Euler angles. The difference be-  
 693 tween the ground-truth quaternion and the estimated quater-  
 694 nion is defined as  $\mathbf{q}_e = \hat{\mathbf{q}}^{-1} \otimes \mathbf{q}_g$ , with  $\otimes$  as the quater-  
 695 nion product. Subsequently, this error quaternion can be  
 696 transformed into an equivalent representation using three  
 697 angles, i.e., roll, pitch, and yaw. In this experiment, the  
 698 quadcopter is moved aggressively in a zigzag pattern in a  
 699 space of 1.8 m in the  $x$ -direction, 0.6 m in  $y$ -direction, and  
 700 0.4 m in  $z$ -direction for about 50 s. The position estimates  
 701 and the resulting errors for this case are illustrated in Fig.  
 702 6 and Fig. 7. The errors obtained from the proposed al-  
 703 gorithm in the  $x$ -,  $y$ - and  $z$ -direction are bounded within  
 704  $\pm 0.06$  m,  $\pm 0.08$  m, and  $\pm 0.02$  m, respectively, while larger  
 705 errors result from the ORB-SLAM2, i.e.,  $p_{e,x} = \pm 0.2$  m,  
 706  $p_{e,y} \in [-0.2, 0.1]^T$  m, and  $p_{e,z} = \pm 0.05$  m. The ori-  
 707 entation errors achieved with the two methods are similar,  
 708 shown at the bottom of Fig. 7. However, the orientation  
 709 errors measured by the proposed system are slightly better  
 710 since the spikes in ORB-SLAM2 are larger. More experi-  
 711 ments can be found in the supplementary material.  
 712



713 Figure 6. Time evolution of the drone trajectory.  
 714

715 0.4 m in  $z$ -direction for about 50 s. The position estimates  
 716 and the resulting errors for this case are illustrated in Fig.  
 717 6 and Fig. 7. The errors obtained from the proposed al-  
 718 gorithm in the  $x$ -,  $y$ - and  $z$ -direction are bounded within  
 719  $\pm 0.06$  m,  $\pm 0.08$  m, and  $\pm 0.02$  m, respectively, while larger  
 720 errors result from the ORB-SLAM2, i.e.,  $p_{e,x} = \pm 0.2$  m,  
 721  $p_{e,y} \in [-0.2, 0.1]^T$  m, and  $p_{e,z} = \pm 0.05$  m. The ori-  
 722 entation errors achieved with the two methods are similar,  
 723 shown at the bottom of Fig. 7. However, the orientation  
 724 errors measured by the proposed system are slightly better  
 725 since the spikes in ORB-SLAM2 are larger. More experi-  
 726 ments can be found in the supplementary material.  
 727

#### 728 4.5.2 Outdoor Flight Experiments

729 Outdoor experiments were conducted to demonstrate the  
 730 capability of the proposed system to detect and track  
 731 motions at very high speeds. In the first scenario, the  
 732 drone is equipped with an ALM and the event-based cam-  
 733 era is mounted vertically on a tripod on the ground, see  
 734

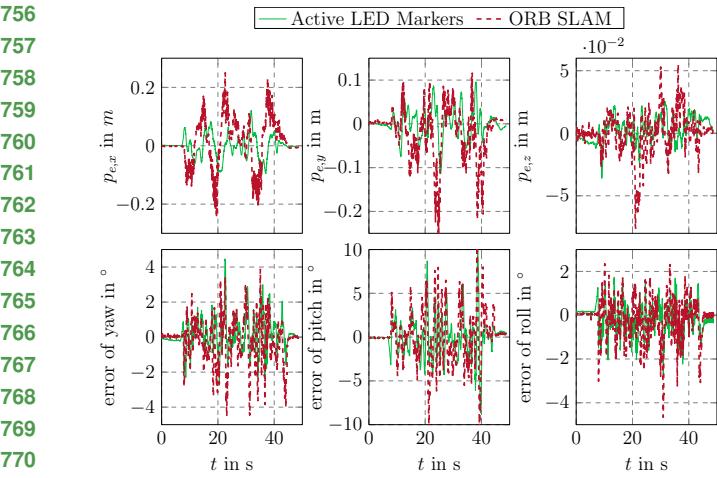


Figure 7. The error plots of the corresponding estimates, depicted in Fig. 6 with respect to the ground-truth measurements.

Fig. 8a. Although the drone moves at a maximum speed of  $4.5 \text{ m s}^{-1}$  and ascends to a height of 9 m, the proposed system is still able to capture the trajectories of the ALM, depicted in Fig. 8a. The position signals  $d_c^M$  indicate a low noise floor. The velocity signals  $\dot{d}_c^M$  are calculated based on the position signals with additional moving average filtering. Unlike in the first scenario, the camera is mounted on the drone and the ALM is static on the ground in the second scenario, as illustrated in Fig. 8b. The captured trajectories are shown in Fig. 8a when the drone is moving with an average speed of  $10 \text{ m s}^{-1}$ . In both scenarios, the velocity in the  $z$  direction is noisy due to the higher noise in the  $z$  estimation by the PnP algorithm. Live videos of the two scenarios are provided in the supplementary material.

#### 4.6. Limitations

The performance of the proposed system is mainly limited by the resolution of the event-based camera and the power of the light source. Based on the sensor resolution, the ground sampling distance (GSD) using the 8 mm lens at 10 m is equal to 0.61 cm. This sets an upper bound on the accuracy of the system, even with sub-pixel tracking precision and precisely calibrated camera. Light intensity decreases with the square of the distance [28]. Hence, LEDs with 1 W at 10% duty cycle were used in the experiments. Altogether, the maximum working distance is around 10m.

In favor of accuracy, ALMs with occluded LEDs are not tracked. However, occlusion of LEDs is detected to determine if the tracking is lost. An ALM is respawned when all LEDs are visible again. As the proposed system features a small marker size intended for outdoor usage, where other positioning systems are not applicable, LED occlusion is not a primary concern.

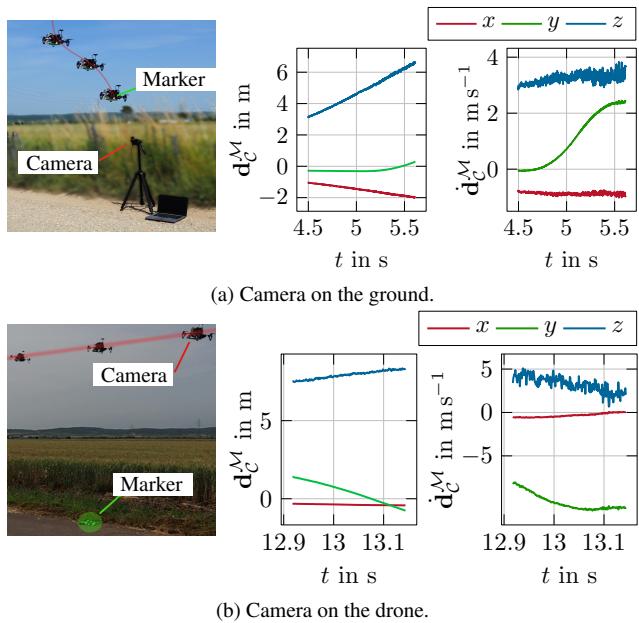


Figure 8. The plot of the drone's trajectory and velocities during experiments. The derivative of the positional signal  $d_c^M$  was filtered with a moving average filter with a window length of 100 samples.

## 5. Conclusion

This paper presents a fast and accurate vision-based localization system using an event-based camera with active led markers. Our proposed method overcomes the limits of traditional marker-based localization systems, *i.e.* low frame rate, motion blur, and high computational costs, by utilizing the advantages of an event-based camera. The proposed algorithm is simple but effective, achieving real-time performance with minimal latency below 0.5 ms and output rates above 3 kHz using a regular PC. The proposed tracking-based approach outperforms detection-based methods, especially in applications with very fast movements. The position error normalized to the distance is constantly below 1.87 % with a mean orientation error of 0.738°. To the best of the authors' knowledge, the combination of the achieved precision at this output rate and latency was not achieved so far. For the applications, where latency is not crucial, output of the system can be filtered and fused with other modalities (IMU, RGB based localisation systems). Proposed system can be used as a cheap relative localization/reference positioning system for outdoor application as data collection where other systems can not be applied (dynamic scenes, fast motion). Also the proposed method opens new possibilities for robotic applications where the high output rates and high precision of 6-DoF pose estimation are important, *e.g.* dynamic handover tasks and pick-and-place tasks.

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860  
861  
862  
863

864

## References

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

- [1] Biases — metavision SDK docs 4.1.0 documentation. 5
- [2] Burak Benligiray, Cihan Topal, and Cuneyt Akinlar. Stag: A stable fiducial marker system. *Image and Vision Computing*, 89:158–169, 2019. 7
- [3] Carlos Campos, Richard Elvira, Juan J Gómez Rodríguez, José MM Montiel, and Juan D Tardós. Orb-slam3: An accurate open-source library for visual, visual-inertial, and multimap slam. *IEEE Transactions on Robotics*, 37(6):1874–1890, 2021. 1, 7
- [4] Andrea Censi, Jonas Strubel, Christian Brandli, Tobi Delbrück, and Davide Scaramuzza. Low-latency localization by active LED markers tracking using a dynamic vision sensor. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 891–898, 2013. 3, 4, 7
- [5] William Chamorro, Joan Solà, and Juan Andrade-Cetto. Event-based line slam in real-time. *IEEE Robotics and Automation Letters*, 7(3):8146–8153, 2022. 7
- [6] Guang Chen, Wenkai Chen, Qianyi Yang, Zhongcong Xu, Longyu Yang, Jörg Conradt, and Alois Knoll. A novel visible light positioning system with event-based neuromorphic vision sensor. *IEEE Sensors Journal*, 20(17):10211–10219, 2020. 3, 7
- [7] Siyuan Chen, Dong Yin, and Yifeng Niu. A Survey of Robot Swarms’ Relative Localization Method. *Sensors*, 22(12):4424, 2022. 2
- [8] Toby Collins and Adrien Bartoli. Infinitesimal plane-based pose estimation. *International Journal of Computer Vision*, 109(3):252–286. 3
- [9] Tobi Delbrück and Manuel Lang. Robotic goalie with 3 ms reaction time at 4event-based dynamic vision sensor. *Frontiers in Neuroscience*, 7:223, 2013. 2
- [10] Guillermo Gallego, Tobi Delbrück, Garrick Orchard, Chiara Bartolozzi, Brian Taba, Andrea Censi, Stefan Leutenegger, Andrew J. Davison, Jorg Conradt, Kostas Daniilidis, and Davide Scaramuzza. Event-Based Vision: A Survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(1):154–180, 2022. 1, 2
- [11] Milad Heydariaan, Hessam Mohammadmoradi, and Omprakash Gnawali. Toward standard non-line-of-sight benchmarking of ultra-wideband radio-based localization. In *IEEE Workshop on Benchmarking Cyber-Physical Networks and Systems*, pages 19–24, 2018. 1
- [12] Javier Hidalgo-Carrió, Guillermo Gallego, and Davide Scaramuzza. Event-aided direct sparse odometry. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5781–5790, 2022. 7
- [13] Michail Kalaitzakis, Brennan Cain, Sabrina Carroll, Anand Ambrosi, Camden Whitehead, and Nikolaos Vitzilaios. Fiducial markers for pose estimation. *Journal of Intelligent & Robotic Systems*, 101(4):71, 2021. 2
- [14] Iman Abaspur Kazerouni, Luke Fitzgerald, Gerard Dooly, and Daniel Toal. A survey of state-of-the-art on visual slam. *Expert Systems with Applications*, 205:117734, 2022. 1
- [15] Oguz Kedilioglu, Tomás Marcelo Bocco, Martin Landesberger, Alessandro Rizzo, and Jörg Franke. Arucoe: En-

- hanced aruco marker. In *International Conference on Control, Automation and Systems*, pages 878–881, 2021. 1
- [16] Beat Kueng, Elias Mueggler, Guillermo Gallego, and Davide Scaramuzza. Low-latency visual odometry using event-based feature tracks. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 16–23, 2016. 2
- [17] Xavier Lagorce, Cédric Meyer, Sio-Hoi Ieng, David Filliat, and Ryad Benosman. Asynchronous event-based multikernel algorithm for high-speed visual features tracking. *IEEE Transactions on Neural Networks and Learning Systems*, 26(8):1710–1720, 2015. 4
- [18] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbrück. A  $128 \times 128$  120 dB  $15\mu s$  Latency Asynchronous Temporal Contrast Vision Sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, 2008. 1
- [19] K. Madsen, H.B. Nielsen, O. Tingleff, and Danmarks tekniske universitet. *Informatik og Matematisk Modellering. Methods for Non-linear Least Squares Problems*. Informatics and Mathematical Modelling, Technical University of Denmark, 2004. 3
- [20] Elias Mueggler, Guillermo Gallego, Henri Rebecq, and Davide Scaramuzza. Continuous-time visual-inertial odometry for event cameras. *IEEE Transactions on Robotics*, 34(6):1425–1440, 2018. 2
- [21] Elias Mueggler, Basil Huber, and Davide Scaramuzza. Event-based, 6-DOF pose tracking for high-speed maneuvers. In *IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 2761–2768, 2014. 2
- [22] Georg R. Muller and Jorg Conradt. A miniature low-power sensor system for real time 2D visual tracking of LED markers. In *IEEE International Conference on Robotics and Biomimetics*, pages 2429–2434, 2011. 3
- [23] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-source slam system for monocular, stereo, and rgbd cameras. *IEEE Transactions on Robotics*, 33(5):1255–1262, 2017. 7
- [24] Anh Nguyen, Thanh-Toan Do, Darwin G. Caldwell, and Nikos G. Tsagarakis. Real-time 6dof pose relocalization for event cameras with stacked spatial lstm networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 1638–1645, 2019. 1
- [25] David Reverter Valeiras, Xavier Lagorce, Xavier Clady, Chiara Bartolozzi, Sio-Hoi Ieng, and Ryad Benosman. An asynchronous neuromorphic event-driven visual part-based shape tracking. *IEEE Transactions on Neural Networks and Learning Systems*, 26(12):3045–3059, 2015. 4
- [26] Nasir Saeed, Haewoon Nam, Tareq Y Al-Naffouri, and Mohamed-Slim Alouini. A state-of-the-art survey on multidimensional scaling-based localization techniques. *IEEE Communications Surveys & Tutorials*, 21(4):3565–3583, 2019. 1
- [27] Wilson Sakpere, Michael Adeyeye Oshin, and Nhlanhla Mlitiwa. A state-of-the-art survey of indoor positioning and navigation systems and technologies. *South African Computer Journal*, 29:145, 2017. 2
- [28] Mohammed Salah, Mohammed Chehadah, Muhammed Hu-mais, Mohammed Wahbah, Abdulla Ayyad, Rana Azzam, 918
- 919
- 920
- 921
- 922
- 923
- 924
- 925
- 926
- 927
- 928
- 929
- 930
- 931
- 932
- 933
- 934
- 935
- 936
- 937
- 938
- 939
- 940
- 941
- 942
- 943
- 944
- 945
- 946
- 947
- 948
- 949
- 950
- 951
- 952
- 953
- 954
- 955
- 956
- 957
- 958
- 959
- 960
- 961
- 962
- 963
- 964
- 965
- 966
- 967
- 968
- 969
- 970
- 971

- 972 Lakmal Seneviratne, and Yahya Zweiri. A neuromorphic  
973 vision-based measurement for robust relative localization in  
974 future space exploration missions. *CoRR*, 2022. 3, 7, 8  
975 1026  
976 [29] Yvonne Stürz, Lukas Affolter, and Roy Smith. Parameter  
977 identification of the KUKA LBR iiwa robot including  
978 constraints on physical feasibility. *IFAC-PapersOnLine*,  
979 1027  
980 50(1):6863–6868, 2017. 6  
981 1028  
982 [30] Minh Nhat Vu, Florian Beck, Christian Hartl-Nesic, Anh  
983 Nguyen, and Andreas Kugi. Machine learning-based frame-  
984 work for optimally solving the analytical inverse kinematics  
985 for redundant manipulators. *Mechatronics*, 91:102970, 2023.  
986 1029  
987 6  
988 [31] Viktor Walter, Martin Saska, and Antonio Franchi. Fast Mu-  
989 tual Relative Localization of UAVs using Ultraviolet LED  
990 Markers. In *International Conference on Unmanned Aircraft  
991 Systems*, pages 1217–1226, 2018. 2  
992 1030  
993 [32] David Weikersdorfer and Jörg Conradt. Event-based particle  
994 filtering for robot self-localization. In *IEEE International  
995 Conference on Robotics and Biomimetics*, pages 866–870,  
996 2012. 2  
997 1031  
998 [33] Henk Wyneersch, Jiguang He, Benoit Denis, Antonio  
999 Clemente, and Markku Juntti. Radio localization and map-  
1000 ping with reconfigurable intelligent surfaces: Challenges,  
1001 opportunities, and research directions. *IEEE Vehicular Tech-  
1002 nology Magazine*, 15(4):52–61, 2020. 1  
1003 1041  
1004 [34] Xudong Yan, Heng Deng, and Quan Quan. Active infrared  
1005 coded target design and pose estimation for multiple objects.  
1006 In *IEEE/RSJ International Conference on Intelligent Robots  
1007 and Systems*, pages 6885–6890, 2019. 2  
1008 1042  
1009 [35] Masaki Yoshino, Shinichiro Haruyama, and Masao Nakagawa.  
1010 High-accuracy positioning system using visible led  
1011 lights and image sensor. In *IEEE Radio and Wireless Sym-  
1012 posium*, pages 439–442, 2008. 2  
1013 1043  
1014 [36] Xiang Zhang, Stephan Fronz, and Nassir Navab. Visual  
1015 marker detection and decoding in ar systems: A comparative  
1016 study. In *International Symposium on Mixed and Augmented  
1017 Reality*, pages 97–106, 2002. 1  
1018 1044  
1019 [37] Yong Zhao, Shibiao Xu, Shuhui Bu, Hongkai Jiang, and  
1020 Pengcheng Han. Gslam: A general slam framework and  
1021 benchmark. In *IEEE/CVF International Conference on Com-  
1022 puter Vision*, pages 1110–1120, 2019. 1  
1023 1045  
1024 1046  
1025 1047  
1026 1048  
1027 1049  
1028 1050  
1029 1051  
1030 1052  
1031 1053  
1032 1054  
1033 1055  
1034 1056  
1035 1057  
1036 1058  
1037 1059  
1038 1060  
1039 1061  
1040 1062  
1041 1063  
1042 1064  
1043 1065  
1044 1066  
1045 1067  
1046 1068  
1047 1069  
1048 1070  
1049 1071  
1050 1072  
1051 1073  
1052 1074  
1053 1075  
1054 1076  
1055 1077  
1056 1078  
1057 1079