# DEVO: Depth-Event Camera Visual Odometry in Challenging Conditions

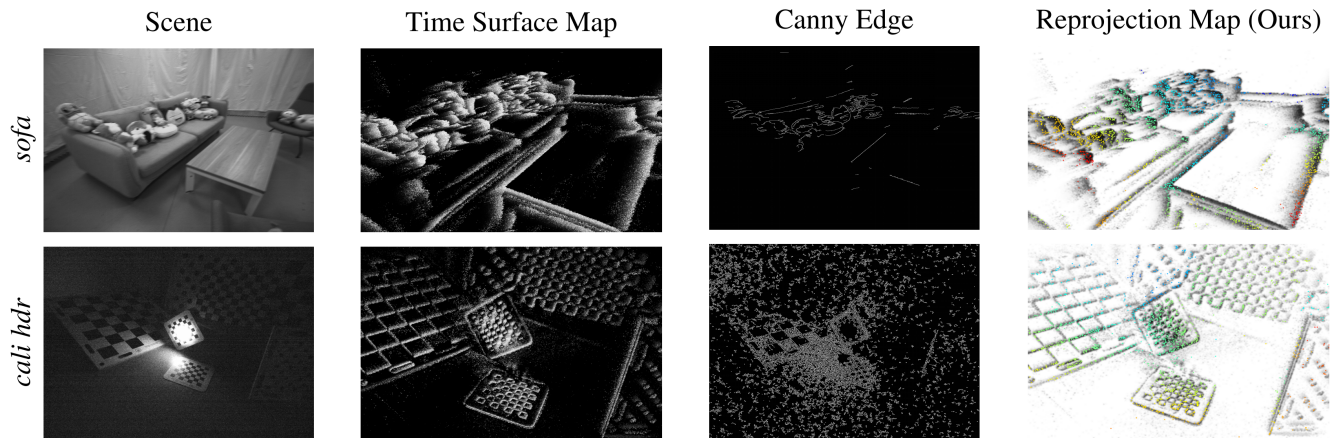Yi-Fan Zuo[1*], Jiaqi Yang[2*], Jiaben Chen[2], Xia Wang[1], Yifu Wang[2†] and Laurent Kneip[2†]

Fig. 1: *Challenging scenarios and results. sofa* is captured under high dynamics. *cali hdr* is a challenging illumination scene with highly self-similar texture. Column 1: Example RGB images. Column 2: Corresponding time-surface map. Column 3: Canny edge detections. Column 4: Reprojected depth points.

*Abstract*—**We present a novel real-time visual odometry framework for a stereo setup of a depth and high-resolution event camera. Our framework balances accuracy and robustness against computational efficiency towards strong performance in challenging scenarios. We extend conventional edge-based semi-dense visual odometry towards time-surface maps obtained from event streams. Semi-dense depth maps are generated by warping the corresponding depth values of the extrinsically calibrated depth camera. The tracking module updates the camera pose through efficient, geometric semi-dense 3D-2D edge alignment. Our approach is validated on both public and self-collected datasets captured under various conditions. We show that the proposed method performs comparable to state-of-the-art RGB-D camera-based alternatives in regular conditions, and eventually outperforms in challenging conditions such as high dynamics or low illumination.**

## I. INTRODUCTION

Real-time localization and 3D mapping are increasingly important tasks to be solved in many emerging technologies such as robotics, intelligent transportation, and intelligence augmentation. Owing to their small scale and affordability, cameras are often considered as an exteroceptive sensor in such applications. Despite being attractive, pure vision-based solutions still lack robustness in more challenging conditions [1], [2], and are therefore often complemented by additional sensors such as an Inertial Measurement Unit

(IMU), wheel encoders, or a depth camera. Especially the addition of the latter has been highly popular in indoor applications since the advent of consumer-grade RGB-D cameras in 2010 (e.g. Microsoft Kinect). RGB-D cameras provide high frequency and high resolution depth images, which significantly improves accuracy and robustness of monocular visual odometry and SLAM methods [3], [4], [5], [6], [7]. However, most RGB-D camera solutions still rely on sparse feature extraction or photometric image alignment, which is why they cannot handle challenging conditions such as highly dynamic motion or low illumination. KinectFusion [8], relies exclusively on depth images, but the method is power hungry and demands high frame-rate depth images as well as GPU resources for the real-time execution of depth fusion and ICP algorithms.

The present work introduces a fresh approach to depth camera-supported indoor visual odometry (VO) for a power-efficient handling of challenging conditions such as high dynamics or low illumination. Our core idea consists of exchanging the RGB camera against a Dynamic Vision Sensor, which pairs high dynamic range (HDR) with high temporal resolution. The basic functionality of a DVS—also called an event camera—as well as its advantages and challenges are well explained in the original work of Brandli et al. [9] or the recent survey by Gallego et al. [10]. Our method relies on an approach similar to Canny-VO [7] and extracts edges from the event stream, assigns depth from the depth camera, and registers subsequent views by 3D-2D edge alignment. The ability of event cameras to see in almost complete darkness paired with their high

* indicates equal contribution, † indicates corresponding author
[1] Key Laboratory of Optoelectronic Imaging Technology and Systems, Ministry of Education, School of Optics and Photonics, Beijing Institute of Technology, Beijing 100081, China. [2]Mobile Perception Lab, ShanghaiTech University; L. Kneip is also with the Shanghai Engineering Research Center of Intelligent Vision and Imaging.

temporal resolution lead to excellent performance in the above mentioned challenging cases. At the same time, our method maintains high computational and energy efficiency owing to potentially low depth camera frame rates as well as semi-dense processing. Our contributions are as follows:

- We present DEVO, a novel visual odometry framework for a hybrid stereo setup of a depth and an event camera.
- The approach relies on thresholded time-surface maps for edge detection and semi-dense depth map extraction.
- Our method handles 6-DoF motion estimation, and we demonstrate high efficiency and successful operation in all conditions.

Our results are obtained on self-collected high-resolution RGB-D-event indoor datasets with ground truth captured by an external motion tracking system. Further tests are conducted on larger scale outdoor datasets where depth is obtained from a LiDAR scanner. A thorough comparison against state-of-the-art event-based and RGB-D based visual odometry frameworks proves that DEVO achieves high quality, continuous visual localization results and eventually outperforms alternative methods in challenging conditions.

## II. RELATED WORK

Next, we review classical RGB-D camera approaches as well as both pure and multi-sensor, event-based visual odometry solutions.

**RGB-D camera-based solutions:** The most straightforward solutions to RGB-D camera-based VO use only depth information [8], [11]. While they may potentially operate in dark environments, they require dense depth image processing at high frame rate, and therefore require high energy and computation resources (e.g. GPU). Approaches that also rely on images [3], [4], [5], [6] often perform dense photometric alignment, and thus still depend on exhaustive parallel computing. They furthermore have the disadvantage of degrading in challenging visual conditions (e.g. blur, low illumination). Most related to our method are approaches relying on sparsified, semi-dense depth maps [7], [12]. They have large convergence basins, stability under illumination changes, and high computational efficiency. Nonetheless, they still depend on intensity images for edge detection, and therefore continue to demonstrate high sensitivity to motion blur and low-illumination conditions.

**Pure event camera-based solutions:** Event cameras offer strong advantages such as high dynamic range, low latency, and low power consumption. However, the complicated nature of event data demands for novel theories and approaches, and full 6-DoF motion estimation with a single event camera remains a challenging problem. Many works rely on simplifying assumptions. Weikersdorfer et al. [13] proposed an event-based 2D SLAM framework for planar motions. Other works rely on a contrast maximization objective that utilizes image-to-image warping, a function that only works if the image transformation is at most a homography (e.g. pure rotation, planar homography) [14], [15], [16], [17], [18]. The first full 6-DoF solution is given by Kim et al. [19], who proposed a complex framework of three decoupled probabilistic filters estimating intensity, depth, and pose, respectively. A geometric solution is given by Rebecq et al. [20], which relies on their earlier ray-density based structure extraction method EMVS [21]. The success of these methods is however limited to small-scale environments and small, dedicated movements on mapping modules. Zhu et al. [22] finally present a promising learning-based approach, which however depends on vast amounts of training data, and provides no guarantees of optimality or generality. ESVO uses a stereo event camera [23], and we compare it against our approach.

**Hybrid event-supported solutions:** Owing to their difficult nature, event cameras are often combined with other sensors such as IMUs or regular cameras. Censi and Scaramuzza [24] present a VO framework that estimates relative poses by fusing events with absolute brightness information. Kueng et al. [25] detect features from grayscale images and track the features using the support of event data. Intensity-based methods do not take full advantage of event cameras and may fail due to motion blur in dark or varying illumination settings. While approaches that process images and events individually [26], [27] may continue to work if no regular image features are perceived, they still suffer from severe robustness issues if such conditions persist over extended time intervals. Another work that is closely related to ours is introduced by Weikersdorfer et al. [28], who extend their previous work [13] and include an RGB-D sensor. However, their method is based on an outdated, low-resolution event camera model and relies on a fully voxelized and thus limited-size environment. The accuracy of their probabilistic approach is highly depend on the frequency of depth updates, which limits the speed of the motion.

## III. DEPTH-EVENT VISUAL ODOMETRY

This section presents the details of our novel stereo depth-event odometry framework, which we denote *DEVO*. We start by seeing an overview of the entire pipeline followed by the details of the event data representation for efficient processing, the semi-dense depth map extraction module, and the final 6 Degree-of-Freedom (DoF) tracking module.

### A. Framework overview

A flowchart of our proposed method detailing all steps is illustrated in Figure 2. We start by generating time-surfaces maps which put our event sets into a suitable representation for efficient and accurate edge extraction and alignment. Details are introduced in Section III-B. The representation is used in both a tracking and a mapping module, which—in analogy to classical visual SLAM architectures—run in two parallel threads. The tracking thread processes only events and incrementally estimates the 6-DoF camera pose by efficient 3D-2D edge alignment. Details of the tracking thread are introduced in Section III-D. The local reference semi-dense depth map is updated at lower framerate inside the mapping thread. It proceeds by extracting the semi-dense edge map from the time-surface maps and assigning depth values from the depth camera readings. The local
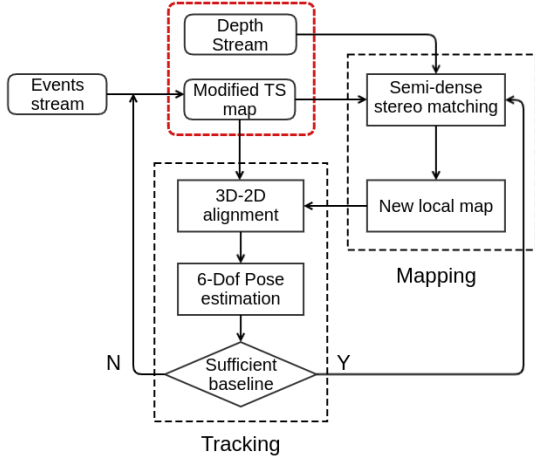
Fig. 2: Overview of our proposed *DEVO* visual localization and mapping pipeline for a stereo depth-event system.

map is updated whenever sufficient displacements between the current and the reference view has been detected. The operations of the map generation and the reference frame selection strategy are detailed in Section III-C.

### B. Event representation

Let us assume that we are given a set of $N$ events $\mathcal{E} = \{e_k\}_{k=1}^N$ occurring over a certain time interval. Each event $e_k = \{\mathbf{x}_k, t_k, b_k\}$ is defined by its image location $\mathbf{x}_k = [x_k \quad y_k]^T$, timestamp $t_k$, and polarity $b_k$. It is common to not process events asynchronously at the very high rate they occur, but aggregate sets of events accumulated during regularly spaced time intervals into one of three possible representations. The first one is given by space-time volumes of events [29], which are often used in conjunction with accurate but more computationally demanding continuous-time motion representations. The second one is given by simply ignoring the temporal nature of the events, and projecting all events along the temporal dimension onto a virtual binary image in which we then perform feature extraction. Though very efficient, the method re-induces motion blur and requires a careful selection of the time interval length. The third representation is given by time-surface maps (**TSM** [29]), which create an interesting balance between accuracy and efficiency. A **TSM** is an image in which a high pixel value denotes a recent event. The value at each pixel location $\mathbf{x}$ is a function of an exponential decay kernel and given by

$$\mathcal{T}(\mathbf{x}, t) = \exp(-\frac{t - t_{last}(\mathbf{x})}{\tau}), \qquad (1)$$

where $t$ is an arbitrary time, and $t_{last}(\mathbf{x}) \leq t$ is the timestamp of the last event triggered at $\mathbf{x}$. $\tau$ denotes the constant decay rate parameter, which requires careful tuning as a function of motion dynamics. The **TSM** visualizes the history of moving brightness patterns at each pixel location and emphasizes on locations in which motion has been more recent. The values in a **TSM** are mapped from $[0, 1]$ to $[0, 255]$ for convenient visualization and processing. We use a modified

**TSM** in which we only consider pixels with a value above a certain threshold $\delta$. Depending on the module (i.e. tracking, or mapping), other pixels are set to 0 or discarded.

### C. Mapping module

The mapping module performs semi-dense point cloud extraction. Let $\mathcal{T}_{\mathrm{ref}}(\cdot) = \mathcal{T}(\cdot, t_{\mathrm{ref}})$ be the **TSM** generated from the set of events $\mathcal{E}$ at time $t_{\mathrm{ref}}$. The semi-dense region $\mathcal{X}^{\mathrm{ref}}$ for which depth values will be extracted is simply given by all pixels for which the value is larger than $\delta$, i.e. $\mathcal{X}^{\mathrm{ref}} = \{\mathbf{x} \text{ s.t. } \mathcal{T}_{\mathrm{ref}}(\mathbf{x}) > \delta\}$. Based on the assumption that events are pre-dominantly triggered by high-gradient edges in the image, a proper choice of the decay rate $\tau$ and threshold $\delta$ will counteract motion blur and encourage the extracted semi-dense region to align tightly with effective appearance contours.

In order to retrieve the depth value for each point in the semi-dense region, we first warp the depth points from the depth camera at time $t_{\mathrm{ref}}$ to the event camera. The location in the event camera is given by

$$\mathbf{x}_k^e = \pi_e(\mathbf{T}_{ed} \cdot D(z_k^d) \cdot \pi_d^{-1}(\mathbf{x}_k^d)), \qquad (2)$$

where $\pi_{e/d}$ and $\pi_{e/d}^{-1}$ represent the known camera-to-image and the image-to-camera transformations of the event and the depth camera, respectively. They are defined as mapping from the 2D image space to 3D homogeneous space and vice-versa. $D(a) = \mathrm{diag}(a, a, a, 1)$ generates a diagonal matrix with elements $a$, $a$, $a$, and 1 along the diagonal. $\mathbf{T}_{ed}$ is the known $4 \times 4$ Euclidean extrinsic transformation matrix from the depth to the event camera. Finally, $\mathbf{x}_k^d$ and $z_k^d$ are a point and its corresponding depth in the depth camera, and $\mathbf{x}_k^e$ is the warped point in the depth frame. The depth $z_k^e$ at the latter point is easily obtained by

$$z_k^e = [0\ 0\ 1\ 0] \cdot \mathbf{T}_{ed} \cdot D(z_k^d) \cdot \pi_d^{-1}(\mathbf{x}_k^d), \qquad (3)$$

Note that the warping maps depth values onto sub-pixel locations rather than event camera pixel centers. It may furthermore induce occlusions or leave pixels with unobserved depths. In order to find a unique depth for each pixel, we create an individual list of nearby warped points from the depth image for each pixel in the semi-dense region. The value of the depth is conditionally set if the pixel is surrounded by warped points from the depth image. A simple depth clustering strategy identifies potential foreground points, and the final value is found by simple interpolation and ray intersection. This ensures that the depth of the pixels in the semi-dense region is always corresponding to foreground points and never affected by occlusions, depth measurement errors, or potential misalignments such as small errors in the extrinsic calibration parameters. The points that have a valid depth assigned to them are renormalized and multiplied by their depth, which finally results in the set $\mathcal{P}^{\mathrm{ref}}$ for our semi-dense 3D point cloud. Note that—in combination with the reference frame poses identified by the tracking module—multiple local maps could be merged into a global map using classical point cloud fusion techniques. The focus of the present work remains however on the localization problem.

## D. 6-Dof Camera tracking

With local 3D semi-dense point clouds from the mapping module in hand, we may now proceed to the details of our continuous, 6-DoF motion tracking module. We use the existing event-based localization strategy of Zhou et al. [23] in order to align subsequent **TSM**s with respect to the local semi-dense point cloud. As shown in Figure 3, the local map (i.e. the reference frame) is furthermore updated by the mapping thread each time the baseline with respect to the previous reference frame exceeds a given threshold. Theoretically, given that events are triggered asynchronously and at very high rate (with temporal resolution in the order of micro-seconds), we could update the pose of the camera with high frequency. Here we choose a rate of 100Hz, which already leads to a strong ability in handling highly dynamic motion.

The tracking proceeds by constructing a potential field in the current view. Based on a hypothesized pose, the reprojected point locations from the semi-dense point cloud then lead to a sampling of this field, and the sum of squares of the sampled values is considered as an energy to be minimized over the pose parameters of the camera. The potential field is constructed by negating and offsetting the **TSM** at the current time $t_{\text{cur}}$, i.e. $\overline{\mathcal{T}}_{\text{cur}}(\cdot) = 1 - \mathcal{T}(\cdot, t_{\text{cur}})$.

The detailed form of the objective to be minimized is as follows. The local semi-dense 3D point cloud at reference time $t_{\text{ref}}$ is still given by $\mathcal{P}^{\text{ref}}$. The absolute pose of the current view is given by

$$\mathbf{T}(\boldsymbol{\theta}) = \begin{bmatrix} \mathbf{R}(\mathbf{q}) & \mathbf{t} \\ \mathbf{0}^{\mathsf{T}} & 1 \end{bmatrix}. \tag{4}$$

where $\boldsymbol{\theta} = [\mathbf{t}^T \ \mathbf{q}^T]^T$ represents a motion parameter vector, $\mathbf{t}$ the position of the camera expressed in a world frame, and $\mathbf{q}$ its orientation as a Rodriguez vector. The relative transformation from the current camera's position to the nearest reference frame is then given as:

$$\mathbf{T}_{\text{rel}}(\boldsymbol{\theta}_{\text{rel}}) = \mathbf{T}_{\text{ref}}(\boldsymbol{\theta}_{\text{ref}})^{-1} \mathbf{T}_{\text{cur}}(\boldsymbol{\theta}_{\text{cur}}).$$

We also define the warping function $W$ that warps a 3D point from the local map to the current frame. It is given by

$$W(\mathbf{x}_k^{\text{ref}}; \boldsymbol{\theta}_{\text{rel}}) = \pi_e(\mathbf{T}^{-1}(\boldsymbol{\theta}_{\text{ref}}) \cdot D(z_k^e) \cdot \pi_e^{-1}(\mathbf{x}_k^e)). \tag{5}$$

The final goal of the tracking module is to find the optimum motion parameters $\boldsymbol{\theta}_{\text{rel}}$ that maximize the alignment of the reprojection of the local map $\mathcal{P}^{\text{ref}}$ and the local minima in our current negated **TSM** $\overline{\mathcal{T}}_{\text{cur}}(\cdot)$. Using the $W$, the objective function to find the optimum $\boldsymbol{\theta}_{\text{rel}}$ can be expressed as

$$\underset{\boldsymbol{\theta}_{\text{rel}}}{\arg\min} \sum_{\mathbf{x}_k^{\text{ref}} \in \mathcal{P}^{\text{ref}}} \rho(\overline{\mathcal{T}}_{\text{cur}}(W(\mathbf{x}_k^{\text{ref}}; \boldsymbol{\theta}_{\text{rel}}))^2), \tag{6}$$

where $\rho$ is a robust loss function. Similar to [23], (6) is reformulated by using a forward compositional Lucas-Kanade method [30], which refines the incremental motion parameters $\Delta\boldsymbol{\theta}_{\text{rel}}$ by minimizing:

$$\underset{\Delta\boldsymbol{\theta}_{\text{rel}}}{\arg\min} \sum_{\mathbf{x}_k^{\text{ref}} \in \mathcal{P}^{\text{ref}}} \rho(\overline{\mathcal{T}}_{\text{cur}}(W(W(\mathbf{x}_k^{\text{ref}}; \Delta\boldsymbol{\theta}_{\text{rel}}); \boldsymbol{\theta}_{\text{rel}})))^2), \tag{7}$$



Semi-dense point cloud

Depth map

$\mathbf{x}_k^e$

Time-surface map
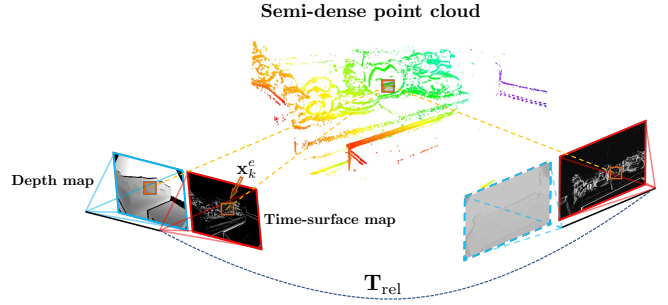
$\mathbf{T}_{\text{rel}}$

Fig. 3: *6-Dof Camera tracking*. Note that the depth image indicated in the dashed frame will not be used in tracking module.

and the new warping function $W(W(\mathbf{x}_k^{\text{ref}}; \Delta\boldsymbol{\theta}_{\text{rel}}); \boldsymbol{\theta}_{\text{rel}})$ is updated in each iteration. The new compositional approach is more efficient than the original method given that the Jacobian of the objective function remains constant at the position of zero increment and can be pre-computed. Smoothness, differentiability and convexity of this method are proven in [23].

## IV. Experimental Evaluation

We evaluate the performance of our novel visual odometry pipeline on both public and self-collected sequences. We start by introducing further details about the implementation and our hardware configuration. Next, we compare our methods against several alternatives on both mild test cases and more challenging scenarios. Alternatives are given by state-of-the-art event-based and RGB-D or depth-only based approaches. Both qualitative and quantitative results are provided, which demonstrate the effectiveness of our method. We conclude with an analysis of the computational performance of all above mentioned VO systems.

### A. Implementation details

Our first experiments are conducted on the Multi-Vehicle-Stereo-Event-Camera dataset (**MVSEC**) presented in [31]. These publicly available sequences include synchronized event streams, intensity images and depth images with ground truth trajectories. Next, in order to put a full stress test onto all methods, we test the methods on several other, self-collected sequences with different types of textures, motion characteristics, and illumination conditions. For different types of scene textures, the sequences are named *cali*, *table* and *sofa*, respectively. *cali* is a scene with many calibration



Fig. 4: Full sensor system, with event camera, industrial camera and RGB-D sensor.

TABLE I: Specifications of sensors used in our experiments.

| Sensor | Exposure Time | Resolution | Frame Rate |
|---|---|---|---|
| PointGrey-GS3 | 10ms | 1224×1024 | 30fps |
| Azure Kinect | 12.8ms | 640×576 | 30fps |
| Prophesee-Gen3 | - | 640×480 | - |

boards, *table* a standard desktop environment, and *sofa* a living room scene. For each texture, we capture datasets under three different motion speeds, denoted *fast*, *mid* and *slow*. More datasets are captured under a variety of illumination conditions, denoted *bright*, *darkish*, *dim*, *dark*, and *hdr*. All sequences are listed in Table IV. The sequences are collected by a custom-designed, hardware-synchronized multi-sensor system (cf. Figure 4), which contains a global-shutter industrial camera (PointGrey-GS3), a high-resolution event camera (Prophesee-Gen3), and an RGB-D sensor (Azure Kinect). Detailed specifications are listed in Table I. The multi-sensor system is intrinsically and extrinsically calibrated, and ground truth for all sequences is captured by a highly accurate external motion capture system.

### B. Comparison against event-based solutions

We first compare our proposed depth-event method **DEVO** against **ESVO**, an open-source event-based stereo visual odometry framework published in [23]. The two methods are evaluated on the public dataset **MVSEC** [31]. We choose both indoor and outdoor sequences, which are captured by a flying drone inside the room, and a stereo event camera mounted on a vehicle, respectively. Note that the depth measurements in **MVSEC** are obtained from a LiDAR, which can easily be considered as a replacement for the depth camera in our method.

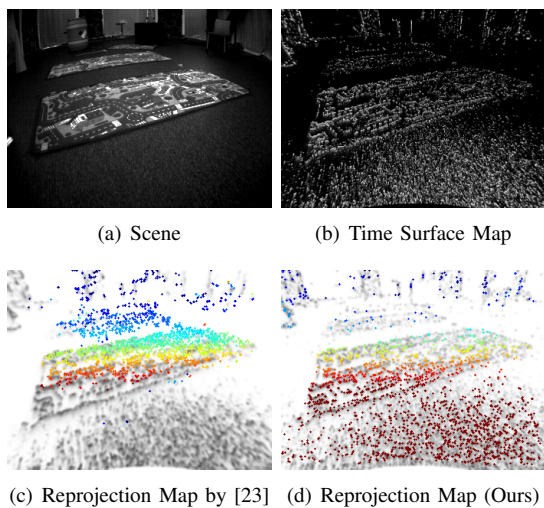Quantitative results are listed in Table II. As can be ob-



(a) Scene      (b) Time Surface Map

(c) Reprojection Map by [23]    (d) Reprojection Map (Ours)

Fig. 5: Visualization of an image (top left) and **TSM** (top right) from sequence *indoor flying2* and corresponding reprojections into a nearby frame for **ESVO** (bottom left) and **DEVO** (bottom right). The coloring indicates the depth of each point.

TABLE II: Comparison on **MVSEC** Datasets
[$\mathbf{R}_{\text{rpe}}$: °/s, $\mathbf{t}_{\text{rpe}}$: cm/s, $\mathbf{t}_{\text{ate}}$: cm]

| Sequence | DEVO | | | ESVO | | |
|---|---|---|---|---|---|---|
| | $\mathbf{R}_{\text{rpe}}$ | $\mathbf{t}_{\text{rpe}}$ | $\mathbf{t}_{\text{ate}}$ | $\mathbf{R}_{\text{rpe}}$ | $\mathbf{t}_{\text{rpe}}$ | $\mathbf{t}_{\text{ate}}$ |
| *upenn indoor flying1* | **0.30** | **0.88** | **20.58** | 0.37 | 1.63 | 21.68 |
| *upenn indoor flying2* | **0.36** | **1.12** | **11.33** | - | - | - |
| *upenn indoor flying3* | **0.53** | **1.21** | **10.60** | 0.54 | 2.14 | 25.40 |
| *upenn indoor flying4* | **0.53** | **1.44** | **13.16** | - | - | - |
| *upenn outdoor day1* | **0.30** | **7.77** | **88.70** | - | - | - |

served, **DEVO** clearly outperforms **ESVO** in all sequences. It should be noted that sequences *indoor flying2* and *indoor flying4* are much more challenging than the other two sequences owing to high noise in the event streams caused by a combination of difficult texture and highly dynamic platform motion. Examples of the sequence are indicated in Figure 5. The ground surface triggers a large number of noisy events for which depth is hard to observe. This severely influences the mapping result of **ESVO** and— due to the highly interleaved tracking and mapping modules—causes tracking failures in this fast exploration scenario. In contrast, **DEVO** directly reads the depth from the depth sensor, the quality of which is less influenced by the noisy nature of the texture. The independent depth readings significantly contribute to the robustness of the entire system when the inputs of the event camera degrade. Furthermore, the ability of any stereo method to perceive depths beyond a certain range is limited by the baseline of the system, which is why **ESVO** is unable to provide competitive results on the outdoor sequence.

### C. Comparison against RGB-D camera-based solutions

In order to analyze robustness under challenging illumination conditions, we compare our method against two classical approaches that rely on RGB-D cameras or depth sensors, only. They are given by **KinectFusion** [8] and **Canny-VO** [7]. We apply all methods to our self-collected datasets. We conduct three types of experiments, and all absolute trajectory errors (ATE) and relative pose errors (RPE) are summarized in Table IV:

- *Variation of light conditions*: We apply all methods on a series of sequences with different illumination conditions denoted *bright*, *darkish*, *dim*, *dark* and *high dynamic range (hdr)*. As summarized in Table III, both **DEVO** and **KinectFusion** are able to continuously track through all sequences, while **Canny-VO** proves to be

TABLE III: Comparison for different light conditions

| Sequence | DEVO | KinectFusion | Canny-VO |
|---|---|---|---|
| *light* | ✓ | ✓ | ✓ |
| *darkish* | ✓ | ✓ | ✗ |
| *dim* | ✓ | ✓ | ✗ |
| *dark* | ✓ | ✓ | ✗ |
| *HDR* | ✓ | ✓ | ✗ |

TABLE IV: Relative pose error and absolute trajectory error on self-collected datasets [$R_{rpe}$: °/s, $t_{rpe}$: cm/s, $t_{ate}$: cm]

| Sequence | DEVO | | | Canny-VO | | | KinectFusion | | |
|---|---|---|---|---|---|---|---|---|---|
| | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ |
| cali_bright_fast | **3.73** | 2.03 | 23.67 | 3.81 | **1.47** | 15.55 | 3.74 | 1.81 | **15.34** |
| cali_bright_mid | 1.44 | 1.77 | **16.90** | 1.42 | **1.26** | 21.23 | **1.35** | 1.77 | 19.22 |
| cali_bright_slow | **0.97** | 0.78 | 11.85 | 1.03 | **0.59** | **7.16** | 0.99 | 0.89 | 14.52 |
| cali_darkish_slow | 1.03 | **0.91** | 18.02 | - | - | - | **1.02** | 0.93 | **11.34** |
| cali_dim_slow | **1.55** | 0.88 | 35.38 | - | - | - | 1.62 | **0.82** | **9.05** |
| cali_dark_fast | **0.58** | **0.87** | 26.43 | - | - | - | 0.63 | 0.92 | **12.61** |
| cali_dark_mid | **0.49** | 0.60 | 17.65 | - | - | - | 0.54 | **0.59** | **12.89** |
| cali_dark_slow | **0.24** | 0.31 | 9.85 | - | - | - | 0.26 | **0.23** | **8.97** |
| cali_hdr_slow | **0.92** | 0.79 | 21.55 | - | - | - | 0.95 | **0.71** | **11.10** |
| table_bright_fast | 1.50 | 2.42 | 46.37 | 1.51 | **2.22** | 30.81 | **1.38** | 2.75 | **27.00** |
| table_bright_mid | 1.16 | 1.53 | 27.83 | 1.18 | **1.25** | 19.68 | **1.10** | 1.79 | 21.71 |
| table_bright_slow | 0.63 | 1.00 | **19.5** | 0.64 | **0.82** | 26.94 | **0.59** | 1.15 | 22.48 |
| sofa_bright_fast | **2.60** | 2.30 | 30.22 | 2.63 | **1.90** | **23.89** | 2.61 | 3.64 | 27.79 |
| sofa_bright_mid | 5.28 | 4.02 | **13.4** | **1.18** | **1.25** | 19.68 | 3.13 | 7.62 | 71.5 |
| sofa_bright_slow | 1.47 | 1.16 | **10.94** | **0.64** | **0.82** | 26.94 | 1.47 | 1.21 | 21.82 |

fragile when applied in poor illumination conditions. The reason is a lack of edge features caused by blur and poor contrast in dark scenarios.

- *Variation of motion characteristics*: We evaluate the performance of all methods for different motion dynamics. The sequences are denoted *fast*, *mid*, or *slow* to indicate the different camera dynamics. As can be observed in Table IV, all methods have a remarkable ability to handle dynamic scenarios for standard depth camera frame rate.
- *Variation of depth camera frame rate*: In order to analyse each method's ability to operate in an energy-saving mode, we finally test all methods for different depth camera frame rates between 30Hz and 1Hz in the *table* environment and for three different camera dynamics. As indicated in Table V, only our method is able to maintain stable tracking for all depth camera frame rates down to 1Hz. While accuracy decreases for more agile motion, it should be noted that the motion on these sequences is highly aggressive.

### D. Computational Performance

As can be observed from the ATE and RPE errors listed in Table IV, **DEVO** has comparable performance with other state-of-the-art methods from the literature. Although **Canny-VO** demonstrates lowest RPE errors in good lighting conditions, it shows degrading performance when the illumination becomes more challenging. We perceive **KinectFusion** as the strongest competitor of our method as it achieves comparable accuracy on all sequences. However, it should be noted that **KinectFusion** results have been obtained by putting the software into a high-performance setting that requires sufficient computing power to run. **KinectFusion** results in this paper have been obtained by a 32 core CPU and two Nvidia RTX 2080Ti. By comparison, the other two methods run on an 8 core CPU, only. It should furthermore be noted that **KinectFusion** depends on sufficiently high

TABLE V: Comparison for different depth frame rates [$R_{rpe}$: °/s, $t_{rpe}$: cm/s, $t_{ate}$: cm]

| Frequency | DEVO | | | Canny-VO | | | KinectFusion | | |
|---|---|---|---|---|---|---|---|---|---|
| Fast | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ |
| 30 | 1.50 | 2.42 | 46.37 | 1.51 | **2.22** | 30.81 | **1.37** | 2.75 | **27.00** |
| 15 | **2.92** | 4.86 | 46.70 | 2.96 | **4.46** | 29.11 | 3.04 | 8.36 | 37.99 |
| 10 | **4.26** | 7.32 | 47.92 | 4.65 | **6.55** | 34.68 | 4.83 | 17.33 | 66.78 |
| 5 | **7.73** | **14.73** | 57.89 | - | - | - | 9.01 | 26.86 | **55.09** |
| 1 | **18.58** | **49.16** | **76.04** | - | - | - | - | - | - |
| Medium | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ |
| 30 | 1.16 | 1.53 | 27.83 | 1.18 | **1.25** | 19.68 | **1.10** | 1.79 | 21.71 |
| 15 | 2.29 | 3.01 | 24.20 | 2.34 | **2.48** | 20.08 | **2.17** | 3.56 | 21.31 |
| 10 | 3.39 | 4.51 | 21.46 | 3.49 | **3.70** | 20.34 | **3.20** | 5.58 | 58.73 |
| 5 | **6.55** | **9.20** | **21.55** | - | - | - | 7.16 | 15.80 | 37.07 |
| 1 | **18.46** | **35.24** | **51.93** | - | - | - | - | - | - |
| Slow | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ | $R_{rpe}$ | $t_{rpe}$ | $t_{ate}$ |
| 30 | 0.63 | 1.00 | **19.50** | 0.64 | **0.82** | 26.94 | **0.59** | 1.15 | 22.48 |
| 15 | 1.21 | 1.97 | **18.44** | 1.24 | **1.65** | 26.01 | **1.13** | 2.29 | 22.22 |
| 10 | 1.76 | 2.98 | **18.46** | 1.82 | **2.46** | 26.32 | **1.66** | 3.41 | 22.10 |
| 5 | **3.28** | 6.09 | **17.61** | 3.71 | **5.54** | 27.97 | - | - | - |
| 1 | **10.66** | **31.15** | **37.88** | - | - | - | - | - | - |

depth camera framerates, which again induces larger energy consumption.

## V. CONCLUSION

We present a novel approach to visual odometry that relies on a stereo depth-event camera. In comparison to depth-only alternatives, it handles faster motion and works more efficiently by requiring lower depth image frame rates and by performing semi-dense image processing. In comparison to RGB-D visual odometry solutions, it successfully handles challenging or low illumination scenarios. In summary, our proposed method handles a large spectrum of challenging situations, and we believe that it could represent a highly interesting approach for intelligent mobile systems that require indoor localization. A pre-condition would however be that event cameras are becoming more affordable in the future.

## REFERENCES

[1] J. Fuentes-Pacheco, J. Ruiz-Ascencio, and J. M. Rendón-Mancha, "Visual simultaneous localization and mapping: a survey," *Artificial intelligence review*, vol. 43, no. 1, pp. 55–81, 2015.

[2] C. Cadena, L. Carlone, H. Carrillo, Y. Latif, D. Scaramuzza, J. Neira, I. Reid, and J. J. Leonard, "Past, present, and future of simultaneous localization and mapping: Toward the robust-perception age," *IEEE Transactions on robotics*, vol. 32, no. 6, pp. 1309–1332, 2016.

[3] F. Steinbrücker, J. Sturm, and D. Cremers, "Real-time visual odometry from dense RGB-D images," in *IEEE International Conference on Computer Vision Workshops (ICCV Workshops)*, 2011.

[4] F. Endres, J. Hess, J. Sturm, D. Cremers, and W. Burgard, "3-d mapping with an rgb-d camera," *IEEE transactions on robotics*, vol. 30, no. 1, pp. 177–187, 2013.

[5] C. Kerl, J. Sturm, and D. Cremers, "Robust odometry estimation for RGB-D cameras," in *Proceedings of the IEEE International Conference on Robotics and Automation (ICRA)*, 2013.

[6] P. Henry, M. Krainin, E. Herbst, X. Ren, and D. Fox, "Rgb-d mapping: Using depth cameras for dense 3d modeling of indoor environments," in *Experimental robotics*. Springer, 2014, pp. 477–491.

[7] Y. Zhou, H. Li, and L. Kneip, "Canny-vo: Visual odometry with rgb-d cameras based on geometric 3-d–2-d edge alignment," *IEEE Transactions on Robotics*, vol. 35, no. 1, pp. 184–199, 2018.

[8] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon, "Kinectfusion: Real-time dense surface mapping and tracking," in *Proceedings of the International Symposium on Mixed and Augmented Reality (ISMAR)*, 2011.

[9] C. Brandli, R. Berner, M. Yang, S.-C. Liu, and T. Delbruck, "A 240× 180 130 db 3 $\mu$s latency global shutter spatiotemporal vision sensor," *IEEE Journal of Solid-State Circuits*, vol. 49, no. 10, pp. 2333–2341, 2014.

[10] G. Gallego, T. Delbruck, G. Orchard, C. Bartolozzi, B. Taba, A. Censi, S. Leutenegger, A. Davison, J. Conradt, and K. Daniilidis, "Event-based vision: A survey," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2020.

[11] T. Whelan, R. F. Salas-Moreno, B. Glocker, A. J. Davison, and S. Leutenegger, "Elasticfusion: Real-time dense slam and light source estimation," *The International Journal of Robotics Research*, vol. 35, no. 14, pp. 1697–1716, 2016.

[12] F. Schenk and F. Fraundorfer, "Reslam: A real-time robust edge-based slam system," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 154–160.

[13] D. Weikersdorfer, R. Hoffmann, and J. Conradt, "Simultaneous localization and mapping for event-based vision systems," in *International Conference on Computer Vision Systems*. Springer, 2013, pp. 133–142.

[14] G. Gallego and D. Scaramuzza, "Accurate angular velocity estimation with an event camera," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 632–639, 2017.

[15] T. Stoffregen and L. Kleeman, "Event cameras, contrast maximization and reward functions: an analysis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 300–12 308.

[16] G. Gallego, M. Gehrig, and D. Scaramuzza, "Focus is all you need: loss functions for event-based vision," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 280–12 289.

[17] D. Liu, A. Parra, and T.-J. Chin, "Globally optimal contrast maximisation for event-based motion estimation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 6349–6358.

[18] X. Peng, L. Gao, Y. Wang, and L. Kneip, "Globally-optimal contrast maximisation for event cameras," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.

[19] H. Kim, S. Leutenegger, and A. J. Davison, "Real-time 3d reconstruction and 6-dof tracking with an event camera," in *European Conference on Computer Vision*. Springer, 2016, pp. 349–364.

[20] H. Rebecq, T. Horstschäfer, G. Gallego, and D. Scaramuzza, "Evo: A geometric approach to event-based 6-dof parallel tracking and mapping in real time," *IEEE Robotics and Automation Letters*, vol. 2, no. 2, pp. 593–600, 2016.

[21] H. Rebecq, G. Gallego, E. Mueggler, and D. Scaramuzza, "Emvs: Event-based multi-view stereo—3d reconstruction with an event cam-

era in real-time," *International Journal of Computer Vision*, vol. 126, no. 12, pp. 1394–1414, 2018.

[22] D. Zhu, Z. Xu, J. Dong, C. Ye, Y. Hu, H. Su, Z. Liu, and G. Chen, "Neuromorphic visual odometry system for intelligent vehicle application with bio-inspired vision sensor," in *2019 IEEE International Conference on Robotics and Biomimetics (ROBIO)*. IEEE, 2019, pp. 2225–2232.

[23] Y. Zhou, G. Gallego, and S. Shen, "Event-based stereo visual odometry," *IEEE Transactions on Robotics*, 2021.

[24] A. Censi and D. Scaramuzza, "Low-latency event-based visual odometry," in *2014 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2014, pp. 703–710.

[25] B. Kueng, E. Mueggler, G. Gallego, and D. Scaramuzza, "Low-latency visual odometry using event-based feature tracks," in *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, 2016, pp. 16–23.

[26] H. Rebecq, T. Horstschaefer, and D. Scaramuzza, "Real-time visual-inertial odometry for event cameras using keyframe-based nonlinear optimization," 2017.

[27] A. R. Vidal, H. Rebecq, T. Horstschaefer, and D. Scaramuzza, "Ultimate slam? combining events, images, and imu for robust visual slam in hdr and high-speed scenarios," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 994–1001, 2018.

[28] D. Weikersdorfer, D. B. Adrian, D. Cremers, and J. Conradt, "Event-based 3d slam with a depth-augmented dynamic vision sensor," in *2014 IEEE international conference on robotics and automation (ICRA)*. IEEE, 2014, pp. 359–364.

[29] X. Lagorce, G. Orchard, F. Galluppi, B. E. Shi, and R. B. Benosman, "Hots: a hierarchy of event-based time-surfaces for pattern recognition," *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 7, pp. 1346–1359, 2016.

[30] S. Baker and I. Matthews, "Lucas-kanade 20 years on: A unifying framework," *International journal of computer vision*, vol. 56, no. 3, pp. 221–255, 2004.

[31] A. Z. Zhu, D. Thakur, T. Ozaslan, B. Pfrommer, V. Kumar, and K. Daniilidis, "The multi vehicle stereo event camera dataset: An event camera dataset for 3d perception," *IEEE Robotics and Automation Letters*, pp. 2032–2039, 2018.