

Master Thesis

Use of a DVS for High Speed Applications

Autumn Term 2018

Declaration of Originality

I hereby declare that the written work I have submitted entitled

Use of a DVS for High Speed Applications

is original work which I alone have authored and which is written in my own words.¹

Author

Weizhen

Huang

Student supervisors

Florian

Tschopp

Ignacio

Alzugaray

Supervising lecturer

Roland

Siegwart

With the signature I declare that I have been informed regarding normal academic citation rules and that I have read and understood the information on 'Citation etiquette' (<https://www.ethz.ch/content/dam/ethz/main/education/rechtliches-abschluesse/leistungskontrollen/plagiarism-citationetiquette.pdf>). The citation conventions usual to the discipline in question here have been respected.

The above written work may be tested electronically for plagiarism.

Place and date

Signature

¹Co-authored work: The signatures of all authors are required. Each signature attests to the originality of the entire piece of written work in its final form.

Contents

Abstract	iii
Symbols	v
1 Introduction	1
2 Per-Frame Optimization	3
2.1 From Events to Frame	3
2.2 Measuring the Sharpness of an Image	5
3 6-DoF Pose Tracking in Planar Scenes	6
3.1 Camera Projection	6
3.2 Planar Homography	7
3.3 From Frames to Map	7
3.3.1 Map	7
3.3.2 Tracking	7
3.3.3 Mapping	10
3.3.4 Relocalization	10
3.4 Experiments	11
4 6DoF Motion Estimation in General 3D Scenes	18
5 Conclusion and Outlook	22
Bibliography	24
A Derivative of the Contrast Metric	25

Abstract

Event cameras are bio-inspired cameras that capture pixel-level intensity changes at a rate of up to 1 MHz. They naturally respond to edges in the scene and are free from motion blur.

We propose a method which performs SLAM in planar scenes and 6-DoF motion estimation in unstructured scenes with a monocular event camera, no additional sensing. This method is based on a contrast maximization pipeline. Events are grouped by a small temporal window and warped with the optimal motion parameters that produce the sharpest image. When performing SLAM, all the images are tracked using image-to-model alignment, but only a selective set of keyframes are used to construct the map. This method performs favorably against the ground truth, and is able to track high speed motions.

Symbols

Symbols

\mathbf{I}_n	$n \times n$ identity matrix
\mathbf{K}	camera projection matrix
\mathbf{A}^\top	matrix transpose
\mathbf{A}^{-1}	matrix inverse
$\mathbf{0}_{m \times n}$	an $m \times n$ -vector whose entries are 0
\mathbb{R}^n	space of real n -vectors, n -dimensional Euclidean space
\mathbb{S}^n	n -sphere
$SO(3)$	3D special orthogonal group, 3D rotation group
$\mathfrak{so}(3)$	the vector space of 3×3 skew-symmetric matrices
$SE(3)$	3D special Euclidean group, the group of rigid body motion
e	event
p	event polarity
\mathcal{E}	a set of events
δ	Dirac delta
\mathbf{x}	2D pixel coordinate
$\bar{\mathbf{x}}$	homogeneous coordinate
$\tilde{\mathbf{x}}$	calibrated coordinate
\mathbf{X}	3D point in space
\mathbf{R}	orientation, rotation
\mathbf{t}	translation
\mathbf{T}	camera pose, including \mathbf{R} and \mathbf{t}
$\boldsymbol{\omega}$	angular velocity
\mathbf{v}	linear velocity
\mathbf{n}	normal direction of a plane
$\boldsymbol{\theta}$	parameter set, might include $\mathbf{R}, \mathbf{t}, \boldsymbol{\omega}, \mathbf{v}$ and \mathbf{n}
\wedge	hat operator
\sim	equality up to a non-zero scale

Acronyms and Abbreviations

DVS	Dynamic Vision Sensor
APS	Active Pixel Sensor
DoF	Degrees of Freedom

SLAM	Simultaneously Localizing and Mapping
BFGS	Broyden-Fletcher-Goldfarb-Shanno Algorithm

Chapter 1

Introduction

Frame-based cameras that widely used in computer vision output images at a pre-set rate, even when the intensity values stay unchanged, resulting in redundant data. When the scene moves too fast, an insufficient frame rate would cause motion blur. To address this problem, Lichtsteiner et al. [1] presents an event-based camera called DVS (dynamic vision sensor) that reports pixel-level log intensity changes at a rate of MHz scale. Unlike traditional cameras that outputs a frame until all the pixels are scanned, an event-based camera outputs are asynchronous: when the intensity change at a pixel reaches a threshold $|\ln I(t) - \ln I(t - \Delta t)| > C$, the sensor outputs an “event” $e = \{x, y, t, p\}$, which includes the pixel coordinate (x, y) , the timestamp t and the polarity $p = \pm 1$ indicating positive or negative intensity change. The output is thus a stream of events instead of frames. Later Brandli et al. [2] presents DAVIS (dynamic and active pixel vision sensor) that has additional APS (active pixel sensor) circuits that provide absolute intensity information, while offering DVS outputs with higher resolution (240×180 against 128×128), higher dynamic range (130 dB against 120 dB) and lower latency ($3 \mu s$ against $15 \mu s$). Newer sensors also provide color information[3, 4] or have higher image resolution[5].

Since for an event-based camera, there is no such thing as frames, feature detection and tracking algorithms that work well for standard cameras cannot be directly applied to event-based cameras. There are several works that try to adopt feature detection and tracking pipelines. Zhu et al. [6] aggregate a small number of events to construct a frame, and perform Harris corner detector [7] on the synthesized frames, then track the features with implicit data association. The work of Tedaldi et al. [8] detects features on the APS outputs, and performs tracking with the DVS outputs. Instead of working with frames, Mueggler et al. [9] detect corners directly in the spatiotemporal event stream. There are also works that perform 3D reconstruction with known poses [10] or 6-DoF tracking with a known map [11], or both with the help of IMUS which is integrated on DAVIS [12]. Besides, [13, 14] perform 6-DoF tracking and 3D reconstruction purely based on event streams, with methods commonly applied in computer vision, for example DSI (disparity space image) or EKF (extended Kalman Filter). Newer works also combine machine learning and event-based camera [15, 16, 17]

Gallego and Scaramuzza [18] first propose an interesting contrast maximization framework that is rather specific for event-based cameras. Without the help of any auxiliary variable such as optical flow or feature, this framework finds the optimal angular velocity that maximizes the contrast of a synthesized event frame via nonlinear optimization. Later in [19] they showed that the same framework can be applied to various important tasks in computer vision, such as motion, depth, and optical flow estimation.

This work is an extension of [18, 19]. In these two works they showed how to estimate angular velocity in general scenes where only rotational motion is present, and 6-DoF motion in planar scenes with contrast maximization framework. In this work we will see that the same idea can be applied to perform SLAM (simultaneously localizing and mapping) in planar scenes, and motion estimation in general scenes with general motion.

In this work we only used the DVS output of DAVIS[2]. The algorithms are tested on the event-camera dataset recorded by Mueggler et al. [20]. The outline of this work is as follows:

Chapter 2

Per-Frame Optimization

The output of a DVS is a stream of events $\{e_k\}$, with $e_k = \{x_k, y_k, t_k, p_k\}$ denoting its pixel coordinate, timestamp and polarity, respectively. Figure 2.1 shows the DVS output interleaved with the APS output. We can clearly see that the events are triggered at the shape boundaries: when the sensor (or the scene) is moving, these pixel locations will receive intensity changes. Thus DVS works naturally as an edge detector.

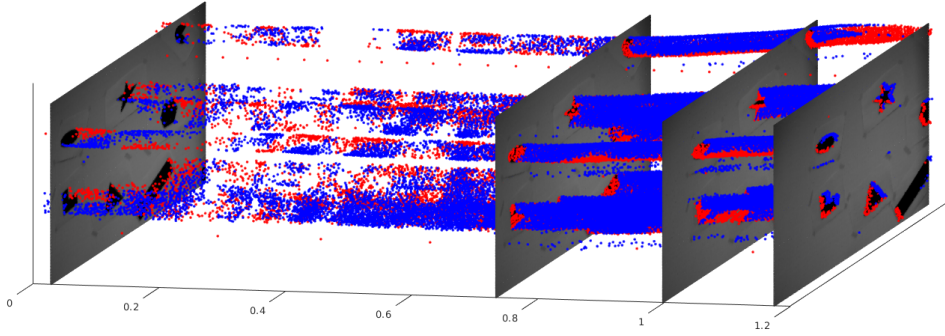


Figure 2.1: Event stream plus grayscale images. During this period 100 000 events and 27 frames are produced. Red and blue points denote the positive and negative events, respectively. DVS has a much higher output rate and is free from motion blur.

Note that events can be triggered either by scene or camera motion. In this work, we assume a static scene and a moving camera.

2.1 From Events to Frame

A single event gives little information. A common strategy would be to aggregate all the events within a small temporal window. Depending on the size of the window we may get different results as in fig. 2.2.

More specifically, we group a set of events $\mathcal{E} \doteq \{e_k\}_{k=1}^N$ happened within a certain period, and synthesize an image by summing up the events. If we simply sum along the time axis, the intensity at each pixel will be the sum of the polarities of all the events that are triggered at this pixel location during this period:

$$\mathcal{I}(\mathbf{x}) = \sum_{k=1}^N \pm_k \delta(\mathbf{x} - \mathbf{x}_k), \quad (2.1)$$

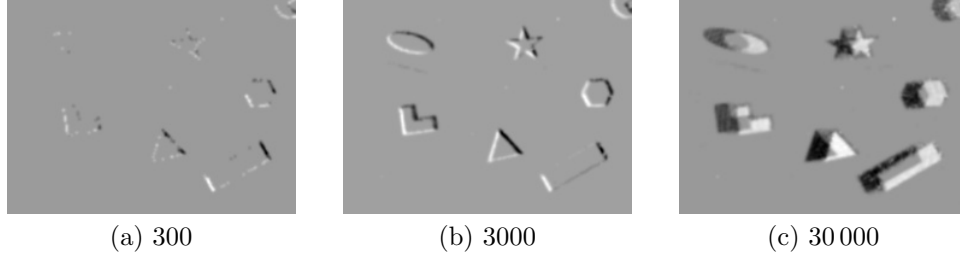


Figure 2.2: The image generated by aggregating 300, 3000 and 30 000 events. White and black dots represent positive and negative events, respectively. Gray indicates that no event is present in this pixel, or positive and negative events cancel each other out.

with \pm_k and $\mathbf{x}_k \in \mathbb{R}^2$ denoting the polarity and pixel coordinates of the k th event, respectively. We call such synthesized results in fig. 2.2 as *frames*.

The temporal resolution of a DVS in [2] is 1 μ s; each event has a distinct timestamp, which also indicates that events happen at distinct camera poses. As a result, a frame as in fig. 2.2 does not represent the scene at a given time. Instead of simply summing up events from different timestamps, a better strategy is to project the events along the motion path, which is clearly visible in fig. 2.1, thus faithfully recover the scene edges.

But how do we compute the motion path? Without additional motion capture devices, we have no access to the ground truth pose; it's also impossible to know the camera pose associated with each event, which happens at such a high rate. To simplify the problem, we assume that the camera moves at a constant velocity within an interval, which is not a bad assumption if this interval is small enough. This group of events is warped by the model $\mathbf{x}'_k = \mathbf{W}(\mathbf{x}_k, t_k; \boldsymbol{\theta})$, after warping the \mathbf{x}_k in the intensity computation eq. (2.1) is substitute by \mathbf{x}'_k . We find the optimal parameters $\boldsymbol{\theta}$ that best describes the events within this window, then shift the window to the next set of events, and repeat this process.

There are several strategies available as to choosing the window size. We may group events by

1. a fixed time interval Δt ,
2. a fixed number of events N per subset,
3. dynamically adjusting the window size according to the current velocity.

The first strategy is used in [16], they also give the comparison on different choices of window size. This might be a good strategy when the motion estimation is needed at a fixed rate. However, it suffers from the same problem as the traditional camera, that is, it keeps producing redundant frames and motion estimation when the camera stops moving. Also, a fixed time interval is not able to well compromise between slow motion and rapid motion. Event cameras are data-driven cameras, the output depends on the motion and the scene complexity, thus the other two strategies are more favorable. Zhu et al. [6] use the *lifetime*[21] of the event to define the window size, which depends on the event's velocity on the image plane. This belongs to the third strategy, and is applicable independent of the scene complexity. Considering our test datasets, each dataset observes a same scene, the scene complexity thus stays almost the same within a certain sequence. For simplicity, we opt for the second strategy, and manually choose the window size for each sequence.

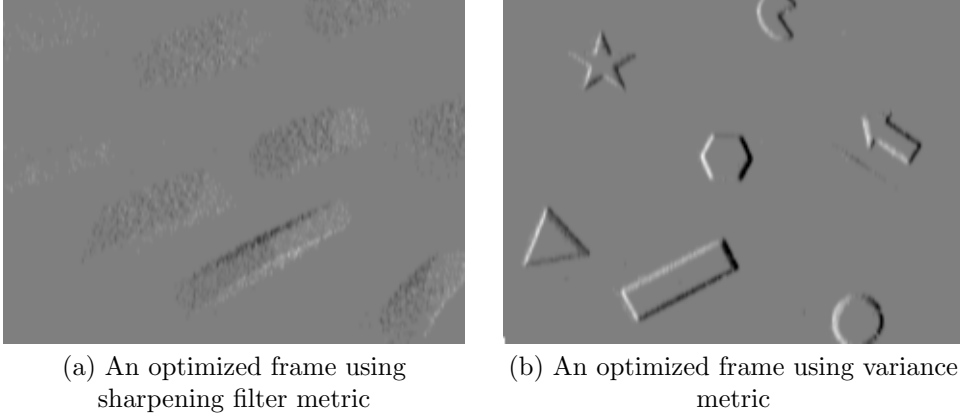


Figure 2.3: Comparison of two different metrics

2.2 Measuring the Sharpness of an Image

Intuitively, the image would appear sharp and have high contrast if the events are warped along its trajectory, otherwise motion-blurred. So we may use the sharpness of the warped image to measure how well a warp \mathbf{W} is.

There are several metrics one could choose from. A local contrast metric could be, for example, to convolute the image with a high pass filter

$$\mathcal{C}_H = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}, \quad (2.2)$$

and sum the pixel value of the filtered image. This metric would favor that each pixel stands out from its neighbors. However, every pixel is only compared with its 8 neighbors, thus an image with scattered events (large noise) is also a valid configuration, as shown in fig. 2.3(a). On the other hand, the Michelson contrast [22], defined as $\mathcal{C}_M = (\mathcal{I}_{\max} - \mathcal{I}_{\min}) / (\mathcal{I}_{\max} + \mathcal{I}_{\min})$, only considers the highest and the lowest luminance in the image and is thus more suitable to quantify contrast for periodic functions.

We choose to measure the contrast by the variance of the image, defined by

$$\text{Var}(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta})) \doteq \frac{1}{|\Omega|} \int_{\Omega} (\mathcal{I}(\mathbf{x}; \boldsymbol{\theta}) - \mu(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta})))^2 d\mathbf{x}, \quad (2.3)$$

where $\mu(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta}))$ is the average image intensity. If the scene has more or less balanced positive and negative events, μ should be close to zero. Actually, since we only want to distinguish the edges with the larger blank background, which indeed has zero intensity, we can also simply substitute $\mu(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta}))$ with 0.

The goal of measuring the contrast of an image is to find the configuration that aligns events triggered from the same visual stimuli; thus the variance eq. (2.3) is a very suitable metric, as a squared metric favors the configuration that projects as many events as possible to the same pixel.

Chapter 3

6-DoF Pose Tracking in Planar Scenes

The warping $\mathbf{x}' = \mathbf{W}(\mathbf{x}, t; \boldsymbol{\theta})$ appears to be done on the 2D image plane, although the 6-DoF rigid body motion is happened in 3D space. In this chapter, we explain the warping process, and how to perform SLAM with the warped events in planar scenes.

3.1 Camera Projection

Throughout this work the following notation is employed: w denotes the world frame, c_1 or c_2 denotes a camera frame. \mathbf{T}_{ab} is the transformation from frame a to frame b , measured in frame a . \mathbf{X} denotes the 3D coordinate of a point, with \mathbf{X}_w and \mathbf{X}_c the coordinate with respect to the world frame and camera frame, respectively. \mathbf{x} denotes a 2D coordinate on the image plane.

Assume the lens distortion has been removed, a point \mathbf{X} is projected on the image plane via

$$\bar{\mathbf{x}} = \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \frac{1}{Z} \begin{bmatrix} f_x & c_x \\ f_y & c_y \\ 1 \end{bmatrix} \begin{bmatrix} X \\ Y \\ Z \end{bmatrix} = \frac{1}{Z} \mathbf{K} \mathbf{X}, \quad (3.1)$$

with \mathbf{K} denoting the camera *intrinsic matrix*, f_x, f_y being its focal lengths and (c_x, c_y) the principal point and $\bar{\mathbf{x}}$ is the homogeneous coordinate of \mathbf{x} . We write $\tilde{\mathbf{x}} = \mathbf{K}^{-1} \bar{\mathbf{x}}$ and call it the *calibrated coordinates*.

Consider warping the events within the timespan $[t_0, t_N]$, and e_k is an event happened in this timespan with $0 < k < N$. Assume the camera pose at t_k relative to that at t_0 is $\mathbf{T}_{0k} = (\mathbf{R}_{0k}, \mathbf{t}_{0k}) \in SE(3)$, then the 3D point associated with e_k is projected on the image plane at t_0 via

$$\mathbf{X}_k = Z_k \tilde{\mathbf{x}}_k \quad (3.2)$$

$$\mathbf{X}_0 = \mathbf{R}_{0k} \mathbf{X}_k + \mathbf{t}_{0k} \quad (3.3)$$

$$\tilde{\mathbf{x}}_0 = \mathbf{X}_0 / Z_0 \quad (3.4)$$

$$\Rightarrow \tilde{\mathbf{x}}_0 \sim \mathbf{R}_{0k} Z_k \tilde{\mathbf{x}}_k + \mathbf{t}_{0k}, \quad (3.5)$$

where \sim means equality up to a non-zero scale factor. Equation (3.5) indicates that a warp \mathbf{W} does not only depend on the motion parameters (\mathbf{R}, \mathbf{t}) , but also the unknown depth of an event. It's intractable to also add the depths of every event to the parameter set $\boldsymbol{\theta}$; we need further simplifications.

3.2 Planar Homography

In the case of a planar scene the warping eq. (3.5) simplifies, since a plane \mathbf{P} can be parameterized by only two sets of parameters (\mathbf{n}, d) , with $\mathbf{n} \in \mathbb{S}^2$ being the unit surface normal of \mathbf{P} with respect to the camera frame, and d the distance from the camera center to \mathbf{P} . A point \mathbf{X} on the plane \mathbf{P} is described by

$$\mathbf{n}^\top \mathbf{X} - d = 0. \quad (3.6)$$

Substitute eq. (3.6) in eq. (3.3) we get

$$\begin{aligned} \mathbf{X}_0 &= \mathbf{R}_{0k} \mathbf{X}_k + \mathbf{t}_{0k} \\ \mathbf{X}_k &= \mathbf{R}_{0k}^\top (\mathbf{X}_0 - \mathbf{t}_{0k}) \end{aligned} \quad (3.7)$$

$$\mathbf{X}_k = \mathbf{R}_{0k}^\top (\mathbf{I}_3 + \mathbf{t}_{0k} \mathbf{n}^\top / d) \mathbf{X}_0. \quad (3.8)$$

We denote

$$\mathbf{H} = \mathbf{R}_{0k}^\top (\mathbf{I}_3 + \mathbf{t}_{0k} \mathbf{n}^\top / d) \quad (3.9)$$

as the planar homography, thus the warping function simplifies to $\tilde{\mathbf{x}}_0 \sim \mathbf{H}^{-1} \tilde{\mathbf{x}}_k$. Under a constant velocity model with linear velocity $\mathbf{v} \in \mathbb{R}^3$ and angular velocity $\boldsymbol{\omega} \in \mathbb{R}^3$, the translation is given by

$$\mathbf{t}_{0k} = \mathbf{v} t_k, \quad (3.10)$$

the rotation matrix is given by the *exponential map* $\exp: \mathfrak{so}(3) \rightarrow SO(3)$:

$$\mathbf{R}_{0k} = \exp(\boldsymbol{\omega}^\wedge t_k), \quad (3.11)$$

where $^\wedge$ is the *hat* operator

$$\boldsymbol{\omega}^\wedge = \begin{bmatrix} \omega_1 \\ \omega_2 \\ \omega_3 \end{bmatrix} = \begin{bmatrix} 0 & -\omega_3 & \omega_2 \\ \omega_3 & 0 & -\omega_1 \\ -\omega_2 & \omega_1 & 0 \end{bmatrix} \in \mathfrak{so}(3). \quad (3.12)$$

3.3 From Frames to Map

The contrast maximization procedure in chapter 2 optimizes the relative pose between successive frames. However, the same idea can be applied to perform global pose tracking in a planar scene by building a *map*. We first explain how the map is defined, and how to track with a known map, then we show how this map is built by selecting a set of keyframes.

3.3.1 Map

A map is a plane associated with a texture; it consists of three components: the normal direction \mathbf{n}_w , the distance d_w to the origin, and the texture, which represents all the edges on the plane. Figure 3.1 shows an example of such a map. Figure 3.1(a) also shows the set of keyframes used to construct the map. We will talk more about keyframes in section 3.3.3. The global coordinate is chosen as the camera coordinate of the first frame.

3.3.2 Tracking

Suppose a map is present, then the normal direction \mathbf{n}_w of the plane and the distance d_w to the origin are known. Also the pose of the current frame $(\mathbf{R}_{wc}, \mathbf{t}_{wc}) \in SE(3)$ is determined by the motion estimation from the last frame (a quick note to the

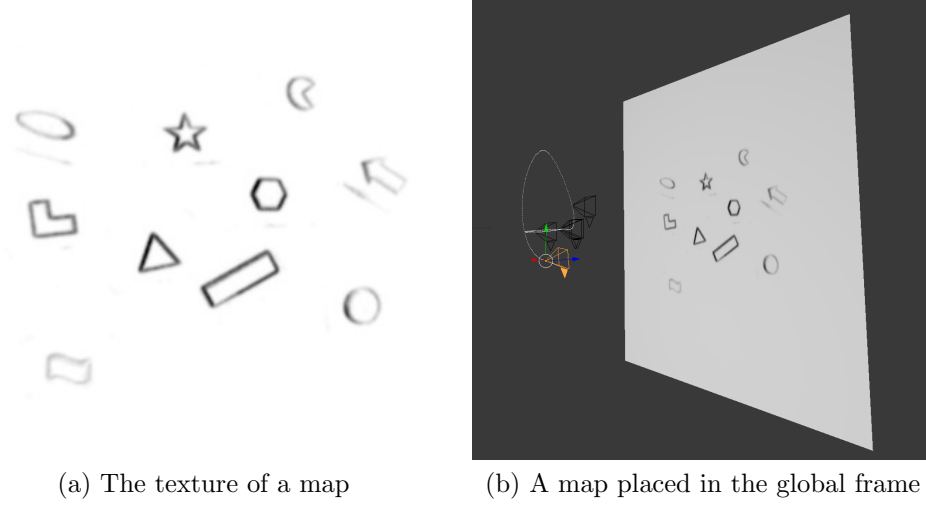


Figure 3.1: Map

terminology we are using: whenever we say the *pose* of a frame, we always refer to the camera pose at which the first event within the frame happens). The parameters left to be estimated for each frame is $\theta = (\omega, v) \in \mathbb{R}^6$. By substituting \mathbf{n} with $\mathbf{n}_c = \mathbf{R}_{cw}\mathbf{n}_w$, and d with $d_c = d_w + \mathbf{t}_{wc} \cdot \mathbf{n}_w$ in eq. (3.9), we get the homography matrix within each frame as

$$\mathbf{H}_1 = \mathbf{R}^\top (\mathbf{I} + \mathbf{t}\mathbf{n}_c^\top/d_c). \quad (3.13)$$

When estimating the motion of the first frame, we use an initial guess of $\omega = \mathbf{0}_{3 \times 1}$, $v = \mathbf{0}_{3 \times 1}$, $\mathbf{n} = (0, 0, -1)$, if no better prior is available. *Zero velocity* should be a good assumption if the camera just starts moving. In the subsequent frames, we use the output from the last frame as the initial guess, assuming that in the short period since the last optimization, the velocity did not change severely.

A nonlinear optimizing problem naturally suffers from the local optima problem. Without a good initialization, the per-frame optimization might come up with a parameter set delivering an image that appears sharp, despite being wrongly estimated (see fig. 3.2). We should also not simply integrate the velocities to obtain the global pose, which would result in drift over time. In order to make sure that the estimated pose from the per-frame contrast maximization also conforms to the global map, we perform another optimization by projecting the events of the current frame to the global map and measure the contrast of the composed image. The parameter set is still (ω^\top, v^\top) , and we use the output from the last procedure (the per-frame optimization) as the initial guess.

The projection from an event to the map is

$$\mathbf{x}_w \sim \mathbf{R}_n \mathbf{H}_2^{-1} \mathbf{H}_1^{-1} \mathbf{x}_c, \quad (3.14)$$

with \mathbf{H}_1 the planar homography as in eq. (3.13), but \mathbf{R} and \mathbf{t} are not necessarily the same as before, since we are refining these parameters, and

$$\mathbf{H}_2 = \mathbf{R}_{cw} (\mathbf{I} + \mathbf{T}_{wc}\mathbf{n}^\top/d_w) \quad (3.15)$$

the projection from the current frame to the global frame. \mathbf{R}_n is the transformation from the orientation of the global frame to the orientation of the map, computed by

$$\mathbf{K} = (\mathbf{n}_w \times \mathbf{z})^\wedge \quad (3.16)$$

$$\mathbf{R}_n = \mathbf{I} + \mathbf{K} + \mathbf{K}^2/(1 + \mathbf{n}_w \cdot \mathbf{z}), \quad (3.17)$$

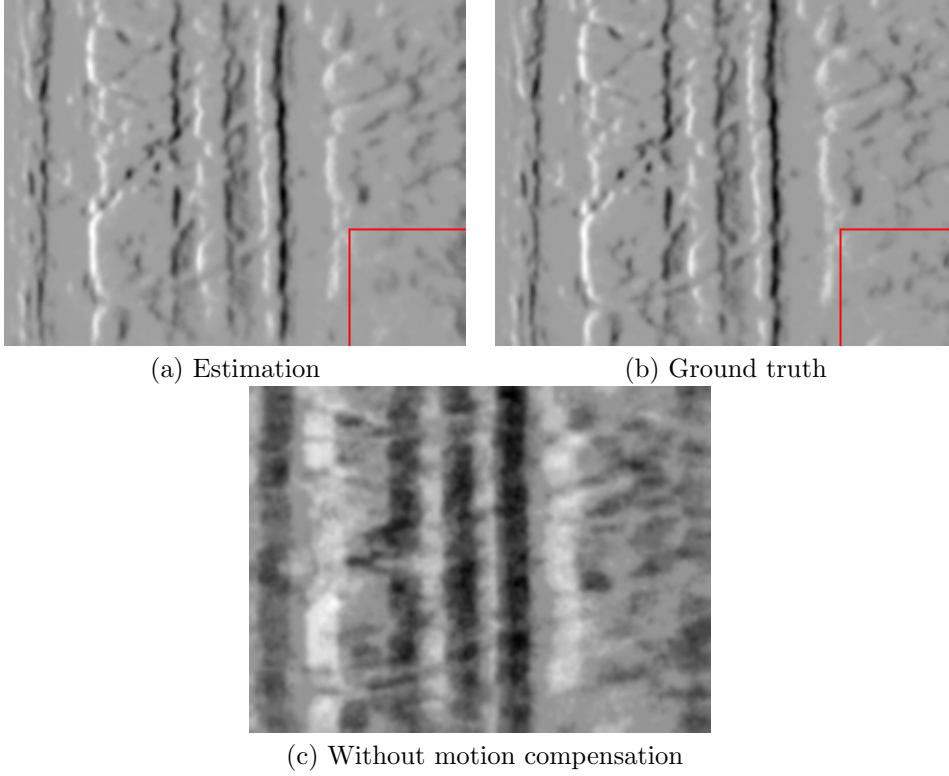


Figure 3.2: An example of the local optima problem. This is the dataset *slider_hdr_close* with a window size of 50000 events. (a) shows the optimized image with linear velocity $\mathbf{v} = (0.231, 0.109, 0.256)\text{m/s}$, angular velocity $\boldsymbol{\omega} = (0.405, -0.130, -0.278)\text{rad/s}$ and plane normal $\mathbf{n} = (-0.579, 0.282, -0.765)$. (b) shows the result using ground truth parameters with $\mathbf{v} = (0.163, 0, 0)$, $\boldsymbol{\omega} = (0, 0, 0)$ and $\mathbf{n} = (0, 0, -1)$. Both images look almost identical, though at the lower right corner, for example, one can still recognize the difference. Both images appear much sharper than the image without motion compensation in (c). It is worth mentioning that the contrast of the estimation is actually slightly larger than that of the ground truth

where $\mathbf{z} = (0, 0, -1)$ denotes the plane frontal-parallel to the camera.

The first per-frame optimizing procedure can be understood as projecting the events on a *blank canvas*. Similarly, in the *projecting-to-map* procedure, we project the events on the *texture* of the map and measure the strength of the composed image with the same variance function as in eq. (2.3), thus finding the set of the parameters that best align the events in the current frame to their correspondences in the texture. Note that in the project-to-map procedure we have to drop the polarity of the event, since the map might include events from the same visual stimuli but triggered when moving in a different direction; when we sum such events together, the polarities might cancel each other out.

The optimized velocity is used for propagating the pose to the next incoming frame via

$$\mathbf{t}_{wc_2} = \mathbf{R}_{wc_1} \mathbf{v} \Delta t + \mathbf{t}_{wc_1} \quad (3.18)$$

$$\mathbf{R}_{wc_2} = \mathbf{R}_{wc_1} \exp(\boldsymbol{\omega}^\wedge \Delta t), \quad (3.19)$$

where c_1 and c_2 denotes the current frame and the next frame, respectively, and Δt is the temporal size of the current frame.

3.3.3 Mapping

After having collected the first N events, we initialize the mapping process. For the first frame we estimate the full parameter set $\boldsymbol{\theta} = (\boldsymbol{\omega}^\top, \mathbf{v}^\top/d_w, \varphi, \psi)^\top \in \mathbb{R}^8$, where (φ, ψ) parametrize the unit vector of the plane normal, and \mathbf{v}^\top/d_w account for the scale ambiguity problem introduced in linear velocity estimation from monocular camera. We can, for example, determine the scale of the scene by setting d_w to 1 m, then the scales of the consecutive frames are determined by $d_c = d_w + \mathbf{t}_{wc} \cdot \mathbf{n}_w$.

From the optimized first frame, we initialize the map by projecting the frame to the planar scene via a rotation \mathbf{R}_n as in eq. (3.17). Then we can track the next frames with this map using the method described in section 3.3.2.

As the camera keeps moving, there might be new information available in the scene, so that the map needs to be expanded. After the optimization for every single frame, we also measure the *per-pixel-overlap* between the frame and the map, that is, how many percent pixels in the current frame can be explained by the map. A perfect match should have a 100% overlap. The overlap might be small when the camera is just exploring a new area so that only part of the frame and the map is overlapping; another possibility is that the map is not accurate enough to explain the current frame, which often happens when only one frame is used until now to construct the map. When the overlap drops below a certain threshold (80% for example), we try to insert a new *keyframe*.

When the system decides that a new keyframe is needed, we collect all the keyframes until now, plus the current frame, which is a keyframe candidate, and optimize the poses and velocities of these frames as well as the map altogether. Suppose there are n frames (including the keyframe candidate), the to be optimized parameter set is

$$\boldsymbol{\theta} = (\mathbf{R}_{(1 \sim n)}, \mathbf{t}_{(1 \sim n)}, \boldsymbol{\omega}_{(0 \sim n)}, \mathbf{v}_{(0 \sim n)}/d_w, \varphi, \psi) \in \mathbb{R}^{12n-4}.$$

Here $1 \sim n$ means from frame 1 to frame n , and we skip the pose optimization of frame 0, which is the first frame, since it defines the global frame. We project all the events of these n frames to the plane parametrize by (φ, ψ) with eq. (3.14), and again optimize the contrast of the synthesized image to obtain the optimal parameter set. For the same reason as before, we also don't use the polarity when fusing keyframes.

After this step, we also measure the *per-pixel-overlap* between the keyframe candidate and the image synthesized by the other keyframes, with the newly optimized parameter set. If this is a good match, we continue the tracking. Otherwise, we consider the pose estimation of the current frame to be inaccurate. In this case, we preserve the former map, and use this map to start the *relocalization* procedure.

3.3.4 Relocalization

A relocalization procedure purely depending on nonlinear optimization is hard, since we always need a good initial guess. If the pose is already too much off, there is nothing much we can help. The suggestion is not to start relocalization only when it's too late. Later in section 3.4 we will talk a little more about the situation where the tracking is lost when there is little texture visible.

When relocalizing, we still project the current events onto the map using eq. (3.14), this time with two additional parameters $(\mathbf{R}_{wc}, \mathbf{t}_{wc})$. As the initial guess, we use the current pose estimation. If the per-pixel overlap after optimization is still low, we restart relocalization with the pose of the last frame. If it still fails, we consider the current tracking as lost.

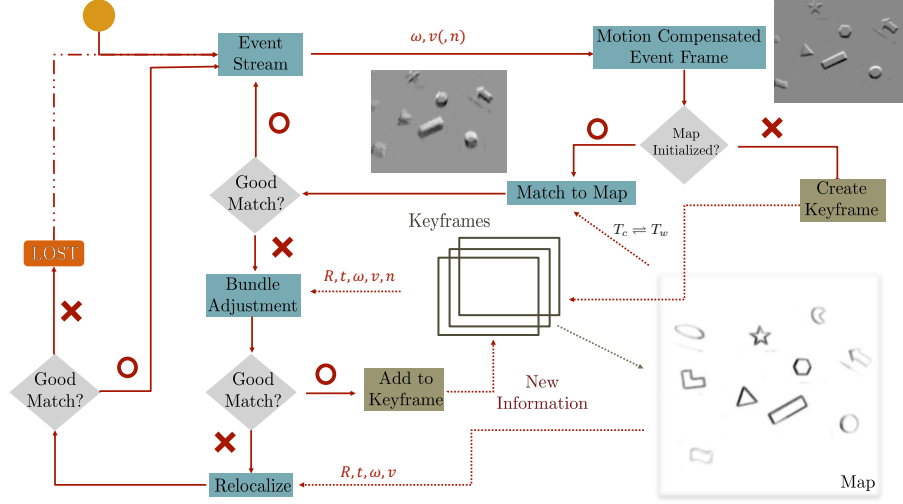


Figure 3.3: The parallel tracking and mapping pipeline. All the frames are used for tracking, but only the keyframes are used for mapping

3.4 Experiments

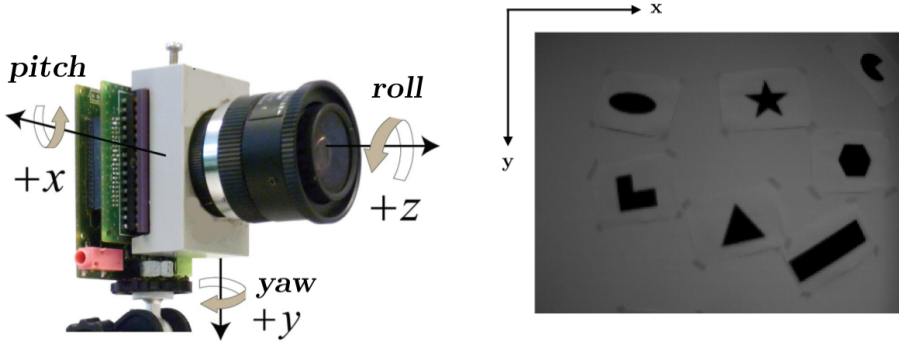
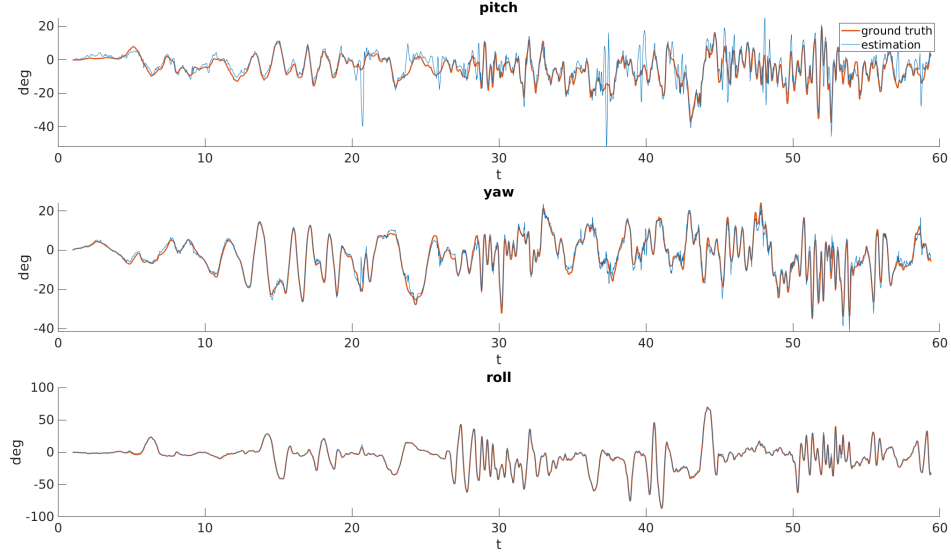


Figure 3.4: Axes definition. Images on the left is adapted from [18]

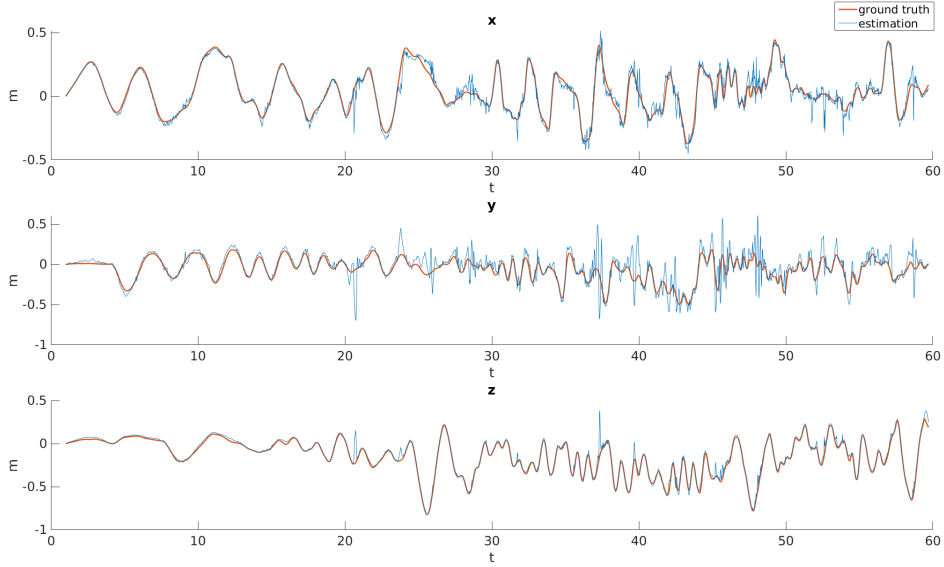
The algorithm is tested on four planar datasets provided by the *Robotics and Perception Group* at UZH [20], which includes *poster_6dof*, *poster_translation*, *shapes_6dof*, *shapes_translation*. The *shapes* dataset contains several simple shapes on the wall, whereas the *poster* contains richer rock texture. A planar assumption is redundant for the rotation datasets and thus the results are not listed for comparison here.

We define the camera axes as in fig. 3.4. The camera coordinate frame of the first synthesized frame is set as the global frame. The x, y and z axes are also pitch yaw, roll axes, respectively.

Figure 3.5 shows the comparison of the results of the proposed method against ground truth for the whole *shapes_6dof* sequence. The orientation and translation estimation along z axis are the most accurate among all the axes; when zoomed in (fig. 3.6), it is clear to see that the roll estimation is almost indistinguishable from the ground truth. We split the quantitative error estimation for orientation along 3 axes (table 3.1), in order to show that this holds true for all the datasets we tested



(a) Orientation



(b) Translation

Figure 3.5: *shapes_6dof* sequence. Comparison of the estimated pose (blue line) against ground truth (red line).

on. In these datasets, the global z axis is close to the plane normal; we believe that the events generated by a motion along the plane normal direction are more obvious on the image plane, and also harder to be explained by a combination of other motions, thus they are more reliable.

Although the algorithm tracks very well in the first 20s, in the latter half of the sequence there are sometimes large fluctuations. We investigated the sequence and found multiple reasons for that:

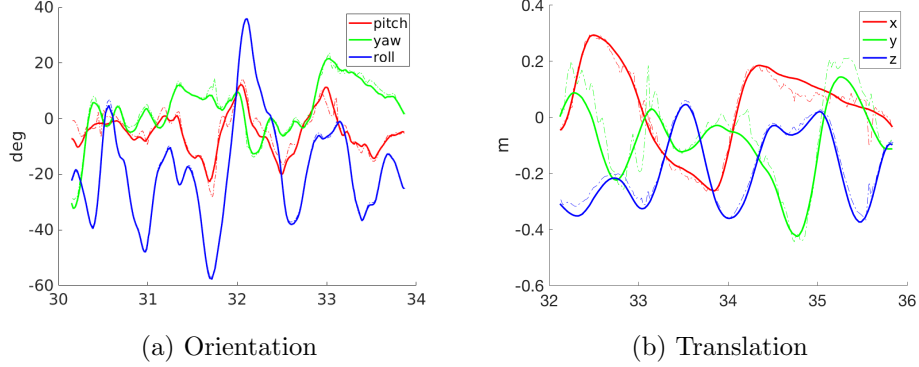


Figure 3.6: *shapes_6dof* sequence. Comparison of the estimated pose (dashed line) against ground truth (full line) within segments of 8 seconds duration.

1. The window size is not suitable

The window size is manually chosen for each dataset. A range of 2000 – 5000 usually works good for the *shapes* datasets, since the scene consists of simple regular shapes. Figure 3.7(a) shows what the camera frustum normally contains throughout the *shapes* sequences. However, in the rapid movement phase, the camera very often moved too far away from the scene, so that other scenes which don't belong to *shapes* also become visible, e.g. at around 47.8 seconds of dataset *shapes_6dof* the scene looks like in fig. 3.7(b); at this moment the *poster* scene on the wall is also visible, which has much more complex textures (when performing pose estimation on the *poster* sequences, we normally choose a window size of 30000–50000). Apparently, a window size of 4000 is no longer suitable in this case, which causes the motion estimation in this period being inaccurate.

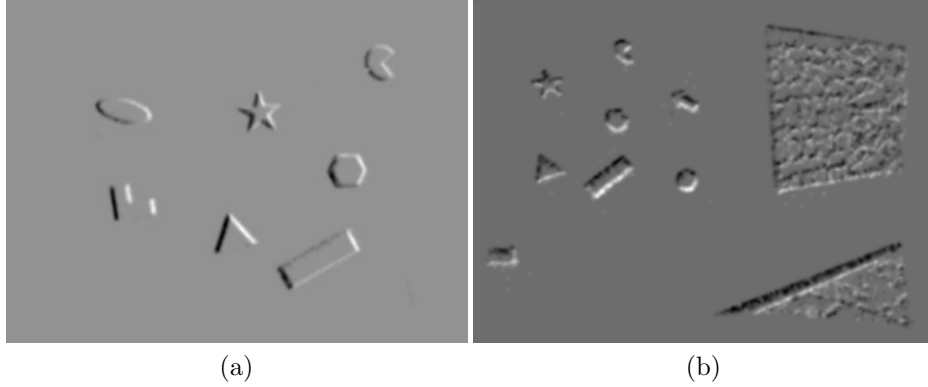
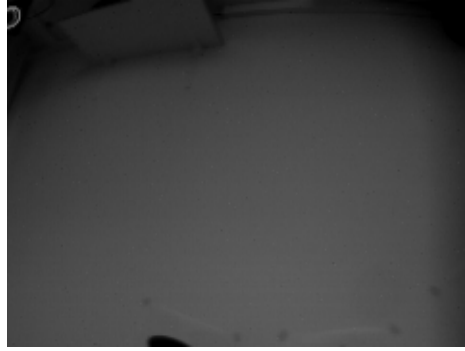


Figure 3.7: A window size of 4000 for a specific dataset *shapes_6dof* but different scene complexity

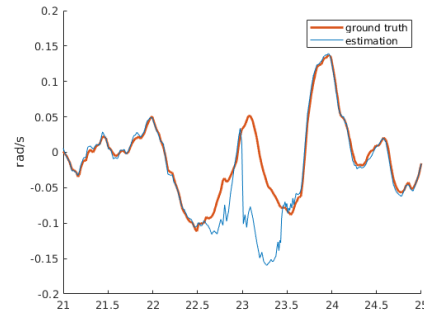
2. There is little texture available

A very representative example is that when tracking the dataset *shapes_translation*, we always get lost at around 23s. Investigating the scene around this moment (fig. 3.8(a)), we found that very little or even no texture is available. It's also not possible to perform a reliable pose estimate during this period. For the *per-frame* velocity estimation proposed in [18, 19], where there is no link between frames, this is usually not a serious problem, since when the texture becomes available again, the optimization pipeline is still able to recover, it just might be harder to start from a bad initialization. However, for the pose estimation where a cumulative result is

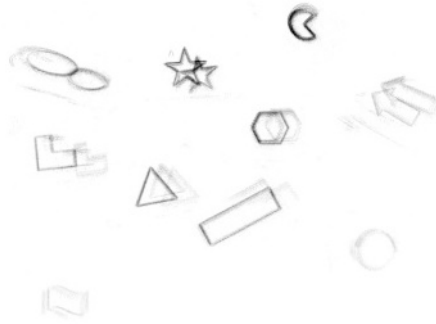
needed, a few badly estimated frames might destroy the whole following sequence. The *bundle adjustment* and *relocalization* parts of the pipeline is designed to deal with such situations. However, a bad initial guess of the pose (\mathbf{R}, \mathbf{t}) is in general harder to start with than a bad initial guess of velocities $(\boldsymbol{\omega}, \mathbf{v})$. When the pose estimation from the last frame is already too off, after optimization there might still be little overlap between the current frame and the pose. This is when the *bundle adjustment* tries to insert a keyframe and track from the new keyframe. The current implementation doesn't reinitialize the map after the tracking is lost. The advantage is that there is still chance to come back to the original map after tracking a wrong "local map" (fig. 3.8(b)); the disadvantage is that the global map will most likely be polluted by the "local map" (fig. 3.8(c))



(a) There is almost no texture visible



(b) Relocalized after lost (roll component)



(c) Polluted map after tracking is lost

Figure 3.8: When tracking the dataset *shapes_translation*, we always get lost at around 23s, but there is still chance to relocalize

A quantitative error estimation is included in table 3.1. Considering the large fluctuations that occasionally happens, which is caused by the inadequate amount of texture, will severely influence the mean error, we choose the median value for quantitative error estimation instead. The window sizes for *shapes* and *poster* sequences are 4000 and 30000, respectively. We use the scientific library *GNU-GSL* to perform the contrast maximization. The *Nelder-Mead Simplex* algorithm, which does not require derivatives, is applied to the *shapes* sequences; the *Fletcher-Reeves conjugate gradient*[24] algorithm with derivatives is applied to the *poster* dataset. Both methods work. The analytic differentiation is faster and more stable, however, when the initial guess is further from the ground truth, the numeric differentiation sometimes works better than the analytic differentiation, since the jacobians used

for analytic differentiation are strictly local (see appendix A), whereas a properly chosen step size for numeric differentiation is able to explore to a further region. We also observe that there is no drift along the trajectory, since we always apply an image-to-model matching. Also, the *6dof* sequences tend to have larger error than the *translation* dataset; so we tried to initialize the map where the plane is around 45° tilted to the camera, and the tracking seems worse. It looks like the estimation is less accurate when the plane is too much tilted.

Dataset	MOE/ $^\circ$			MTE/m	TD/m	Note
	pitch	yaw	roll			
shapes_6dof	1.89	1.01	0.85	0.047	46.89	Tracked all 60 s
shapes_translation	0.85	0.45	0.27	0.019	14.08	Lost after 25.0 s
poster_6dof	1.14	1.22	0.63	0.054	12.50	Lost after 22.8 s
poster_translation	0.55	0.57	0.25	0.022	11.14	Lost after 26.0 s

Table 3.1: Quantitative evaluation on planar sequences. *MOE*: median orientation error; *MTE*: median translation error; *TD*: traveled distance

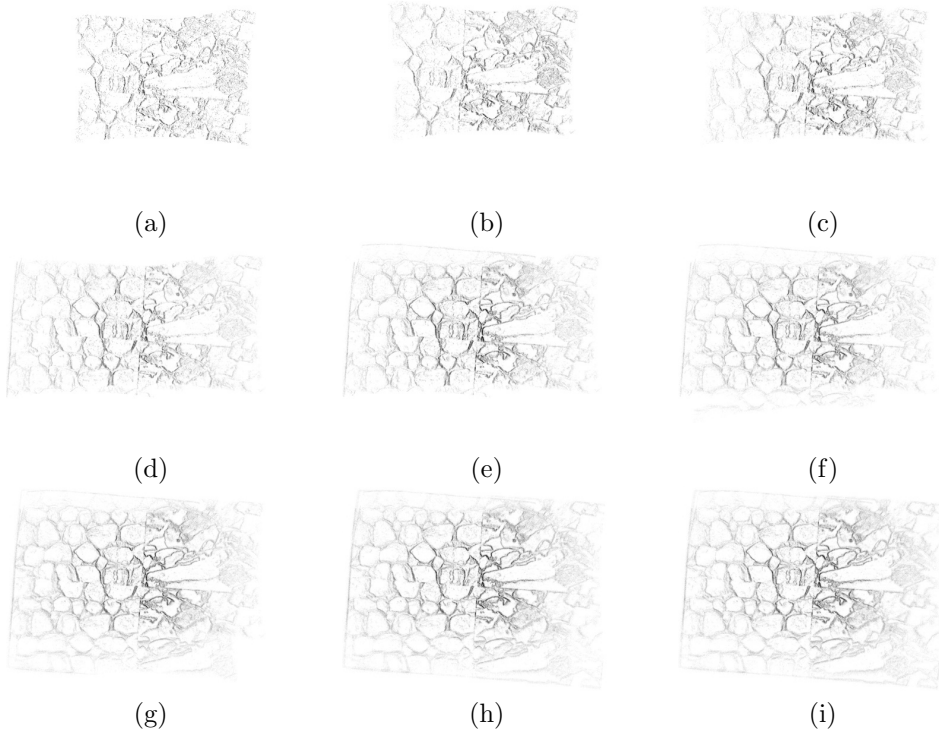


Figure 3.9: Map evolution of *poster_translation*. With a window size of 30 000, the scene is fully explored after 18 s, 274 frames, among which 10 are keyframes

Despite the rapid movements, in the above datasets the camera actually only moves in a relatively small region in front of the textured plane, as depicted in fig. 3.10. However, this method also applies when the camera travels a longer distance with respect to the scene depth, as shown in fig. 3.11. Also, although this method is designed for planar scenes, in a complex scene with rotation only movement, where no scene parameters are needed, this method can still be applied to build a global map, delivering a result similar to *image stitching*. Figure 3.12 is such an example.

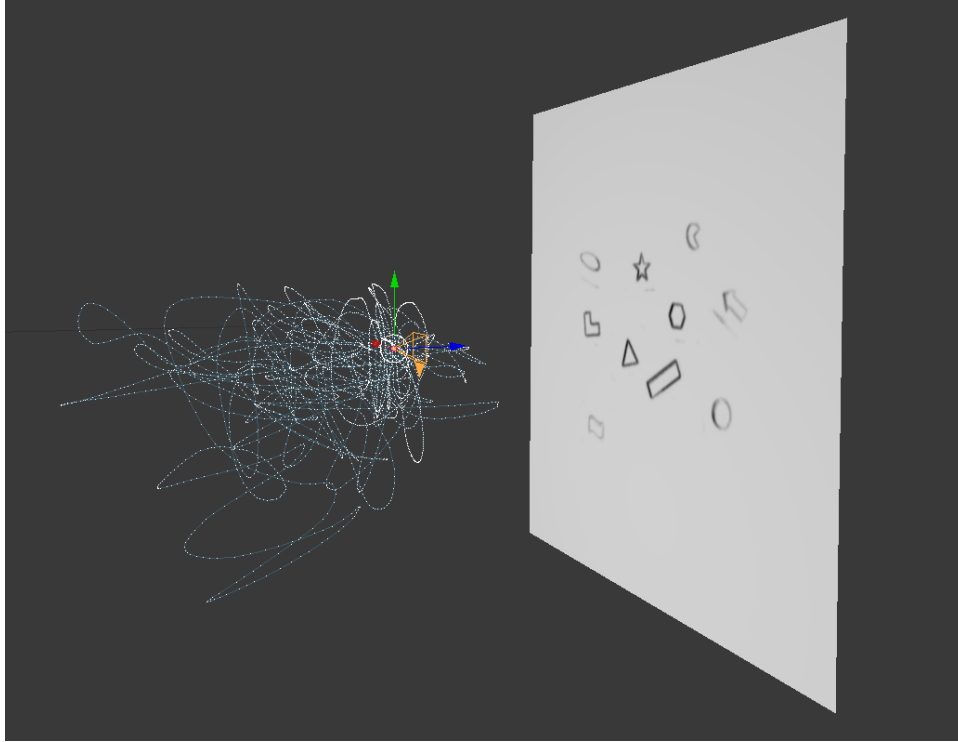


Figure 3.10: The motion path of the dataset *shapes_6dof*. Despite a trajectory length of about 50 m, in the whole 60 s the camera actually stays in a relative small region compared to the scene depth, and there is no observable increase of drift using the method in this work.

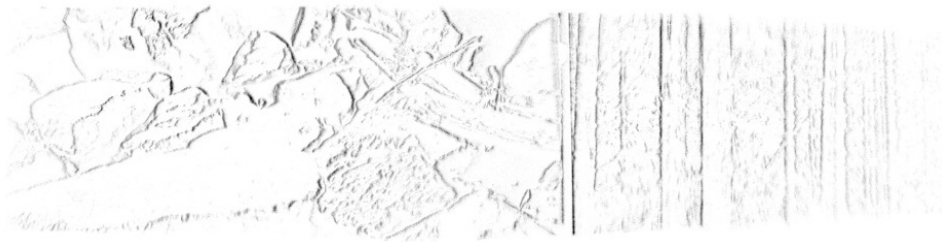


Figure 3.11: The dataset *slider_hdr_far*, with a scene depth of 0.584 m. Note that this is a much wider map than what we have seen before. The camera is frontal-parallel to the plane and moves along the positive x . As a result, a correct map should remain the same height. For our estimation, this is true for the first 4.8 s within the whole range of 6.3 s; afterwards, when tracking in the forest, almost all the textures are vertical, causing severe local optima problem. Note how the map tilted while the edges remain sharp. If we constrain the motion estimation to translation only, it delivers a much better result.



Figure 3.12: The algorithm is able to build a map from a complex scene as in dataset *boxes_rotation*. The motion is rotation-dominated

Chapter 4

6DoF Motion Estimation in General 3D Scenes

For general 3D scenes with complex scene structure, we may consider dividing the scene into multiple patches, and assume each patch is a plane.

Without knowing which parts of the scene correspond to real planes in space, we simply divide the image into $(4 \times 3 =) 12$ regular grids, each corresponds to a 30×30 patch with a different set of parameters (2 for plane normal and 1 for depth). We assign each event to the patch it belongs to according to the pixel coordinate, and optimize the frame with the same motion model. Accounting for the scale ambiguity, this amounts to $(3 \times 12 + 5 =) 41$ parameters for each *per-frame* optimization. The scale ambiguity is solved by parametrizing the linear velocity with 2 parameters $(\phi, \psi) \in \mathbb{S}^2$ indicating its direction, the same as the parametrization of a plane normal. There is a singularity at zero velocity; however, for the hand-held datasets we tested on this doesn't appear to be a problem. Note that we don't need to optimize each patch individually and sum the costs together, since the intensity computation eq. (2.1) naturally sums all the events together; we only need one cost function per image. Also, the regions with fewer events have more pixels with zero intensity, which means a lower contribution to the cost function eq. (2.3), thus no extra weighting is needed.

In fig. 4.1, we compare the synthesized images before and after motion correction of some selective images in various sequences. After motion correction, the edges appear much sharper.

Since each frame has a different scale, without a global map as in section 3.3, it's hard to measure the accuracy of the linear velocity estimation. We give the estimation of angular velocity in table 4.1. Comparing line 2 with line 3, we can see that for general 3D scenes, planar patches assumption delivers a better result than assuming the whole scene as one single plane, which is a very strong assumption. We didn't include the test result of planar assumption on *dynamic_dof* (line 1) and *outdoors_running*, because the trials got lost very soon. Notice that planar patches assumption for general scenes (line 1 and 2) has a similar error magnitude as a planar assumption for planar scenes (line 5 and 7). However, the angular velocity estimated by either planar or planar patches assumption cannot achieve the same accuracy when estimating scenes with rotation dominated movements using 3 rotational parameters (line 4 and 6), because nonlinear optimization with more parameters suffers more from local optima problem.

Without a known scale for each frame, we can compare the angles between ground truth and estimated linear velocities. It turns out that either for a planar assumption on planar scenes or planar patches assumption on general scenes, the mean

angle is

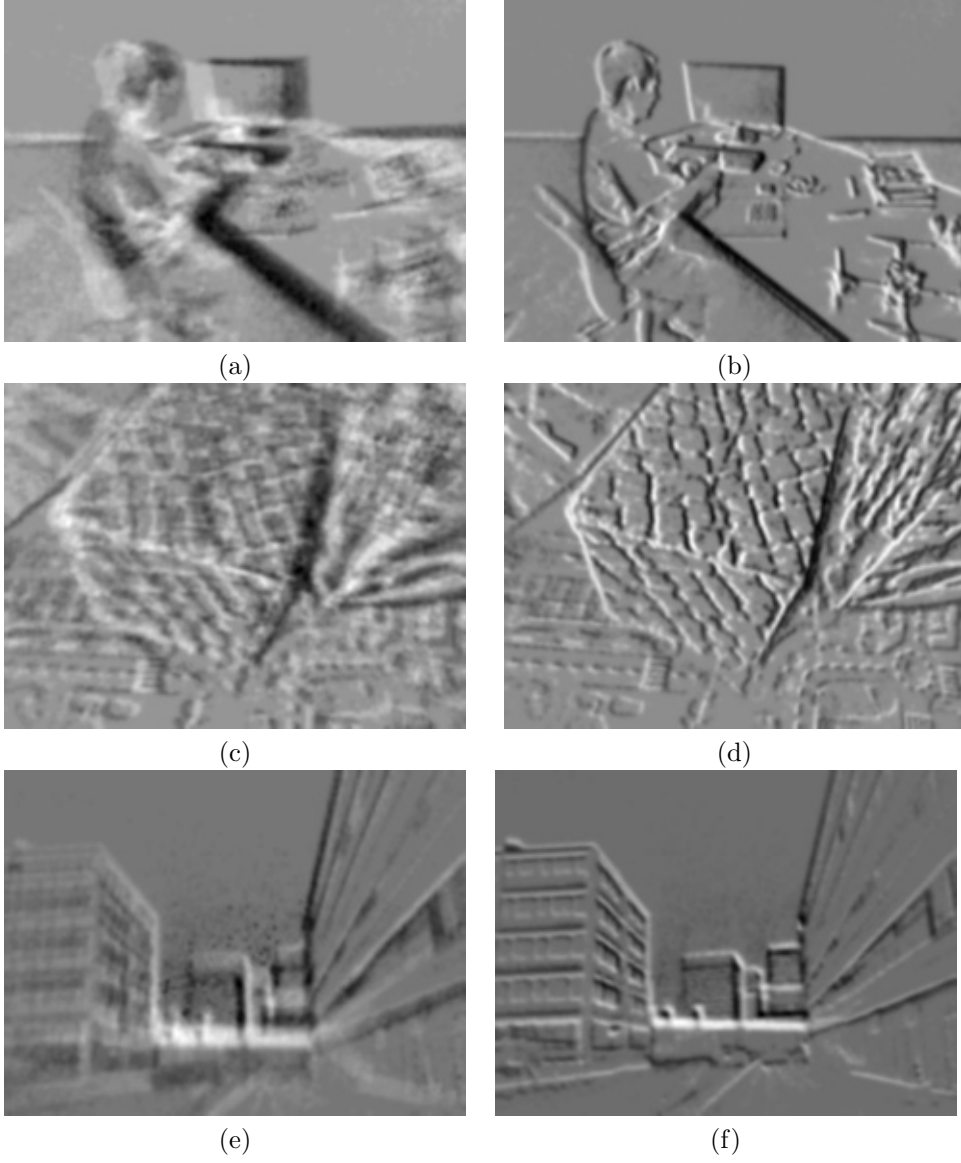


Figure 4.1: Test results on 3 sequences with 6DoF motion. *Left:* Accumulated events without motion compensation. *Right:* Images after motion compensation via dividing the scene into 12 regular grids. (a)(b) *dynamic_6dof* sequence, with ground truth velocities $\omega = (-0.145; -3.761; 1.925)\text{rad/s}$, $\mathbf{v} = (0.313; 0.855; 0.326)\text{m/s}$. (c)(d) *boxes_6dof* sequence, with ground truth velocities $\omega = (-1.357; 0.981; 4.577)\text{rad/s}$, $\mathbf{v} = (2.805; 0.538; -0.676)\text{m/s}$. (e)(f) *outdoors_running* sequence, with no ground truth information available.

around 50° , which is much larger than that of angular velocity, which is around 20° . We guess that the temporal window is not large enough to provide a sufficient baseline for linear velocity estimation. The window size cannot be chosen too large, otherwise, the constant velocity assumption would fail.

The short baseline problem can be solved by tracking over a longer distance, as is shown in chapter 3. Similarly, we could try to randomly initialize the patch

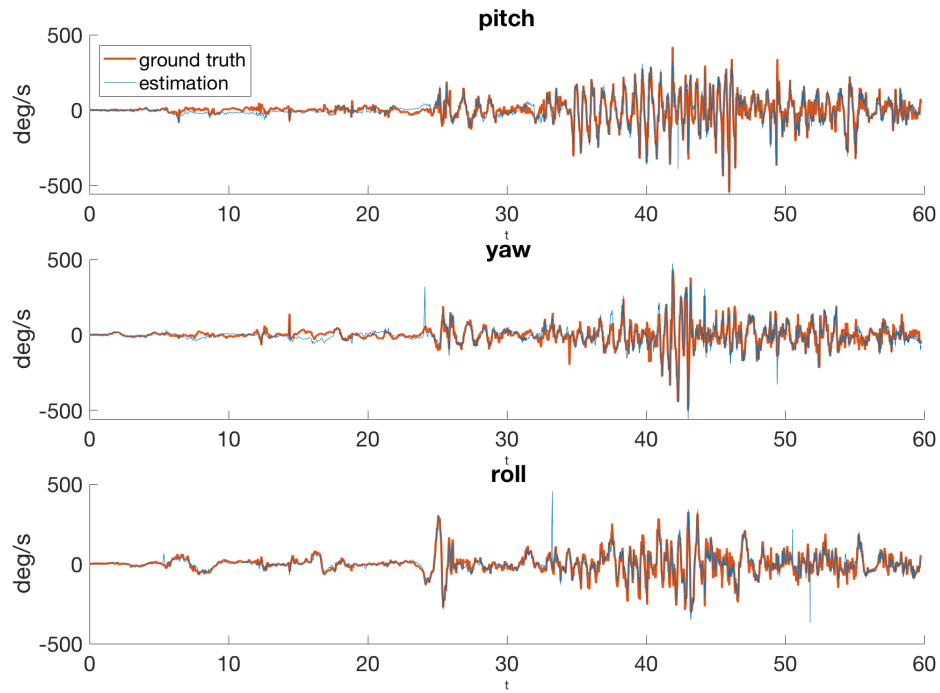


Figure 4.2: Test result on *boxes_6dof*. Comparison of the estimated angular velocity against the ground truth. In fast motion phase, the estimation appears more accurate

locations and track with the patches. However, if the selected patch is not a plane, it can not be tracked reliably over a long distance, and we need good strategies to delete unreliable patches while adding new patches. Resolving this problem is part of our future works.

No.	Dataset	Median Angular Velocity Error°/s			Maximal Velocity °/s	Parametrization
		pitch	yaw	roll		
1	dynamic_6dof	20.88	21.60	18.08	509.37	6DoF, 12 patches
2	boxes_6dof	34.69	27.49	20.21	428.19	6DoF, 12 patches
3	boxes_6dof	60.16	31.33	21.11	428.19	6DoF, planar assumption
4	boxes_rotation	9.39	9.84	12.55	679.78	3DoF (angular velocity)
5	poster_6dof	32.23	31.55	22.37	905.58	6DoF, planar assumption
6	poster_rotation	11.62	8.47	13.23	934.88	3DoF (angular velocity)
7	poster_rotation	36.45	27.00	26.14	905.58	6DoF, planar assumption

Table 4.1: Comparison between planar assumption, planar patches assumption and 3DoF angular velocity estimation. The maximal velocity in the table is the maximal velocity among all three axes

Chapter 5

Conclusion and Outlook

In this work, we developed a system that performs SLAM in planar scenes with high accuracy, with only events and no additional sensing. In unstructured scenes with 6-DoF motion, the planar patches method also delivers sharp images, even under high speed motions.

There is still a large space for improvement. For example, we do not yet have a reliable strategy to handle the situation where the tracking is lost. Also, because of the time limit, we haven't yet made the SLAM system to also work in the unstructured scenes.

We found that the calibration matrix provided in [20] is not very accurate, edges close to the upper left corner of the image always appear curved. This is observable at the upper left corner of fig. 3.9(i) and fig. 3.12. We checked the provided calibration dataset, and found that the checkerboard has barely any corners in the upper left part of the image, which should be the reason why the undistortion in that area is not working well. This might be one of the causes why sometimes the image-to-model matching is not performing well. We also tried to cropped part of the image, but the resolution of DAVIS is relatively low (240×180), cropping reduces the amount of information and does not work better.

Usually, the maps of *shapes* and *poster* sequences are done after 6 and 10 keyframes, respectively; as the number of keyframes increases, the size of the parameter set also increases (116 parameters in the bundle adjustment phase with 10 keyframes), which makes the problem hard to optimize. For a least square optimizing problem, which is often encountered in visual odometry, the Levenberg-Marquardt algorithm[25] is a good iteration method. However, the cost function eq. (2.3) used in this work is a more general non-linear function, and thus harder to optimize. The Nelder-Mead Simplex algorithm[26] used in bundle adjustment phase is better suitable for problems with small parameter size, and the choice of the initial simplex size is important. In the per-frame optimization and frame-to-model matching phase we make use of the derivatives; both BFGS and conjugate gradient approximate the objective function by a quadratic hypersurface, thus work well only when the evaluation point is close enough to the optimum. For a more reliable tracking over a longer distance, we should consider if better strategies are available, for example optimize only a subset of the keyframes that are most important for the current frame or the most recent frames, or try multiple starting points.

Bibliography

- [1] P. Lichtsteiner, C. Posch, and T. Delbruck, “A 128×128 120 db 15 μ s latency asynchronous temporal contrast vision sensor,” *IEEE journal of solid-state circuits*, vol. 43, no. 2, pp. 566–576, 2008.
- [2] C. Brandli, R. Berner, M. Yang, S.-C. Liu, and T. Delbruck, “A 240×180 130 db 3 μ s latency global shutter spatiotemporal vision sensor,” *IEEE Journal of Solid-State Circuits*, vol. 49, no. 10, pp. 2333–2341, 2014.
- [3] C. Li, C. Brandli, R. Berner, H. Liu, M. Yang, S.-C. Liu, and T. Delbruck, “Design of an rgbw color vga rolling and global shutter dynamic and active-pixel vision sensor,” in *Circuits and Systems (ISCAS), 2015 IEEE International Symposium on*. IEEE, 2015, pp. 718–721.
- [4] D. P. Moeys, F. Corradi, C. Li, S. A. Bamford, L. Longinotti, F. F. Voigt, S. Berry, G. Taverni, F. Helmchen, and T. Delbruck, “A sensitive dynamic and active pixel vision sensor for color or neural imaging applications,” *IEEE transactions on biomedical circuits and systems*, vol. 12, no. 1, pp. 123–136, 2018.
- [5] B. Son, Y. Suh, S. Kim, H. Jung, J.-S. Kim, C. Shin, K. Park, K. Lee, J. Park, J. Woo *et al.*, “A 640×480 dynamic vision sensor with a 9 μ m pixel and 300meps address-event representation,” in *Solid-State Circuits Conference (ISSCC), 2017 IEEE International*. IEEE, 2017, pp. 66–67.
- [6] A. Z. Zhu, N. Atanasov, and K. Daniilidis, “Event-based feature tracking with probabilistic data association,” in *Robotics and Automation (ICRA), 2017 IEEE International Conference on*. IEEE, 2017, pp. 4465–4470.
- [7] C. Harris and M. Stephens, “A combined corner and edge detector.” in *Alvey vision conference*, vol. 15, no. 50. Citeseer, 1988, pp. 10–5244.
- [8] D. Tedaldi, G. Gallego, E. Mueggler, and D. Scaramuzza, “Feature detection and tracking with the dynamic and active-pixel vision sensor (davis),” in *Event-based Control, Communication, and Signal Processing (EBCCSP), 2016 Second International Conference on*. IEEE, 2016, pp. 1–7.
- [9] E. Mueggler, C. Bartolozzi, and D. Scaramuzza, “Fast event-based corner detection,” in *28th British Machine Vision Conference (BMVC)*, 2017.
- [10] H. Rebecq, G. Gallego, and D. Scaramuzza, “Emvs: Event-based multi-view stereo,” in *British Machine Vision Conference (BMVC)*, no. EPFL-CONF-221504, 2016.
- [11] G. Gallego, J. E. Lund, E. Mueggler, H. Rebecq, T. Delbruck, and D. Scaramuzza, “Event-based, 6-dof camera tracking from photometric depth maps,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.

- [12] H. Rebecq, T. Horstschaefer, and D. Scaramuzza, “Real-time visualinertial odometry for event cameras using keyframe-based nonlinear optimization,” in *British Machine Vis. Conf.(BMVC)*, vol. 3, 2017.
- [13] 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.
- [14] H. Rebecq, T. Horstschaefer, 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, 2017.
- [15] G. Orchard, C. Meyer, R. Etienne-Cummings, C. Posch, N. Thakor, and R. Benosman, “Hfirst: a temporal approach to object recognition,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 37, no. 10, pp. 2028–2040, 2015.
- [16] A. I. Maqueda, A. Loquercio, G. Gallego, N. Garcia, and D. Scaramuzza, “Event-based vision meets deep learning on steering prediction for self-driving cars,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 5419–5427.
- [17] A. Z. Zhu, L. Yuan, K. Chaney, and K. Daniilidis, “Ev-flownet: Self-supervised optical flow estimation for event-based cameras,” *arXiv preprint arXiv:1802.06898*, 2018.
- [18] 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.
- [19] G. Gallego, H. Rebecq, and D. Scaramuzza, “A unifying contrast maximization framework for event cameras, with applications to motion, depth, and optical flow estimation,” in *IEEE Int. Conf. Comput. Vis. Pattern Recog.(CVPR)*, vol. 1, 2018.
- [20] E. Mueggler, H. Rebecq, G. Gallego, T. Delbruck, and D. Scaramuzza, “The event-camera dataset and simulator: Event-based data for pose estimation, visual odometry, and slam,” *The International Journal of Robotics Research*, vol. 36, no. 2, pp. 142–149, 2017.
- [21] E. Mueggler, C. Forster, N. Baumli, G. Gallego, and D. Scaramuzza, “Lifetime estimation of events from dynamic vision sensors,” in *Robotics and Automation (ICRA), 2015 IEEE International Conference on*. IEEE, 2015, pp. 4874–4881.
- [22] A. A. Michelson, *Studies in optics*. Courier Corporation, 1995.
- [23] B. Gough, *GNU scientific library reference manual*. Network Theory Ltd., 2009.
- [24] R. Fletcher, *Practical methods of optimization*. John Wiley & Sons, 2013.
- [25] W. H. Press, S. A. Teukolsky, W. T. Vetterling, and B. P. Flannery, “Numerical recipes in c,” *Cambridge University Press*, vol. 1, p. 3, 1988.
- [26] J. A. Nelder and R. Mead, “A simplex method for function minimization,” *The computer journal*, vol. 7, no. 4, pp. 308–313, 1965.

Appendix A

Derivative of the Contrast Metric

Since we use bilinear voting to evaluate the *dirac delta*, the derivatives of the cost function eq. (2.3) can be analytically computed as

Let

$$\rho(\mathbf{x}; \boldsymbol{\theta}) = \mathcal{I}(\mathbf{x}; \boldsymbol{\theta}) - \mu(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta})) \quad (\text{A.1})$$

with

$$\mu(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta})) = \frac{1}{|\Omega|} \int_{\Omega} \mathcal{I}(\mathbf{x}; \boldsymbol{\theta}) d\mathbf{x}, \quad (\text{A.2})$$

then

$$\frac{\partial}{\partial \boldsymbol{\theta}} \text{Var}(\mathcal{I}(\mathbf{x}; \boldsymbol{\theta})) = \frac{2}{|\Omega|} \int_{\Omega} \rho(\mathbf{x}; \boldsymbol{\theta}) \frac{\partial \rho(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} d\mathbf{x} \quad (\text{A.3})$$

The derivatives of the warped image are

$$\frac{\partial \mathcal{I}(\mathbf{x}; \boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = - \sum_{k=1}^N \pm_k \nabla \delta(\mathbf{x} - \mathbf{x}'_k(\boldsymbol{\theta})) \frac{\partial \mathbf{x}'_k(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad (\text{A.4})$$

In the case of planar homography the events warping is computed as

$$\boldsymbol{\theta} = (\boldsymbol{\omega}^\top, \mathbf{v}^\top, \varphi, \psi)^\top \in \mathbb{R}^8 \quad (\text{A.5})$$

$$\mathbf{x}'_k(\boldsymbol{\theta}) = [x'_{im} \quad y'_{im}]^\top = [x'/z' \quad y'/z']^\top \quad (\text{A.6})$$

$$\begin{aligned} \bar{\mathbf{x}}'_k(\boldsymbol{\theta}) &= [x' \quad y' \quad z']^\top \\ &= \mathbf{P} (\mathbf{R}_t^\top (\mathbf{I} + \mathbf{t}_t \mathbf{n}^\top / d))^{-1} \bar{\mathbf{x}}_k \\ &= \mathbf{P} (\mathbf{I} + \mathbf{v} \mathbf{t} \mathbf{n}^\top)^{-1} \mathbf{R}_t \bar{\mathbf{x}}_k, \end{aligned} \quad (\text{A.7})$$

where \mathbf{P} is the projection matrix from the camera frame to the world frame or the map, which does not have to correspond with the camera matrix provided by the data set. Also, for simplicity of notation we assume $d = 1$, and denote

$\mathbf{P}(\mathbf{I} + \mathbf{v}t\mathbf{n}^\top)^{-1}$ as \mathbf{P}_v . Since t usually spans a small temporal window, we can simplify the derivative with respect to angular velocity as

$$\mathbf{R}_t = \exp(\boldsymbol{\omega}^\wedge t) \approx \mathbf{I} + \boldsymbol{\omega}^\wedge t \quad (\text{A.8})$$

$$\frac{\partial \bar{\mathbf{x}}'_k}{\partial \boldsymbol{\omega}} = -\mathbf{P}_v t \bar{\mathbf{x}}_k^\wedge. \quad (\text{A.9})$$

The above equation makes use of the equivalence $\boldsymbol{\omega} \times \bar{\mathbf{x}}_k = -\bar{\mathbf{x}}_k \times \boldsymbol{\omega}$. The derivative with respect to the linear velocity is

$$\frac{\partial \bar{\mathbf{x}}'_k}{\partial \mathbf{v}} = -\frac{t\mathbf{n}^\top \mathbf{R}_t \mathbf{x}_k}{\mathbf{n}^\top \mathbf{v}t + 1} \mathbf{P}_v \quad (\text{A.10})$$

Similarly we have

$$\frac{\partial \bar{\mathbf{x}}'_k}{\partial \mathbf{n}} = -\frac{t\mathbf{v}^\top \mathbf{R}_t \mathbf{x}_k}{\mathbf{n}^\top \mathbf{v}t + 1} \mathbf{P}_v \quad (\text{A.11})$$

and

$$\frac{\partial \bar{\mathbf{x}}'_k}{\partial(\varphi; \psi)} = \frac{\partial \bar{\mathbf{x}}'_k}{\partial \mathbf{n}} \frac{\partial \mathbf{n}}{\partial(\varphi; \psi)}, \quad (\text{A.12})$$

with

$$\frac{\partial \mathbf{n}}{\partial(\varphi; \psi)} = \begin{bmatrix} -\sin \varphi \sin \psi & \cos \varphi \cos \psi \\ \cos \varphi \sin \psi & \sin \varphi \cos \psi \\ 0 & -\sin \psi \end{bmatrix} \quad (\text{A.13})$$

With the above quantities

$$\frac{\partial \bar{\mathbf{x}}'_k}{\partial \boldsymbol{\theta}} = \begin{bmatrix} \frac{\partial \bar{\mathbf{x}}'_k}{\partial \boldsymbol{\omega}} & \frac{\partial \bar{\mathbf{x}}'_k}{\partial \mathbf{v}} & \frac{\partial \bar{\mathbf{x}}'_k}{\partial \varphi} & \frac{\partial \bar{\mathbf{x}}'_k}{\partial \psi} \end{bmatrix} \in \mathbb{R}^{3 \times 8}$$

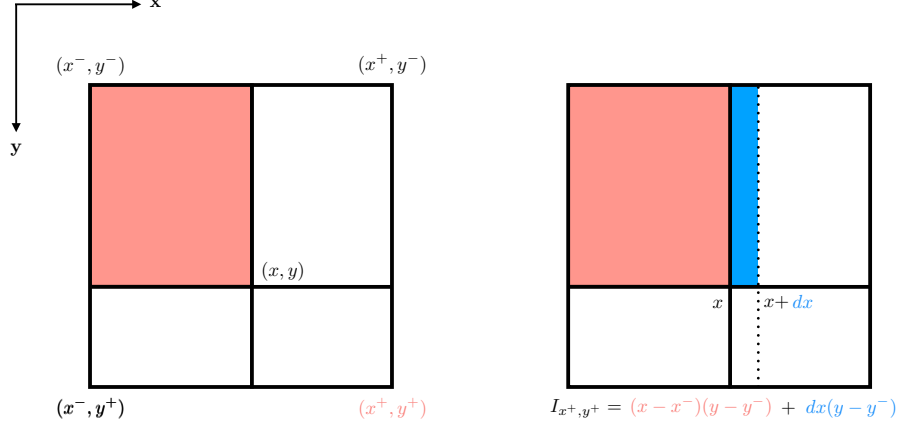
we can compute the gradient $\frac{\partial \mathbf{x}'_k}{\partial \boldsymbol{\theta}}$ as

$$\frac{\partial \mathbf{x}'_k}{\partial \boldsymbol{\theta}} = \left[\frac{\partial \bar{\mathbf{x}}'_k}{\partial \boldsymbol{\theta}} \frac{1}{z'} - \bar{\mathbf{x}}'_k \frac{\partial z'}{\partial \boldsymbol{\theta}} \frac{1}{z'^2} \right]_{1:2}, \quad (\text{A.14})$$

and the third component of the derivative within the bracket is equal to zero. Theoretically, to evaluate the derivative of the cost function, the dirac delta $\nabla \delta(\mathbf{x})$ should be evaluated at each pixel location (240×180 for DAVIS). However, since we are using bilinear voting, a infinitesimal change $d\mathbf{x}$ at an event location \mathbf{x} will only affect the derivative evaluated at the four neighboring pixels of \mathbf{x} . An illustration is shown in fig. A.1

The detailed algorithm for the evaluation of the Dirac delta as well as its derivative is shown in algorithm. 1. The behavior of the bilinear voting for derivatives is undefined when \mathbf{x}'_k locates exactly at a pixel location. However, the possibility of this to happen is zero except when $t = 0$, and can be solved by simply discarding the event or disturbing with a small random noise. When projecting back to the map, The author uses a project matrix with different focal lengths than the ones in the camera calibration matrix, so that there won't be event at integer pixel locations after warping.

For better convergence behavior, we smooth $\mathcal{I}(\mathbf{x}; \boldsymbol{\theta})$ by a Gaussian with a sigma of 1 as in[18]. Although for each event the jacobian only considers the four neighbors, the smoothing operation spreads this to broader pixels, leading to a smoother cost function which is easier to optimize.



(a) The intensity at a pixel location is the area of the rectangle spanned by the opposite pixel and the event location

(b) Intensity change after an infinitesimal movement of the event

Figure A.1: Bilinear voting

Algorithm 1: Bilinear Voting

Input : A warped event $e'_k = \{x'_k, y'_k, t_k\}$

Output: The *Dirac delta* $\delta(\mathbf{x} - \mathbf{x}'_k)$ as well as its derivative $\frac{\partial \delta(\mathbf{x} - \mathbf{x}'_k)}{\partial \mathbf{x}'_k}$

Let

$$x^- = \lfloor x'_k \rfloor, x^+ = \lceil x'_k \rceil, y^- = \lfloor y'_k \rfloor, y^+ = \lceil y'_k \rceil$$

Then

$$\delta(x^- - x'_k, y^- - y'_k) = (x - x^+)(y - y^+)$$

$$\delta(x^- - x'_k, y^+ - y'_k) = -(x - x^+)(y - y^-)$$

$$\delta(x^+ - x'_k, y^- - y'_k) = -(x - x^-)(y - y^+)$$

$$\delta(x^+ - x'_k, y^+ - y'_k) = (x - x^-)(y - y^-)$$

$$\frac{\partial \delta(x^- - x'_k, y^- - y'_k)}{\partial x'_k} = y - y^+, \frac{\partial \delta(x^- - x'_k, y^- - y'_k)}{\partial y'_k} = x - x^+$$

$$\frac{\partial \delta(x^- - x'_k, y^+ - y'_k)}{\partial x'_k} = -(y - y^-), \frac{\partial \delta(x^- - x'_k, y^+ - y'_k)}{\partial y'_k} = -(x - x^+)$$

$$\frac{\partial \delta(x^+ - x'_k, y^- - y'_k)}{\partial x'_k} = -(y - y^+), \frac{\partial \delta(x^+ - x'_k, y^- - y'_k)}{\partial y'_k} = -(x - x^-)$$

$$\frac{\partial \delta(x^+ - x'_k, y^+ - y'_k)}{\partial x'_k} = y - y^-, \frac{\partial \delta(x^+ - x'_k, y^+ - y'_k)}{\partial y'_k} = x - x^-$$
