



Visual Odometry for Navigating in GPS Denied Environments

Final Report

Alex Kreimer

Department of Computer Science

Abstract

In this work we revisit the problem of visual odometry. Visual odometry is the process of estimating the motion of the camera by examining the changes that the motion induces on the images made by it. This work consists of two parts. In the first part we propose a novel visual odometry algorithm. In the second part we describe a new data set.

The algorithm we propose exploits a scene structure typical for that seen by a moving car and is suitable for use in either the stereo or the monocular setting. We recover the rotation and the translation separately, thus dealing with two separate, smaller problems. The rotation is estimated by means of the infinite homography. The rotation estimation algorithm operates on distant image points using the 3-D to partition them into the distant and the near-by ones. We start with an initial estimate and then refine it using an iterative procedure. After the rotation is compensated for, the translation is found by means of the 1-point algorithm in the stereo setting and epipole computation for pure translational motion in the monocular setting.

We also created a new dataset that contains stereo video tracks synchronized with the DGPS locations of the vehicle. The dataset is suitable for future visual odometry experiments and is interesting because it has some challenging sequences (e.g. significant scene occlusions by a moving vehicle, fast camera motion, rural scenery, dusky sequences).

Contents

1	Introduction	2
2	The Infinite Visual Odometry Algorithm	4
2.0.1	Summary of Our Work	4
2.1	Preliminaries and Notation	5
2.1.1	Image Point Mapping Related to Camera Motion	5
2.2	Motion Estimation	6
2.2.1	Stereo	6
2.3	Experimental Results	9
2.3.1	The Choice of Features	9
2.3.2	Experimental Results	9
2.4	Conclusions and Discussion	10
3	The Dataset	11
3.1	Data Acquisition	11
3.1.1	Overview	11
3.1.2	Physical Setup	11
3.2	Post-processing	13
3.2.1	Image sequences	13
3.2.2	GPS Data	13
3.2.3	Missing GPS Data	13
3.3	VO results evaluation	13

Chapter 1

Introduction

Visual odometry refers to the problem of recovering camera motion based on the images taken by it. This problem naturally occurs in robotics, wearable computing, augmented reality and automotive.

Wheel odometry, recovers the motion of the vehicle by examining and integrating the wheel turns over time. In a similar manner, visual odometry operates by estimating relative motion of the camera between subsequent images by observing changes in them. Later, these estimates are combined into a single trajectory. Just as in wheel odometry, visual odometry is subject to error accumulation over time. Contrary to wheel odometry, visual odometry is not affected by wheel slip in a rough terrain. Visual odometry is able to produce motion estimates with errors that are lower than those of the wheel odometry. Another advantage of visual odometry is that cameras are low cost and low weight sensors. All these make visual odometry a viable supplement to other motion recover methods such as global positioning systems (GPS) and inertial measurement units (IMUs).

Visual odometry becomes a harder problem as the amount of detail in the images diminishes. The images should have sufficient overlap and the scene needs to be illuminated. In the stereo setup, the scene must be static or the images taken at the same time. Also, the video processing incurs a computational burden.

Visual odometry is an active fields of research with a large amount of published work. We review only the most pertinent works. [23] provides a more complete survey.

Similar to [21] we partition visual odometry algorithms by four traits:

1. Feature-based vs direct
2. Global vs local
3. Filter based vs bundle adjustment based
4. Monocular vs stereo

Visual odometry algorithms use large number of corner detectors (e.g., Moravec [17], Harris [7], Shi-Tomasi [24], Fast [22]) and blob detectors (e.g., SIFT [15], SURF [3]). Corners are faster to compute and usually are better localized, while blobs are more robust to scale change. The choice of a specific feature point depends mainly on the images at hand. Motion estimation results for different feature points are presented in [6]. In this work we choose Harris [7] corners, but this choice is not crucial. We view the feature point choice as a parameter, which needs to be determined from the data (e.g., by cross-validation).

The features are either tracked [10] or matched [5] (i.e., freshly detected in each new frame) between subsequent images. While the early works chose to track features, most of the current works detect and match them. The output of this stage are pairs of the image features, which are the projections of the same 3-D point.

Matched features are used as an input for a motion estimation procedure. Whether the features are specified in 2-D or 3-D, the estimation procedures, may be classified into 3-D-to-3-D [16], 3-D-to-2-D [5] and 2-D-to-2-D [19]. Most of the early works were of the 3-D-to-3-D type. More recent works [19] claim that this approach is inferior to the latter two. Popular techniques that participate in most algorithms in some way are essential matrix estimation and (possibly) its subsequent decomposition [19], perspective 3-point algorithm [13], and re-projection error minimization [5].

Global methods [1], [18] keep the map of the environment and make sure that motion estimates are globally consistent with this map, while local methods do not. Some local methods [2] also keep track of a (local) map, but the underlying philosophy is different: global vs local. Global methods usually more accurate since they make use of a vast amount of information (which, of course, comes at a computational price). Note that accuracy does not imply robustness, since outliers that made their way into the map may greatly skew subsequent pose estimates.

Methods that explicitly model system state uncertainty tend to use filtering mathematical machinery, e.g., [14], [20], [12]. Another alternative to maintain map/pose estimate consistency is to use the bundle adjustment approach [26]. Monocular systems [25] make use of a single camera, while stereo systems [5] rely on a calibrated stereo rig. In the monocular setup the translation of the camera may only be estimated up to scale, while in stereo all six motion parameters may be recovered. An additional advantage of the stereo setup is that more information is available at each step, which may be one of the reasons why stereo algorithms perform better.

Chapter 2

The Infinite Visual Odometry Algorithm

2.0.1 Summary of Our Work

In this work we present a novel algorithm for camera motion estimation. The novelty of the algorithm is in camera rotation estimation procedure. We rely on the fact that for scene points that are infinitely far from the camera, the motion of the projected (image) points may be described by an homography (the infinite homography). For distant points this assumption is nearly true. Our algorithm starts by partitioning the scene points into two sets: distant and near-by. Then, camera rotation is estimated from the distant points and, subsequently, the translation is recovered from the near-by points.

With respect to the classification of the visual odometry methods given in the introduction, our work is local, feature based, stereo odometry. We do not use bundle adjustment, however the results of our algorithm may be subsequently improved with some form of bundle adjustment.

The outline of the our method:

1. Feature detection. We use Harris [7] corners.
2. Feature matching. The matching is done both across the stereo pair images as well as previous vs. current pair. We enforce epipolar constraint, chierality and use circle heuristics similar to [5] to reject outliers.
3. Partition the scene points into two sets: distant and near-by.
4. Estimate the rotation of the camera from the distant points.
5. Estimate the translation of the camera from the near-by points.

We choose the work [5] as our baseline (our implementation of their work). The results in the Section 2.3 show that on the KITTI dataset our rotation estimation method outperforms the baseline.

We also build a dataset that contains a number of sequences that may be used to benchmark visual odometry algorithms. We install synchronized stereo pair along a (D)GPS receiver on a car that travels in an urban as well as rural areas. Our goal is to cover some challenging conditions that happen in day-to-day driving, e.g., field of view occlusions, poor lighting conditions and low texture areas.

Note, that budget constraints led us to a simplified system that is capable of producing 3DOF data (i.e., location only) being unable to track vehicle orientation.

2.1 Preliminaries and Notation

2.1.1 Image Point Mapping Related to Camera Motion

Suppose the camera matrices are those of a calibrated stereo rig P and P' with the world origin at the first camera

$$P = K[I \mid 0], \quad P' = K'[R \mid t]. \quad (2.1)$$

Consider the projections of a 3D point $\mathbf{X} = (X, Y, Z, 1)^T$ into the image planes of both views:

$$\mathbf{x} = P\mathbf{X}, \quad \mathbf{x}' = P'\mathbf{X}. \quad (2.2)$$

If the image point is normalized as $\mathbf{x} = (x, y, 1)^T$ then

$$\mathbf{x}Z = P\mathbf{X} = K[I \mid 0]\mathbf{X} = K(X, Y, Z)^T.$$

It follows that $(X, Y, Z)^T = K^{-1}\mathbf{x}Z$, and:

$$\mathbf{x}' = K'[R \mid t](X, Y, Z, 1)^T \quad (2.3)$$

$$= K'R(X, Y, Z)^T + K't \quad (2.4)$$

$$= K'RK^{-1}\mathbf{x}Z + K't. \quad (2.5)$$

We divide both sides by Z to obtain the mapping of an image point \mathbf{x} to image point \mathbf{x}'

$$\mathbf{x}' = K'RK^{-1}\mathbf{x} + K't/Z = H_\infty\mathbf{x} + K't/Z = H_\infty\mathbf{x} + \mathbf{e}'/Z. \quad (2.6)$$

H_∞ is the infinite homography that transfers the points at infinity to the points at infinity. If $R = I$ (e.g., pure translation) the point \mathbf{x} will undergo a motion along a corresponding epipolar line:

$$\mathbf{x}' = \mathbf{x} + K't/Z = \mathbf{x} + \mathbf{e}'/Z. \quad (2.7)$$

If $\mathbf{t} = \mathbf{0}$ the motion of the point may be represented by a homology:

$$\mathbf{x}' = H_\infty\mathbf{x}. \quad (2.8)$$

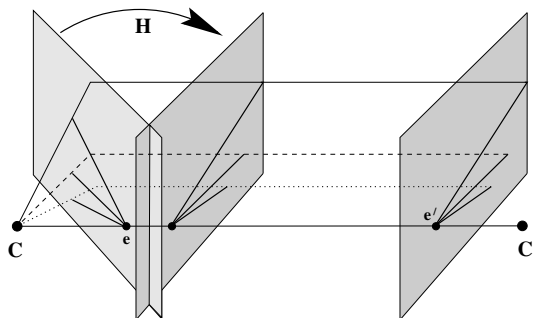


Figure 2.1: (Adapted from [9]) The effect of the camera motion on the image points may be viewed as a two-step process: mapping by a homography H_∞ followed by a motion along the corresponding epipolar lines.

In a general case the mapping of an image point \mathbf{x} into \mathbf{x}' may be viewed as a two step process: transformation by a homology (a specialization of homography which has two equal eigenvalues) H_∞ which simulates a pure rotational motion of the camera followed by an offset along the epipolar line which simulates a pure translational motion of the camera, see Figure 2.1.

2.2 Motion Estimation

Our strategy to attack the problem is to separate it into two smaller sub-problems: rotation estimation and translation estimation. The algorithm relies on the ability to partition the scene points into two sets: the distant and the near-by ones. The distant points are used for rotation estimation while the near-by ones take part in the translation estimation.

First, the stereo algorithm is presented, followed by the monocular one. The main difference is in the translation estimation part. While it is possible to implement a stereo-like algorithm in the monocular setting as well, it suffers from the scale drift. Thus, we propose a different technique.

2.2.1 Stereo

Partitioning the points To partition the points in the stereo case we hard-threshold their Z -coordinates (the threshold is a parameter of the algorithm). The depth of the points was computed by a stereo triangulation.

Rotation Estimation: We use distant points to estimate rotation R (i.e., near-by points do not take part in rotation estimation). As Eq. (2.6) states:

$$\mathbf{x}' = \mathbf{K} \mathbf{R} \mathbf{K}^{-1} \mathbf{x} + \mathbf{K} \mathbf{t} / Z. \quad (2.9)$$

The total motion of the feature point in the image plane may be viewed as a two-step process (the order is not important): transformation by homography

$H_\infty = KRK^{-1}$ followed by displacement along the line defined by the epipole e' and the point $H_\infty \mathbf{x}$. The magnitude of the displacement along the epipolar line depends on the camera translation and the inverse depth of the point, see Figure 2.2.

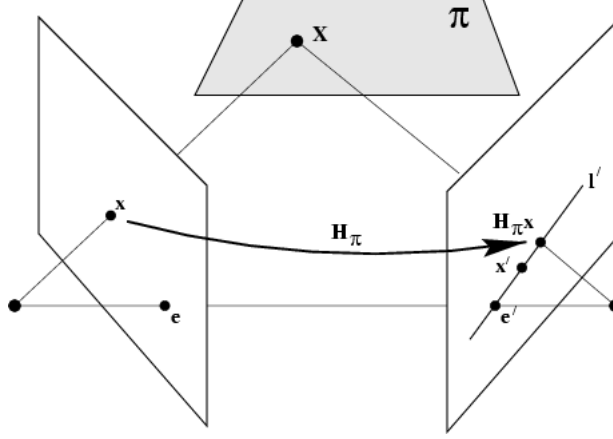


Figure 2.2: (Adapted from [9]) The homography H_π transfers the point onto a corresponding epipolar line. The displacement along the epipolar line depends on the inverse depth of the point and the camera translation magnitude.

Our estimation algorithm consists of initialization and non-linear refinement.

Initialization: to compute the initial estimate of the rotation parameters we assume that for the distant points (s.t., $\|\mathbf{t}\|/Z \ll \|H_\infty \mathbf{x}\|$):

$$\mathbf{x}' \approx KRK^{-1}\mathbf{x}. \quad (2.10)$$

This assumption is justified by the fact that for the distant points the displacement along the epipolar line is small. We multiply both sides of Eq. (2.10) by K^{-1} and denote $\mathbf{u}' = K^{-1}\mathbf{x}'$ and $\mathbf{u} = K^{-1}\mathbf{x}$:

$$\mathbf{u}' = K^{-1}\mathbf{x}' \approx RK^{-1}\mathbf{x} = R\mathbf{u}. \quad (2.11)$$

Since \mathbf{u} and \mathbf{u}' are projective quantities, only their directions are of importance, we normalize them to unit length and denote normalized quantities by $\tilde{\mathbf{u}}$ and $\tilde{\mathbf{u}}'$ respectively. We choose a sample of n points ($n = 3$) and stack them as columns of matrices $\tilde{\mathbf{U}}$ and $\tilde{\mathbf{U}}'$ respectively. We search for R that solves the following minimization problem:

$$\underset{R}{\operatorname{argmin}} \|\tilde{\mathbf{U}}' - R\tilde{\mathbf{U}}\|_2. \quad (2.12)$$

Eq. (2.12) is known as the absolute orientation problem (see e.g., [11]) and its solution provides an initial estimate for the subsequent non-linear optimization problem.

Refinement: The idea of the refinement is this: the residual vector $H_\infty \mathbf{x} - \mathbf{x}'$ may be viewed as a sum of a vector orthogonal to the epipolar line and the vector parallel to it. We search for the camera rotation that ignores the parallel component (we view it as a “legal” bias) while trying to minimize the orthogonal one. We define the point residual as the orthogonal distance to the corresponding epipolar line and minimize the sum of squared residuals for all points. We do so, because, as Eq. (2.6) suggests, after we compensate for a rotation, the point is still allowed to move along the epipolar line.

Consider the objective:

$$R(v, \theta) = \underset{v, \theta}{\operatorname{argmin}} \sum_{i=1}^N r_i^2 \quad \text{s.t.} \quad r_i = n_i \cdot (\mathbf{x}'_i - H_\infty(v, \theta) \mathbf{x}_i) \quad (2.13)$$

where $n_i = (\mathbf{F} \mathbf{x}_i)_\perp$

$H_\infty(v, \theta) = \mathbf{K} \mathbf{R}(v, \theta) \mathbf{K}^{-1}$ is the homography that transforms the image points in Eq. (2.6). It depends on a known camera intrinsic matrices \mathbf{K} and a rotation matrix \mathbf{R} . We choose to parameterize the rotation by an angle θ and an axis v . \mathbf{F} denotes the fundamental matrix that corresponds to two subsequent stereo rig poses and is computed elsewhere. $\mathbf{F} \mathbf{x}_i$ denotes the epipolar line that corresponds to \mathbf{x}_i in the second image and $(\mathbf{F} \mathbf{x}_i)_\perp$ is the normal to this line. We solve the minimization problem (2.13) by means of the Levenberg-Marquardt optimization algorithm.

To make the estimation robust we wrap the initialization procedure into the RANSAC iterations. We choose the strongest consensus estimate and its support set as an input for the solution of the Eq. (2.13).

Translation Estimation (1-point algorithm) To estimate the translation we use only the near-by points. First, we triangulate these points in the previous stereo pair to obtain their 3-D locations, and then iteratively minimize the sum of reprojection errors into the current frame.

The reprojection of point $\mathbf{X} = (X, Y, Z, 1)^T$ into the current left image is given by:

$$\pi^{(l)}(\mathbf{X}; \mathbf{t}) = \mathbf{K} \begin{bmatrix} \mathbf{R} & | & \mathbf{t} \end{bmatrix} \mathbf{X}. \quad (2.14)$$

and the reprojection of the same point into the current right image (b is the baseline of the stereo-rig) is given by:

$$\pi^{(r)}(\mathbf{X}; \mathbf{t}) = \mathbf{K} \begin{bmatrix} \mathbf{R} & | & \mathbf{t} \end{bmatrix} (\mathbf{X} - (b, 0, 0, 0)^T), \quad (2.15)$$

We use the Levenberg-Marquardt algorithm to iteratively minimize the sum of squared reprojection errors (starting from $\mathbf{t} = \mathbf{0}$):

$$\|\mathbf{x}' - \pi^{(l)}(\mathbf{X}; \mathbf{t})\|^2 + \|\mathbf{x}' - \pi^{(r)}(\mathbf{X}; \mathbf{t})\|^2. \quad (2.16)$$

There are three unknown parameters, since $\mathbf{t} = (t_x, t_y, t_z)^T$, thus a single 3-D point provides enough constraints to determine \mathbf{t} .

2.3 Experimental Results

2.3.1 The Choice of Features

We chose to evaluate our algorithm on the KITTI dataset [4], which is a de-facto standard for the visual odometry research works.

Feature Detector/Descriptor: We use Harris [8] corner detector. It is fast, well localized and (most important) Harris corners are abundant in urban scenes we work with. We detect corners in each new image and then match them to obtain putative matches. We tune sensitivity threshold of the detector in such a way, that we are left with about five hundred putative matches after matching and pruning. We extract a square patch of 7×7 pixels centered at the corner point and use this vector as feature descriptor.

We would like to point out that our method may be used with any feature detector that would allow to match features across images. The choice of feature detector should be viewed as a parameter to the algorithm and mainly depends on the images at test.

Feature Matching: We use sum-of-square differences (SSD) of feature descriptors as a metric function when matching features. For each feature we choose a single best match w.r.t. the metric function in the other image. We employ a number of heuristics to prune outliers:

- Reciprocity: features a and b match only if a matches b and b matches a
- Epipolar constraint: we work with calibrated stereo pair. When we match features across images of stereo pair, the search is one-dimensional, i.e., along the horizontal epipolar line. This heuristic is not used when matching features across subsequent frames.
- Chierality (visibility): also used when matching features across stereo pair images. We triangulate the features to obtain the 3-D point and keep the match only if the 3-D point is visible in both cameras.
- Circular match: similar to [5] we keep only those matches that form a circle.

2.3.2 Experimental Results

The Tables 2.1 and 2.2 present the results of the experiments for the KITTI dataset. The columns denote the number of the sequence. The rows denote the algorithm: SS is the baseline, HX is the stereo version of the algorithm. The numbers are the mean error for the corresponding sequence with the last column is the mean error for the dataset. The error computation method is described in [4].

Table 2.1: Rotation errors for the KITTI sequences [deg/m]

	00	01	02	03	04	05	06	07	08	09	10	mean
SS	3.95e-04	1.98e-04	4.11e-04	1.07e-03	8.03e-04	3.49e-04	4.72e-04	2.96e-04	3.69e-04	3.44e-04	4.82e-04	4.71e-04
HX	2.70e-04	1.75e-04	4.10e-04	6.51e-04	6.04e-04	3.95e-04	3.77e-04	2.37e-04	3.23e-04	3.22e-04	5.66e-04	3.93e-04

Table 2.2: Translation errors for the KITTI sequences %

	00	01	02	03	04	05	06	07	08	09	10	mean
SS	4.40e+00	9.25e+00	4.03e+00	1.22e+01	5.06e+00	2.80e+00	4.37e+00	2.21e+00	4.12e+00	5.25e+00	5.60e+00	5.39e+00
HX	3.07e+00	1.08e+01	3.80e+00	7.94e+00	3.82e+00	4.06e+00	3.99e+00	1.67e+00	3.28e+00	3.77e+00	5.65e+00	4.72e+00

Stereo In this set of experiments we ran our algorithm in its stereo mode as described in the Section 2.2.1. Table 2.1 and 2.2 present the rotation and the translation errors respectively in the row HX. The columns marked bold are those that our algorithm outperforms the baseline (on 9 of 11 sequences). The results show that our algorithm improves the rotation results over the benchmark algorithm and successfully competes with it in the translation estimation.

Mono In an additional set of experiments we ran our algorithm in a monocular mode. Monocular motion estimation lacks a scale parameter. In order to compare the results we set the scale of the translation to be that of the stereo algorithm. Feature point selection/partition was done without using any stereo information and the motion estimation was done as explained in the Section ??.

2.4 Conclusions and Discussion

This paper presents a novel visual odometry algorithm. The novelty of the algorithm is in its rotation estimation method. The rotation is estimated by means of the infinite homography. The algorithm may be used both in the stereo and in the monocular setting.

The strengths of the presented algorithm are in its ability to split the motion estimation problem into two smaller problems and to operate directly on the image points instead of on the computed 3-D quantities. Splitting the problem helps because each sub-problem is easier to solve. The ability to partition the points into the distant and the near-by ones is what allows us to separate the rotation and the translation estimation.

The stereo version of the algorithm shows better performance, but the monocular version has the advantage of being a more practical one. Indeed the authors in [4] report that they re-calibrate the cameras before each drive, which is hardly possible in real world installations.

Chapter 3

The Dataset

3.1 Data Acquisition

3.1.1 Overview

Our system consist of a synchronized stereo pair and a GPS receiver synced to the cameras. During the recording process GPS device produces a clock signal of 10 Hz that triggers the cameras. The data is written to PC that is connected to both the cameras and the GPS.

3.1.2 Physical Setup

The cameras We use a pair of 3Mp Ethernet IDS-Imaging cameras UI-5240SE. Hardware synchronization is made by external vendor and uses built-in camera-flash sync option. It allows us to capture images at about the same time (empirically measured time delta is around single digit in ms). Low skew is essential for the stereo algorithms. If the time delta is significant, static scene assumption breaks rendering stereo pair useless. The sync quality depends on the vehicle velocity as well. We built a system suitable for regular driver usecase.

Another rig parameter is the baseline (i.e., distance between the cameras). We tend to be close to KITTI setup in this sense.

Another area of concern is the amount of data such stereo rig produces. We shoot at 10 FPS which makes it about $3 \times 2 \times 10 = 60$ Mb/s. On top of this the system needs to work at a moving vehicle, which rules out non-SSD hard drives.

The cameras are fixed focus and auto-exposed. Auto exposure is critical for outside scenes, especially on a high-contrast sunny day.

The GPS Consumer grade GPS measurements expected accuracy are about 3-5 m. This is not the same order of magnitude as state of the art visual odometry algorithms. On the other hand VO suffers from dead-reckoning (e.g. error accumulation) while GPS error is absolute and constant. In our work we use Differential GPS correction to improve the expected accuracy of the

GPS measurements. An enabler for such correction is the raw GPS data (e.g., pseudoranges, satellites availability, signal strengths, etc.)

After evaluation we chose UBLOX Neo-7P as an inexpensive device that suites our needs. It is able to both produce the raw measurement data and to trigger the cameras. The schematic drawing of the system is presented in Figure 2. UBLOX is connected to the PC through USB and to the cameras by means of electric wire. It triggers the cameras at 10 Hz and produces GPS measurements synced to the images at 1 Hz.

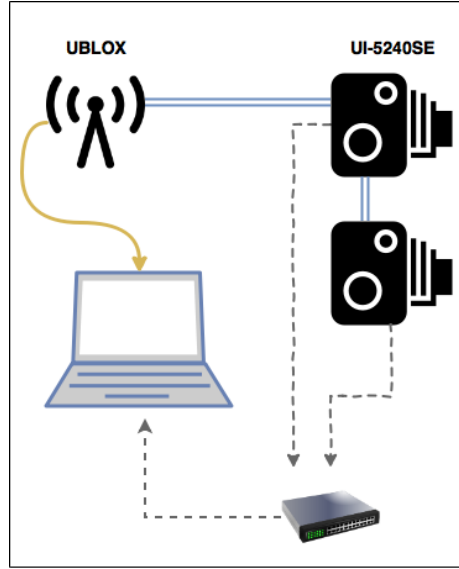


Figure 3.1: Schematic depiction of the system

The Software Our system uses a number of software packages for data recording and post-processing.

We use IDS Imaging software suite to record the image sequences. The sequences recorded as video files at 10 FPS. This software allows us to track the quality of the recording, e.g., if there are lost frames, which improves the quality of the data.

We use goGPS open source positioning software to record GPS data. Our setup is a bit unusual and thus it is rather hard to find out of the box software that is capable to synchronize the GPS and the cameras. We use MATLAB version of the software. We modified out of the box version to communicate with UBLOX device. The modification we made allows us to tell UBLOX to issue a sync pulse to the cameras, which is essential to us.

The experiment day We install the camera rig and the GPS antenna on the roof our lab vehicle. The whole system is powered by the vehicle. At the

beginning of each day we re-calibrate the stereo rig to make sure the calibration data is correct.

We usually record a few minute length video tracks along with the GPS tracks. The data is stored on the PC for further post-processing.

3.2 Post-processing

3.2.1 Image sequences

At this stage we: Compute the calibration data from the calibration board images Split video files to produce image sequences Rectify image sequences.

We use OpenCV calibration toolbox for these tasks. We bundle this code with the project.

3.2.2 GPS Data

First we perform the DGPS correction. The idea of this correction is simple: use a known-location base-station to correct the vehicle position measurement. While the idea is simple, there are a lot of technicalities (e.g., which satellites should be taken into consideration and what the weight of their pseudo-ranges should be in the computation). This may be done in real-time or offline. Budget constraint did not allow us to obtain equipment that is capable of real-time correction, thus we do it offline.

We obtain Israel Mapping Center data for the experiment day. There is a base-station on Technion campus. This is important, since the reliability of the correction decays as the distance to the base-station. It is common wisdom, that the correction is reliable for ≤ 10 Km baselines (e.g. the distance between the vehicle and the base-station). We tried to respect this rule when doing our experiments.

All the correction heavy-lifting is done by the RTKLib software package. RTKLib correction outputs measurement timestamp, WGS84 location and its covariance, along with some additional data. The example is in the figure below:

3.2.3 Missing GPS Data

RTKLib may disregard a certain measurement because of its quality (e.g., if RTKLib is sure that the measurement is wrong). In this case the file will have a gap in time stamps. We provide missing data information as a part of the final dataset.

3.3 VO results evaluation

There is a number of ways to compare VO and DGPS data, but all of the assume that both the GPS track and the VO track are in the same reference frame, which is not the case with our data.

While GPS locations are provided in Earth Centered WGS84 coordinate frame, Visual Odometry data is usually computed relative to the initial position/orientation of the camera. Thus, before the results may be compared to the ground-truth they need to be represented in the same coordinate frame.

To represent both tracks in the same reference frame we use the following procedure: Subtract the location of first GPS measurement from all GPS measurements. This effectively brings the origin of the GPS measurements into the location of the first camera instead of being earth centered. Rotate the VO track to coincide with the GPS data.

To compute the rotation we propose the following procedure: let A be the matrix with VO measurements stacked as its rows, and let B be the matrices with GPS measurements stacked as its rows. We search for the rotation matrix R , s.t. $\arg\min \|RA(:,i) - B(:,i)\|_2$ This optimization problem is called absolute orientation problem in the literature and has an SVD and quaternion based solutions (see e.g., [1] for reference).

	Distance Travelled [m]	Total Number of Measurements	Number of Valid Measurements
02	349.042	104	104
03	2369.983	152	152
04	1507.547	145	145
05	2551.071	227	227
06	4696.472	369	368
07	2565.551	357	289
08	1886.901	235	235
09	1887.352	241	241
10	3022.444	269	269
12	2466.176	379	326
13	3521.839	262	212
14	4856.919	437	361
15	4137.280	332	332
16	3190.590	343	328
17	2032.794	321	320
18	1190.792	166	165
19	1432.072	216	214
total	43664.824	4555	4288

Table 3.1: Dataset details

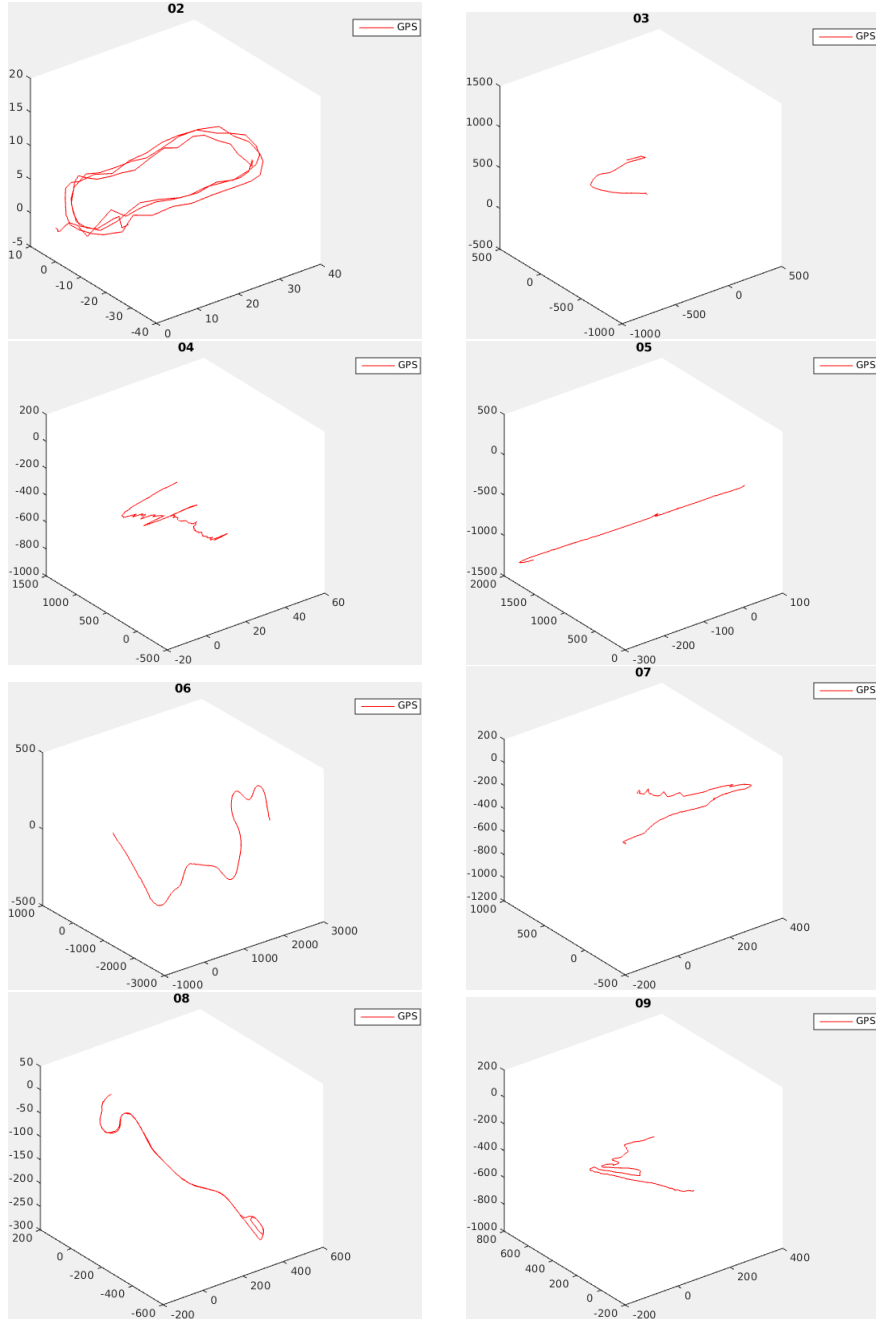


Figure 3.2: Eight sequences, note the difference in axis scales

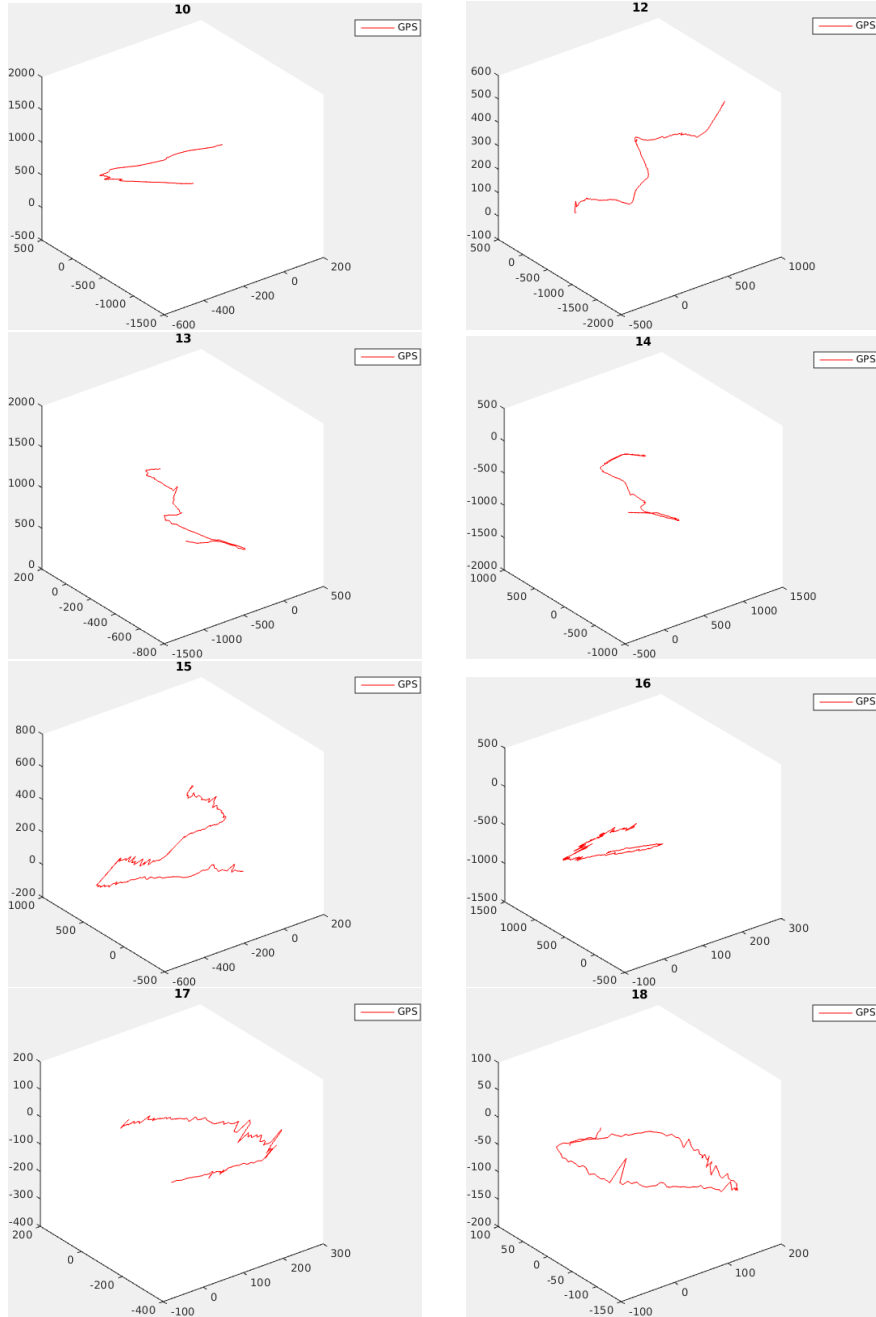


Figure 3.3: Eight more sequences

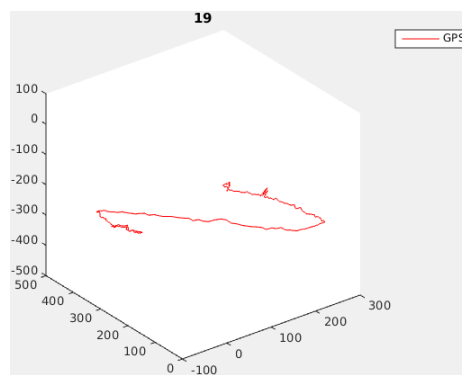


Figure 3.4: One more sequence

Bibliography

- [1] Parallel tracking and mapping for small AR workspaces. *2007 6th IEEE and ACM International Symposium on Mixed and Augmented Reality, ISMAR*, 2007.
- [2] Hernán Badino, Akihiro Yamamoto, and Takeo Kanade. Visual odometry by multi-frame feature integration. *Proceedings of the IEEE International Conference on Computer Vision*, (Cvad 13):222–229, 2013.
- [3] Herbert Bay, Tinne Tuytelaars, and Luc Van Gool. SURF: Speeded up robust features. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 3951 LNCS:404–417, 2006.
- [4] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the KITTI vision benchmark suite. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3354–3361, 2012.
- [5] Andreas Geiger, Julius Ziegler, and Christoph Stiller. StereoScan: Dense 3d reconstruction in real-time. *IEEE Intelligent Vehicles Symposium, Proceedings*, pages 963–968, 2011.
- [6] Natasha Govender. Evaluation of Feature Detection Algorithms for Structure from Motion. *Csir*, 2009.
- [7] C G Harris and J M Pike. 3D Positional Integration from Image Sequences. *Proceedings of the Alvey Vision Conference 1987*, pages 233–236, 1987.
- [8] Chris Harris and Mike Stephens. A Combined Corner and Edge Detector. *Proceedings of the Alvey Vision Conference 1988*, pages 147–151, 1988.
- [9] R. I. Hartley and A. Zisserman. *Multiple View Geometry in Computer Vision*. Cambridge University Press, ISBN: 0521540518, second edition, 2004.
- [10] Johan Hedborg, Pe Forssén, and Michael Felsberg. Fast and accurate structure and motion estimation. *5th International Symposium on Advances in Visual Computing: Part I*, pages 211–222, 2009.

- [11] Berthold K P Horn. Closed-form solution of absolute orientation using unit quaternions. *Journal of the Optical Society of America A*, 4(4):629, 1987.
- [12] Michael Kaess, Ananth Ranganathan, and Frank Dellaert. iSAM: Incremental smoothing and mapping. 24(6):1365–1378, 2008.
- [13] Laurent Kneip, Davide Scaramuzza, and Roland Siegwart. A novel parametrization of the perspective-three-point problem for a direct computation of absolute camera position and orientation. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 2969–2976, 2011.
- [14] Kurt Konolige, Motilal Agrawal, and Joan Solà. Large-scale visual odometry for rough terrain. *Springer Tracts in Advanced Robotics*, 66(STAR):201–212, 2010.
- [15] David G Lowe. Distinctive image features from scale invariant keypoints. *Int’l Journal of Computer Vision*, 60:91–110, 2004.
- [16] Annalisa Milella and Roland Siegwart. Stereo-based ego-motion estimation using pixel tracking and iterative closest point. *Proceedings of the Fourth IEEE International Conference on Computer Vision Systems, ICVS’06*, 2006(Icvs):21, 2006.
- [17] Hans Peter Moravec. Obstacle avoidance and navigation in the real world by a seeing robot rover. *tech. report CMU-RI-TR-80-03*, page 175, 1980.
- [18] R A Newcombe, S J Lovegrove, and A J Davison. {DTAM}: Dense Tracking and Mapping in Real-Time. *Int. Conf. on Computer Vision (ICCV)*, pages 2320–2327, 2011.
- [19] David Nistér. An efficient solution to the five-point relative pose problem. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 26(6):756–770, 2004.
- [20] Clark F. Olson, Larry H. Matthies, Marcel Schoppers, and Mark W. Maimone. Rover navigation using stereo ego-motion. *Robotics and Autonomous Systems*, 43(4):215–229, 2003.
- [21] Mikael Persson, Tommaso Piccini, Michael Felsberg, and Rudolf Mester. Robust stereo visual odometry from monocular techniques. *IEEE Intelligent Vehicles Symposium, Proceedings*, 2015-Augus:686–691, 2015.
- [22] Edward Rosten and Tom Drummond. Machine learning for high-speed corner detection. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 3951 LNCS:430–443, 2006.
- [23] Davide Scaramuzza and Friedrich Fraundorfer. Visual Odometry. (June), 2011.

- [24] Jianbo Shi Jianbo Shi and Carlo Tomasi. Good features to track. *Computer Vision and Pattern Recognition, 1994. Proceedings CVPR '94., 1994 IEEE Computer Society Conference on*, (December):593–600, 1994.
- [25] Shiyu Song and Clark C Guest. Parallel , Real-Time Monocular Visual Odometry.
- [26] Bill Triggs, Philip F. Mclauchlan, Richard I. Hartley, and Andrew W. Fitzgibbon. Bundle Adjustment - a Modern Synthesis. *Vision algorithms: theory and practice. S*, 34099:298–372, 2000.