

000
001
002054
055
056

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

057
058
059003
004
005
006
007060
061
062
063008
009
010
011064
065
066012
013
014067
068
069

Anonymous WACV Applications Track submission

070
071
072

Paper ID 282

073
074
075

Abstract

076
077
078

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 for this task 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: <https://almpose.github.io/>

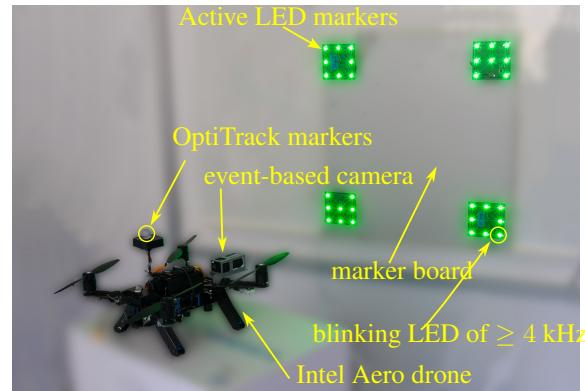
079
080
081

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 with respect to the event-based camera's frame of reference. An OptiTrack system is used for assessing the accuracy of the proposed method.

most common approaches 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 [10, 32]. Vision-based approaches are, however, computationally expensive and typically require more than one sensor, e.g. infrared-based systems [25]. To reduce the computational overhead, classical markers [13, 35] which serve 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 which is determined by the frame rate of the utilized cameras.

Event-based vision is an emerging field that has attracted much attention in recent years [9, 23]. Compared to conventional cameras, event-based cameras capture changes in light intensity asynchronously together with high temporal resolution and dynamic range. This feature makes them ideal for applications that require fast and accurate

1. Introduction

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 the safe, reliable and robust execution of corresponding tasks. Vision-based methods [2, 12, 36] are the

099
100
101102
103
104105
106
107

*These authors contributed equally to this work.

108 detection. Starting from the early years of event-based vi-
 109 sion development [8], advantages given by those sensors for
 110 robotics applications are noticeable. Due to the mediocre
 111 signal-to-noise ratio and low-resolution capabilities, applica-
 112 tions are limited and fall far behind standard RGB cam-
 113 eras in terms of performance. However, recent advan-
 114 cements in event-based sensor development [9] have enabled
 115 them to compete with the precision of other localization
 116 methods [6] due to increased resolution and reduced noise.
 117 Avoiding accumulated event representations (*i.e.* frames),
 118 markers can be tracked online utilizing the event-based
 119 camera's high temporal resolution of up to 1 μ s.
 120

In this paper, we propose a fast and simple method em-
 121 ploying an event-based camera together with active IED
 122 markers (ALMs) for simultaneous detection and tracking
 123 of the 6 degrees-of-freedom (DoF) pose of a rigid object in
 124 the 3D space. An overview of our proposed approach is de-
 125 picted in Fig. 1 with four ALMs attached to a marker board.
 126 The event-based camera is mounted on the drone, which is
 127 utilized to estimate the pose of the marker board with re-
 128 spect to the camera's base frame. To estimate the pose, the
 129 blinkings of the ALM's LEDs are logged with an event-
 130 based camera to identify the corresponding frequencies of
 131 each LED in the ALM. These blinking frequencies are uti-
 132 lized to identify each individual LED and match it with the
 133 known geometry of the ALM. With this mapping of the in-
 134 dividual points on the camera's sensor plane and the known
 135 geometry of the ALM, the pose of the ALM can be com-
 136 puted by utilizing a Perspective-n-Point (PnP) algorithm. In
 137 the presented approach, by tuning the biases, *i.e.*, parame-
 138 ters for tuning the analog front-end of the event-based cam-
 139 era, and using a priori knowledge about timing, the com-
 140 plexity of the ALM tracking can be simplified. This aids in
 141 reducing the tracking latency. During the tracking, the ini-
 142 tial detection is continuously refined, resulting in subpixel
 143 resolution. Such an approach can still precisely estimate
 144 the pose even under fast rotational and linear motion. Note
 145 that the output rate is independent of the tracking per-
 146 formance and is limited only by the Perspective-n-Point (PnP)
 147 algorithm [17] performance. Furthermore, the proposed ap-
 148 proach allows simultaneous tracking of multiple markers
 149 and is computationally undemanding, allowing it to be run
 150 on a wide range of devices. The proposed approach was
 151 tested and verified extensively using an external infrared-
 152 based positioning system. Our contributions are listed in
 153 the following.

- We propose the fast event-based pose estimation sys-
 tem 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 dy-
 namic scenarios for several in-depth aspects, *e.g.*, ab-
 solute accuracy, static noise, and latency. Translational

162 error achieved was equal to $34.5 \text{ mm} \pm 16 \text{ mm}$ and
 163 $0.74^\circ \pm 0.15^\circ$ orientation error at distances of 2.1 m to
 164 4.8 m between camera and marker, was achieved. To-
 165 gether with the fast computing speed, this proves that
 166 the proposed algorithm is promising for real-time ap-
 167 plications.

- 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 m s^{-1} .

The paper is organized as follows. Section 2 presents the
 177 related work in the field of pose estimation with event-based
 178 cameras and active markers. Section 3 describes the pro-
 179 posed method for marker detection and tracking. In Section
 180 4, we present the experimental setup and results. Finally,
 181 we conclude the paper in Section 5 with a summary of our
 182 contributions and suggestions for future work.

2. Related Work

Visual localization systems show improved accuracy
 186 compared to systems based on other physical principles
 187 [26], [6]. Fiducial marker-based systems [11] constitute the
 188 most common choice for robotic applications. Due to the
 189 limited range and the dependence on the lighting conditions,
 190 some studies proposed LED-based solutions based on stan-
 191 dard RGB cameras [34], infrared [33], or ultraviolet [30]
 192 spectrum. The pose can be retrieved much faster and more
 193 accurately due to the easy-to-extract anchor points visible in
 194 the image. However, the latency cannot be reduced beyond
 195 the camera's frame rate. This limitation can be avoided by
 196 using event-based cameras with their intrinsic high tempo-
 197 ral resolution.

One of the first works in the direction of localiza-
 199 tion based on event-based sensors was the 2D localiza-
 200 tion method [31]. The known shape (contours) was tracked and
 201 the relative localization was determined by event-based vi-
 202 sion. The high temporal resolution of the event-based sen-
 203 sors was used in [20] to localize an Unmanned Aerial Ve-
 204 hicle (UAV) during high-speed maneuvers. The full pose
 205 information was retrieved using a black square as a known
 206 shape. In [14] a visual odometry method was proposed
 207 based on the feature tracking algorithm. In this direction
 208 multiple methods were developed [14], [19], which show
 209 a significant improvement compared to the RGB-based ap-
 210 proach for high-speed applications.

The utilization of ALMs was proposed first in [21],
 212 where the authors tracked the 2D position of the LED and
 213 used it as a feedback signal for the robot homing and a pan-
 214 tilt system. Later, the first method for full pose estimation

216 using ALMs was presented in [3]. Therein, ALMs were
 217 used to detect and estimate the position of a flying quadrocopter.
 218 Taking advantage of the high temporal resolution of
 219 the camera, LEDs were recognized and detected using event
 220 polarity changes in the event stream. In [3], the authors
 221 used an accumulated event representation to decode the frequency
 222 and estimate the pose. In [5], a Gaussian mixture
 223 probability hypothesis density filter was proposed to localize
 224 the camera with respect to the active marker. Therein,
 225 the idea of online tracking was presented to increase the robustness
 226 and reliability of the pose estimation. The achieved
 227 results indicate a localization error lower than 3 cm in scenarios
 228 where the camera was within 1 m relative to the active marker.
 229 It shall moreover be noted that in [5] standard deviation or mean values are not reported.
 230

231 Most recent works using ALMs propose the additional
 232 fusion of inertial measurements [27]. The error in the predicted
 233 relative position is in the subcentimeter range. However, utilizing
 234 only the vision-based approach increases the error by the order of one magnitude. Compared to previous
 235 methods, the marker size is significantly larger. The LEDs
 236 are placed 1 m apart. Current work in the field of active
 237 marker-based solutions also focuses on the visual communication
 238 aspect of modulated light.

239 Different from other approaches in the literature, in the
 240 proposed approach, the complexity can be simplified by
 241 tuning the biases and using a priori knowledge about timing.
 242 This helps to reduce the tracking latency. To the best
 243 of the authors' knowledge, this work achieves the lowest
 244 latency compared to other methods in the literature. During
 245 tracking, initial detection is continuously refined, resulting
 246 in subpixel resolution. Such an approach can still precisely
 247 estimate the pose even under fast rotational and linear motion.
 248 The output rate is independent of the tracking performance
 249 and is limited only by the Perspective-n-Point (PnP)
 250 algorithm [17]. The method allows simultaneous tracking
 251 of multiple markers and is computationally undemanding,
 252 allowing it to be run on a wide range of devices. The pro-
 253 posed approach was tested and verified extensively using an
 254 external infrared-based positioning system.

257 3. Active Marker Tracking and Pose Estima- 258 tion

260 An event-based camera consists of an array of independent
 261 pixels measuring changes in luminosity $L = \log(I)$, based on the photocurrent I [16]. A change in the continuous
 262 luminosity signal

$$263 \Delta L(\mathbf{u}_k, t_k) = L(\mathbf{u}_k, t_k) - L(\mathbf{u}_k, t_k - \Delta t_k) > p_k C \quad (1)$$

264 triggers an event $\mathbf{e}_k = (\mathbf{u}_k, t_k, p_k)$ at pixel location $\mathbf{u}_k =$
 265 (u_k, v_k) due to a temporal contrast threshold $\pm C$, $p_k \in \{+1, -1\}$ being its polarity and Δt_k its time since the last

266 event (at \mathbf{u}_k) occurred at t_k [9]. The current generation
 267 of event-based sensors can produce up to 1.2 Giga events
 268 per second (Geps) [9] with a microsecond range timestamp
 269 accuracy. Compared to frame-based cameras that deliver
 270 periodically dense (*i.e.* full-frame) information, the event
 271 stream is sparse and contains information that relates only to
 272 changes in the scene. Using ALMs, the periodic and dense
 273 signals can be generated as a projection of the LED on the
 274 camera's sensor plane.

275 To reduce the computational complexity and required
 276 bandwidth, the biases of the sensor are tuned to generate a
 277 single event per pixel on every LED blink while suppressing
 278 all other background events to increase the signal-to-noise
 279 ratio, as presented in Figure 2. While [3] uses events of
 280 both polarities and relatively low frequencies (1-2kHz), the
 281 amount of noise can be reduced by using higher frequencies
 282 and disabling one polarity.

283 The ALM's structure is an arrangement of high-
 284 frequency blinking LEDs, where a unique frequency of
 285 the blinking pattern can individually recognize each LED
 286 (*e.g.* different blinking frequencies). The arrangement of
 287 the LEDs has to be fixed and determined in 3D coordinate
 288 space. However, it can also be arranged on a plane,
 289 as utilized in this work. Based on the 2D projection of the
 290 LEDs, knowing their 3D arrangement and camera intrin-
 291 sics, the relative pose of the marker with respect to the sen-
 292 sor can be reconstructed using a Point-n-Perspective (PnP)
 293 algorithm [18]. In this work, the IPPE PnP algorithm is
 294 used [7].

295 The proposed approach is divided into four parts, as il-
 296 lustrated in Fig. 2. To reduce the noise in the signal the bias
 297 settings are tuned to produce a single event per pixel on ev-
 298 ery blink of a LED. Next, events are accumulated over a pe-
 299 riod of $\frac{2}{f_{\min}}$ and for the accumulated event clusters, frequen-
 300 cies are recognized. These frequencies are used to iden-
 301 tify newly appearing ALMs. For each of the ALM's LEDs,
 302 trackers are spawned that keep track of the LEDs' center
 303 points based on single events. The tracking of the LEDs is
 304 independent of the detection loop. The pose of the ALM is
 305 estimated utilizing the trackers of the ALM. The accuracy
 306 of the pose estimation can be obtained using the reprojec-
 307 tion error. When an ALM leaves the field of view or the
 308 reprojection error exceeds the defined maximal value, the
 309 corresponding trackers are deleted. If the ALM enters the
 310 field of view again, the detection respawns it. Such a two-
 311 folded approach reduces latency and increases the accuracy
 312 of the solution.

313 3.1. Detection

314 The detection of an ALM includes the recognition of an
 315 individual ALM and the estimation of the blinking frequen-
 316 cies for each LED. For this, the geometrical arrangement,
 317 as well as the blinking frequency of the ALM LEDs, has to

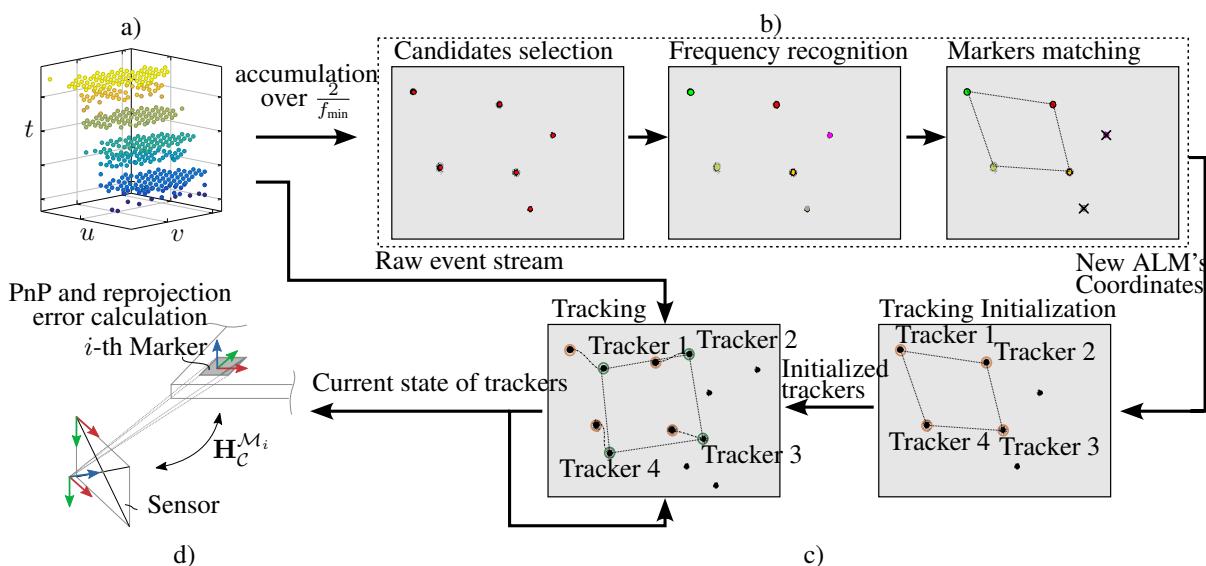


Figure 2. Overview of the proposed approach. To reduce the latency and ensure real-time processing, the pipeline consists of four asynchronous parts. *First*, illustrated in (a), to reduce the noise in the signal, biases are tuned to produce a single event per pixel on every blink. The proposed detection algorithm accumulates the events over a short period of time (two times the period of the minimal frequency LED to ensure at least two blinks are visible for every LEDs). *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.

be provided a priori.

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

Due to the limited noise, detection can be simplified by using only single types of events. While in [27] and [3], the detections rely on the transitions between event polarities. Instead, in the proposed approach, we rely only on the timing information between consecutive events generated by a single pixel.

For detection, an event frame generated over the period T_d is used. T_d has to be larger than $\frac{2}{f_{\min}}$ of the minimal used frequency f_{\min} (typ. $f_{\min} \approx 2\text{ kHz}$) to ensure that at least two blinks are visible for every LED. Candidates for the blinking LEDs can be retrieved by selecting the connected regions where more than $T_d f_{\min}$ events per pixel are generated. Each region with an area larger than a defined minimal area is selected as a potential candidate. Due to the short accumulation period, even under fast motion, the LEDs' center points can be calculated by computation of the center of mass on a 2D plane. The error introduced by the relative movement of the LED is refined by the tracking procedure.

For the frequency estimation of the LEDs, a histogram of

the time differences between events in a given area is used. 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 [3].

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 $c_i = [x_i, y_i]$ and radius r_i . In comparison to the assumptions of [15] and [24], 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{c}_i is calculated using all events within its current radius r_i . The update term

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

introduces low-pass filtering, where every new event \mathbf{u}_k updates the current solution directly with an update factor $\tilde{\beta}$ of

432 typically 0.02. The radius r_i is updated every N events and
 433 set to twice the average distance of the events with respect
 434 to the center point of the tracker.
 435

436 3.3. Pose Estimation

437 The 6-DoF pose is estimated asynchronously, using the
 438 current center points of the ALM's trackers. To increase
 439 the update rate of the algorithm, a PnP algorithm is started
 440 whenever the previous iteration is done. Due to the simplicity
 441 of the tracking, the PnP calculation is decisive in terms
 442 of latency and output rate.
 443

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

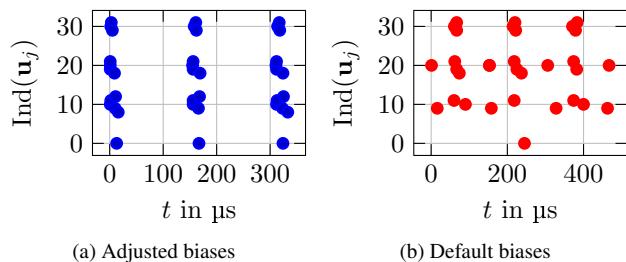
451 4. Experiments

452 4.1. Setup

453 For the experimental setup, the EVK4 HD evaluation kit
 454 from Prophesee is used. It includes the event-based vision
 455 sensor IMX636ES providing HD resolution (1280×720
 456 pixels) and the Soyo SFA0820-5M lens. The ALM consists
 457 of printed circuit boards with 8 LEDs arranged in a square
 458 of 9 cm side length. Each ALM has a base frequency (first
 459 LED), and the remaining frequencies are selected to match
 460 integer microsecond period times. For the experiments, four
 461 markers are arranged in a square resulting in a marker board
 462 with a side length of 59 cm. From this, 8 outermost LEDs
 463 were chosen to create a single marker. The camera's bias
 464 settings are set to minimize noise while maintaining a high
 465 event rate caused by the ALMs. The processing of the event
 466 stream is done on a PC running Ubuntu 20.04 with an Intel
 467 i9-12900K processor and 32GB of RAM and Intel Aero
 468 Compute Board. The camera is connected via USB 3 and
 469 the camera's external trigger input is used for time synchro-
 470 nization between the PC and the camera. As ground truth,
 471 the commercial infrared 3D tracking system OptiTrack is
 472 used, comprising nine Prime 17W cameras. The OptiTrack
 473 recordings are triggered with the same external trigger sig-
 474 nal as the event-based camera. In order to use the event-
 475 based camera simultaneously with the OptiTrack system,
 476 an infrared filter is used.
 477

478 4.2. Bias Adjustment

479 The bias adjustment of the event camera is essential for
 480 the proposed system's performance. The IMX636ES sensor
 481 biases [1] allow control over analog pixel gate thresholds to
 482 achieve the desired sensor response. The adjustment goal
 483 is to minimize the number of activated pixels between two
 484



486 Figure 3. Visualization of bias adjustment in an event camera us-
 487 ing the IMX636ES sensor. The vertical axis represents the flat-
 488 tened indices of the pixels in the region of interest (ROI) around an
 489 LED light. The left plot demonstrates an optimal bias adjustment,
 490 where event fronts - closely packed clusters of events triggered by
 491 a rapid change in the scene (like a sudden LED blink) - are clearly
 492 distinct. No unwanted spurious events occur between these event
 493 fronts. The right plot, in contrast, displays the event distribution
 494 with the camera's default bias settings, showing a less distinct sep-
 495 aration of event fronts.
 496

497 LED blinks, thereby reducing processing complexity. The
 498 proposed method employs a single event polarity for sim-
 499 plicity.
 500

501 By adjusting the refractory period setting, a pixel should
 502 be rendered insensitive to subsequent changes in LED
 503 brightness. An optimal value during adjustments should fil-
 504 ter out all events between two consecutive LED-triggered
 505 events. By utilizing high-pass and low-pass filter setups, the
 506 number of environment-generated events (excluding those
 507 by LEDs) can be limited to prevent sensor overflow and
 508 maintain manageable event blob density, as shown in Figure
 509 Fig. 3.
 510

511 A detailed explanation of sensor biases is provided in [1].
 512 The description of the bias adjustment procedure is given in
 513 our supplementary document. Although the focus here is
 514 on the IMX636ES sensor, the setup procedure should be
 515 applicable to cameras from other manufacturers.
 516

517 4.3. Absolute Accuracy

518 To evaluate the absolute accuracy, the pose estimation of
 519 the ALMs and the marker board are compared with the syn-
 520 chronized measurements of the OptiTrack system (ground
 521 truth). For this experiment, the marker board is placed stat-
 522 ically in the scene. The camera moves from close to far,
 523 covering the working distance of the setup, which is lim-
 524 ited by the OptiTrack setup. The kinematic relations of the
 525 experimental setup are described in detail in the appendix.
 526

527 The magnitude of the absolute position error e_t and the
 528 orientation error Θ_r with the distance between the marker
 529 and the camera $\|\mathbf{d}_c^M\|$, ranging from 2.1 m to 4.8 m, as
 530 well as the orientation Θ_c^M of the marker with respect to
 531 the camera, is depicted in Fig. 4. The difference of using
 532 a single ALM with a side length of 9 cm for pose estima-
 533

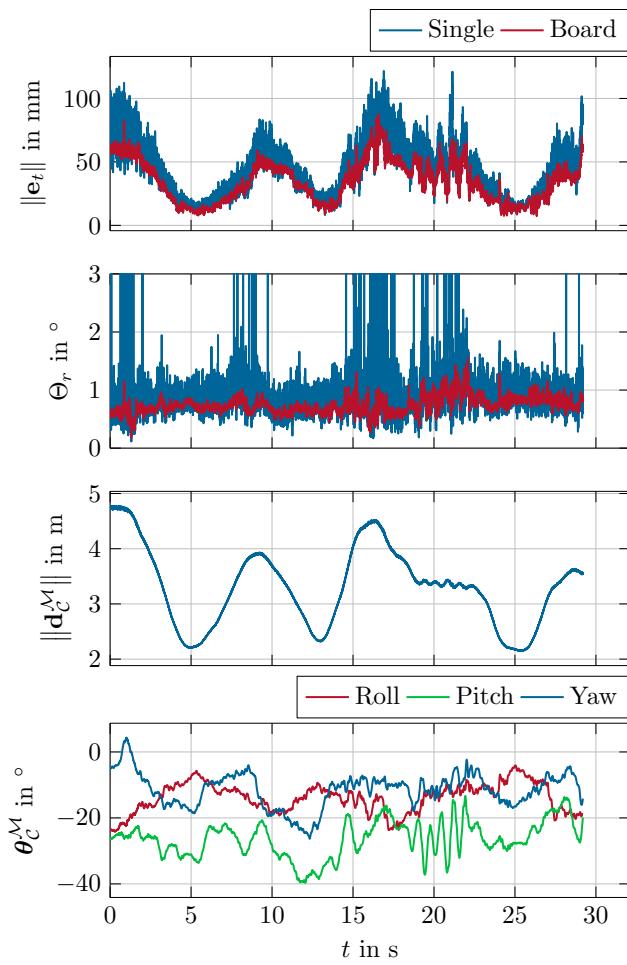


Figure 4. The ALM and the marker board are placed statically in the scene, and the camera moves from close to far. The first plot shows the absolute translational error of a single ALM and the whole marker board. The second plot displays the orientation error, with the third and fourth plot indicating the distance $\|\mathbf{d}_c^M\|$ ($\min \|\mathbf{d}\|_c^M = 2.1$ m, $\max \|\mathbf{d}\|_c^M = 4.8$ m) and the orientation Θ_c^M between the camera and the marker.

tion compared to larger marker board with a side length of 59 cm is illustrated. The plot of the position error $\|\mathbf{e}_t\|$ in Fig. 4 shows less noise but comparable error magnitude. The second plot displaying the orientation error Θ_r shows significant spikes for the single ALM curve. 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 Θ_c^M in Fig. 4 shows fast orientation changes beginning at 20 s. This demonstrates the ability of the proposed method to estimate pose information even in highly dynamic scenes accurately.

In order to compare the performance between the detection-based pose estimation (Sec. 3.1) and the tracking-

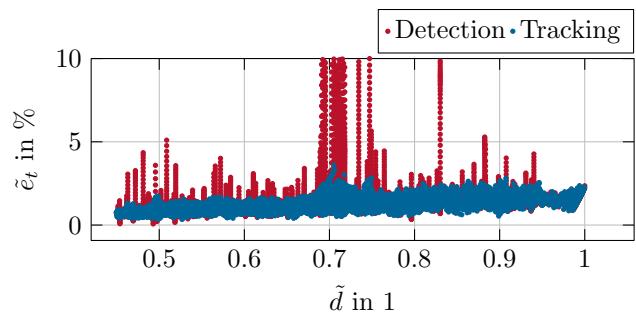


Figure 5. Comparison between the detection and tracking algorithm for pose estimation. The relative position error \tilde{e}_t is plotted over the normalized distance \tilde{d} .

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 with respect to 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 distance is displayed as well.

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, indicate the robustness of our method.

Table 1. Statistical values of the absolute accuracy measurements for the detection-based and the tracking-based approach.

	Tracking		Detection	
	$\ \mathbf{e}_t\ $	Θ_r	$\ \mathbf{e}_t\ $	Θ_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°

4.4. Static Noise

In Tab. 2, the noise floor of the proposed method is characterized at 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 utilized to tune Bayesian filters (*e.g.* Kalman filters) in the subsequent data processing in applications.

4.5. 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

648 Table 2. Statistics of the noise in static scenes at different distances
 649 between marker and camera.
 650

$\ \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

655
 656 disabled. This avoids unintentional interrupts and stalls during
 657 the execution of the pose estimation.
 658

659 The latency and output rate values are listed in Tab. 3 and
 660 Tab. 4. The output rate of the tracking-based approach out-
 661 performs the detection-based method while achieving com-
 662 parable latency results. As shown in Table Tab. 4, the pro-
 663 posed method is capable of running even on an embedded
 664 pc with an Intel Atom processor on an Intel Aero Ready to
 665 Fly Drone. While maintaining real-time performance, we
 666 can notice a reduced output rate (limited by the PnP com-
 667 putation time) as well as an increased average delay (limited
 668 by a number of concurrent threads) compared to a desktop
 669 PC. Latency is measured using a precise synchronization
 670 trigger signal and is equal to the time difference between
 671 the trigger and the time when pose estimation for this time-
 672 stamp is available.

673 Our proposed method achieves lower mean latency com-
 674 bined with low standard deviation compared to the state-of-
 675 the-art. To further reduce the latency it is possible to imple-
 676 ment this solution on different hardware architectures, *e.g.*
 677 FPGAs. A large part of the resulting latency is due to com-
 678 munication overhead.

679 Table 3. Latency and output rates for the detection-based and the
 680 tracking-based approach. (Desktop PC)
 681

	Tracking		Detection	
	latency	rate	latency	rate
Mean	354 μ s	3.805 kHz	699 μ s	0.670 kHz
Std. Dev.	92 μ s	0.064 kHz	35 μ s	0.288 kHz

688 Table 4. Latency and output rates for the detection-based and the
 689 tracking-based approach. (Intel Aero Compute Board)
 690

	Tracking		Detection	
	latency	rate	latency	rate
Mean	1232 μ s	1.32 kHz	1953 μ s	0.223 kHz
Std. Dev.	194 μ s	0.14 kHz	240 μ s	0.094 kHz

697 4.6. Application to 6-DoF position estimation for a 698 quadcopter 699

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

4.6.1 Indoor flight experiments

702 Compared to stationary robots, e.g. articulated manipula-
 703 tors [28, 29], which are equipped with high-precision en-
 704 coders to monitor their state, flying robots mainly rely on
 705 IMUs, barometers, and vision-based systems to estimate
 706 their state. While the pose estimation module equipped with
 707 only IMUs and barometers often suffers from the problem
 708 of drift, vision-based systems ensure a more reliable mea-
 709 surement. In this experiment, the Intel® Aero Ready to
 710 Fly (RTF) Drone, shown in Fig. 1, is employed as it of-
 711 fers enough computational power for on-board processing
 712 of all algorithms.

713 For absolute drift-free pose information, the ORB-
 714 SLAM2 [22] algorithm is utilized. The ORB-SLAM2 was
 715 chosen for its impressive performance and open-source im-
 716 plementation. Note that ORB-SLAM3 [2], as the suc-
 717 cessor of ORB-SLAM2, is a more robust version compared to
 718 ORB-SLAM2. However, the accuracy of these approaches
 719 on stereo and RGB-D cameras is still comparable since the
 720 key concepts of the estimation module and the relocaliza-
 721 tion method, remain unchanged. We are aware that compar-
 722 ing the proposed system with other SLAM algorithms, e.g.,
 723 feature-based SLAM [22] and event-based SLAM [4], may
 724 not be a fair comparison because the key concept is differ-
 725 ent. However, this comparison could provide a qualitative
 726 guide for choosing the right methods in a given situation
 727 and contextualize results. Similar to previous subsections,
 728 the OptiTrack serves as the source of ground truth.

729 For the translational errors, the metric for comparison
 730 is the difference between a ground truth position and the
 731 estimated 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
 732 the ORB-SLAM2 algorithm [22] or with the event-based
 733 marker, respectively, the subindex $(\cdot)_g$ stands for the three-
 734 dimensional ground truth quantities, and the subindex $(\cdot)_e$
 735 for the resulting three-dimensional errors. The orientation
 736 errors are represented in Euler angles. The difference be-
 737 tween the ground truth quaternion and the estimated quater-
 738 nion is defined as $\mathbf{q}_e = \hat{\mathbf{q}}^{-1} \otimes \mathbf{q}_g$ with \otimes is the quater-
 739 nion product. Subsequently, this error quaternion can be
 740 transformed into an equivalent representation using three
 741 angles, i.e., roll, pitch, and yaw. In this experiment, the
 742 quadcopter is moved aggressively in a zigzag pattern in a
 743 space of 1.8 m in the x -direction, 0.6 m in y -direction, and
 744 0.4 m in z -direction for about 50 s. The position estimates
 745 and the resulting errors for this case are illustrated in Fig.
 746 6 and Fig. 7. The errors obtained from the proposed al-
 747 gorithm in the x -, y - and z -direction are bounded within
 748 ± 0.06 m, ± 0.08 m, and ± 0.02 m, respectively, while larger
 749 errors result from the ORB SLAM, i.e., $p_{e,x} = \pm 0.2$ m,
 750 $p_{e,y} \in [-0.2, 0.1]^T$ m, and $p_{e,z} = \pm 0.05$ m. The ori-
 751 entation errors achieved with the two methods are similar,
 752 shown at the bottom of Fig. 7. However, the orientation
 753

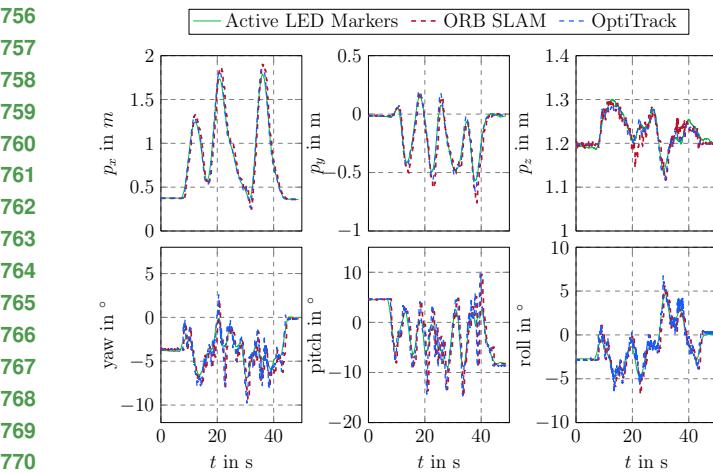


Figure 6. Time evolution of the drone trajectory. The estimated poses from the proposed method, the ORB SLAM, and the ground truths from OptiTrack are illustrated in green, red, and blue, respectively.

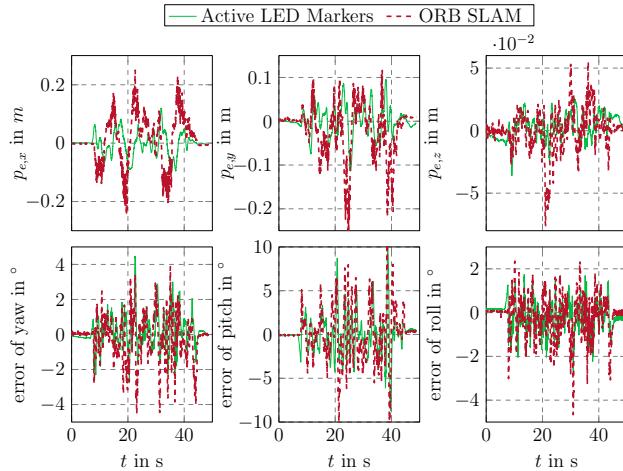


Figure 7. The error plots of the corresponding estimates, depicted in Fig. 6 with respect to the ground truth measurements from the OptiTrack.

errors measured by the proposed system are slightly better since the spikes in ORB SLAM are larger. More experiments can be found in our supplementary material.

4.6.2 Outdoor Flight Experiments

To illustrate the applicability and the robustness of the proposed method to real-world applications, outdoor experiments were conducted. In the first scenario, the drone was equipped with an ALM and the camera on the ground facing vertically. The trajectory during a vertical ascend of the drone is depicted in Fig. 8. The position signals d_c^M indicate a low noise floor. The velocity signals \dot{d}_c^M are calcu-

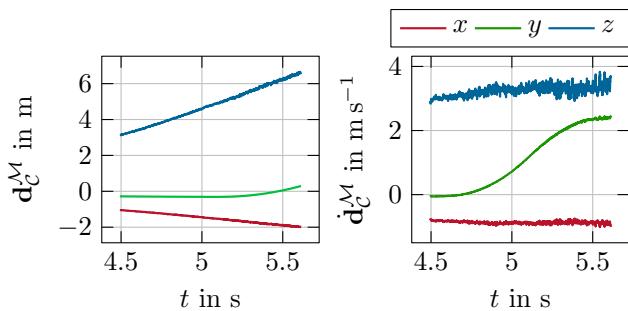


Figure 8. The plot of the drone's trajectory during a vertical ascend. The left plot illustrates the position d_c^M with the corresponding velocity \dot{d}_c^M shown in the right plot. The derivative of the positional signal was filtered with a moving average filter with window length of 500 samples to obtain the plotted velocity signal.

lated based on the position signals with additional moving average filtering. The velocity in the z -direction is noisier due to the higher noise in the z estimation by the PnP algorithm. In the second scenario, the camera is mounted on the drone, and the ALM is static on the ground. The results of these two scenarios are illustrated in the supplemental video.

TODO: we run into the problem here by giving the plot without expression! can you include two teaser picture of the two scenarios here (in the video, before each scenario, I put a teaser picture). I suggest to remove Fig. 8. Just only with two teaser images, you just explain that during the 5-10m/s flight, you can still capture the pose. If you want to explain Fig. 8, there should be the FIg.9 for the second scenario. But also you need to have 2 teaser picture, because the reader does not understand what we are doing here. SO.. I think the best way is to provide two teaser figures for scenario 1 and 2 and remove fig. 8 and 9

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 achieves superior performance over detection-based methods, especially under fast motion. 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

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 was not achieved so far. The proposed method enables new
865 possibilities for robotic applications, where the high output
866 rates of the 6-DoF pose with high precision are crucial, e.g.
867 dynamic handover tasks and pick-and-place tasks with high
868 precision.
869

References

- 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
- [1] Biases — metavision SDK docs 4.1.0 documentation. 5
 - [2] 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
 - [3] 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
 - [4] 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
 - [5] 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
 - [6] Siyuan Chen, Dong Yin, and Yifeng Niu. A Survey of Robot Swarms’ Relative Localization Method. *Sensors*, 22(12):4424, 2022. 2
 - [7] Toby Collins and Adrien Bartoli. Infinitesimal plane-based pose estimation. 109(3):252–286. 3
 - [8] 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
 - [9] 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, 3
 - [10] 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
 - [11] Michail Kalitzakis, 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
 - [12] 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
 - [13] Oguz Kedilioglu, Tomás Marcelo Bocco, Martin Landesberger, Alessandro Rizzo, and Jörg Franke. Arucoe: Enhanced aruco marker. In *International Conference on Control, Automation and Systems*, pages 878–881, 2021. 1
 - [14] 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
 - [15] Xavier Lagorce, Cédric Meyer, Sio-Hoi Ieng, David Filliat, and Ryad Benosman. Asynchronous event-based multikernel algorithm for high-speed visual features tracking. *IEEE*

- 972 *Transactions on Neural Networks and Learning Systems*,
 973 26(8):1710–1720, 2015. 4
- 974 [16] Patrick Lichtsteiner, Christoph Posch, and Tobi Delbrück.
 975 A 128×128 120 dB $15\mu s$ Latency Asynchronous Temporal
 976 Contrast Vision Sensor. *IEEE Journal of Solid-State Circuits*, 43(2):566–576, 2008. 3
- 977 [17] Xiao Xin Lu. A review of solutions for perspective-n-point
 978 problem in camera pose estimation. In *Journal of Physics: Conference Series*, volume 1087, page 052009, 2018. 2, 3
- 979 [18] K. Madsen, H.B. Nielsen, O. Tingleff, and Danmarks
 980 tekniske universitet. Informatik og Matematisk Modellering.
 981 *Methods for Non-linear Least Squares Problems*. Informatics
 982 and Mathematical Modelling, Technical University of
 983 Denmark, 2004. 3
- 984 [19] Elias Mueggler, Guillermo Gallego, Henri Rebucq, and Da-
 985 vide Scaramuzza. Continuous-time visual-inertial odome-
 986 try for event cameras. *IEEE Transactions on Robotics*,
 987 34(6):1425–1440, 2018. 2
- 988 [20] Elias Mueggler, Basil Huber, and Davide Scaramuzza.
 989 Event-based, 6-DOF pose tracking for high-speed maneu-
 990 vers. In *IEEE/RSJ International Conference on Intelligent*
 991 *Robots and Systems*, pages 2761–2768, 2014. 2
- 992 [21] Georg R. Muller and Jorg Conradt. A miniature low-power
 993 sensor system for real time 2D visual tracking of LED mark-
 994 ers. In *IEEE International Conference on Robotics and*
 995 *Biomimetics*, pages 2429–2434, 2011. 2
- 996 [22] Raul Mur-Artal and Juan D Tardós. Orb-slam2: An open-
 997 source slam system for monocular, stereo, and rgbd cam-
 998 eras. *IEEE Transactions on Robotics*, 33(5):1255–1262,
 999 2017. 7
- 1000 [23] Anh Nguyen, Thanh-Toan Do, Darwin G. Caldwell, and
 1001 Nikos G. Tsagarakis. Real-time 6dof pose relocalization
 1002 for event cameras with stacked spatial lstm networks. In
 1003 *IEEE/CVF Conference on Computer Vision and Pattern*
 1004 *Recognition Workshops*, pages 1638–1645, 2019. 1
- 1005 [24] David Reverter Valeiras, Xavier Lagorce, Xavier Clady,
 1006 Chiara Bartolozzi, Sio-Hoi Ieng, and Ryad Benosman. An
 1007 asynchronous neuromorphic event-driven visual part-based
 1008 shape tracking. *IEEE Transactions on Neural Networks and*
 1009 *Learning Systems*, 26(12):3045–3059, 2015. 4
- 1010 [25] Nasir Saeed, Haewoon Nam, Tareq Y Al-Naffouri, and
 1011 Mohamed-Slim Alouini. A state-of-the-art survey on mul-
 1012 tidimensional scaling-based localization techniques. *IEEE*
 1013 *Communications Surveys & Tutorials*, 21(4):3565–3583,
 2019. 1
- 1014 [26] Wilson Sakpere, Michael Adeyeye Oshin, and Nhlanhla Mi-
 1015 litwa. A state-of-the-art survey of indoor positioning and nav-
 1016 igation systems and technologies. *South African Computer*
 1017 *Journal*, 29:145, 2017. 2
- 1018 [27] Mohammed Salah, Mohammed Chehadah, Muhammed Hu-
 1019 maïs, Mohammed Wahbah, Abdulla Ayyad, Rana Azzam,
 1020 Lakmal Seneviratne, and Yahya Zweiri. A neuromorphic
 1021 vision-based measurement for robust relative localization in
 1022 future space exploration missions. *CoRR*, 2022. 3, 4
- 1023 [28] Yvonne Stürz, Lukas Affolter, and Roy Smith. Param-
 1024 eter identification of the KUKA LBR iiwa robot includ-
 1025 ing constraints on physical feasibility. *IFAC-PapersOnLine*,
 50(1):6863–6868, 2017. 7
- 1026 [29] Minh Nhat Vu, Florian Beck, Christian Hartl-Nesic, Anh
 1027 Nguyen, and Andreas Kugi. Machine learning-based frame-
 1028 work for optimally solving the analytical inverse kinematics
 1029 for redundant manipulators. *Mechatronics*, 91:102970, 2023.
 1030 7
- 1031 [30] Viktor Walter, Martin Saska, and Antonio Franchi. Fast Mu-
 1032 tual Relative Localization of UAVs using Ultraviolet LED
 1033 Markers. In *International Conference on Unmanned Aircraft*
 1034 *Systems*, pages 1217–1226, 2018. 2
- 1035 [31] David Weikersdorfer and Jörg Conradt. Event-based particle
 1036 filtering for robot self-localization. In *IEEE International*
 1037 *Conference on Robotics and Biomimetics*, pages 866–870,
 2012. 2
- 1038 [32] Henk Wyneersch, Jiguang He, Benoit Denis, Antonio
 1039 Clemente, and Markku Juntti. Radio localization and map-
 1040 ping with reconfigurable intelligent surfaces: Challenges,
 1041 opportunities, and research directions. *IEEE Vehicular Tech-*
 1042 *nology Magazine*, 15(4):52–61, 2020. 1
- 1043 [33] Xudong Yan, Heng Deng, and Quan Quan. Active infrared
 1044 coded target design and pose estimation for multiple objects.
 1045 In *IEEE/RSJ International Conference on Intelligent Robots*
 1046 *and Systems*, pages 6885–6890, 2019. 2
- 1047 [34] Masaki Yoshino, Shinichiro Haruyama, and Masao Nak-
 1048 agawa. High-accuracy positioning system using visible led
 1049 lights and image sensor. In *IEEE Radio and Wireless Sym-*
 1050 *posium*, pages 439–442, 2008. 2
- 1051 [35] Xiang Zhang, Stephan Fronz, and Nassir Navab. Visual
 1052 marker detection and decoding in ar systems: A comparative
 1053 study. In *International Symposium on Mixed and Augmented*
 1054 *Reality*, pages 97–106, 2002. 1
- 1055 [36] Yong Zhao, Shibiao Xu, Shuhui Bu, Hongkai Jiang, and
 1056 Pengcheng Han. Gslam: A general slam framework and
 1057 benchmark. In *IEEE/CVF International Conference on Com-*
 1058 *puter Vision*, pages 1110–1120, 2019. 1