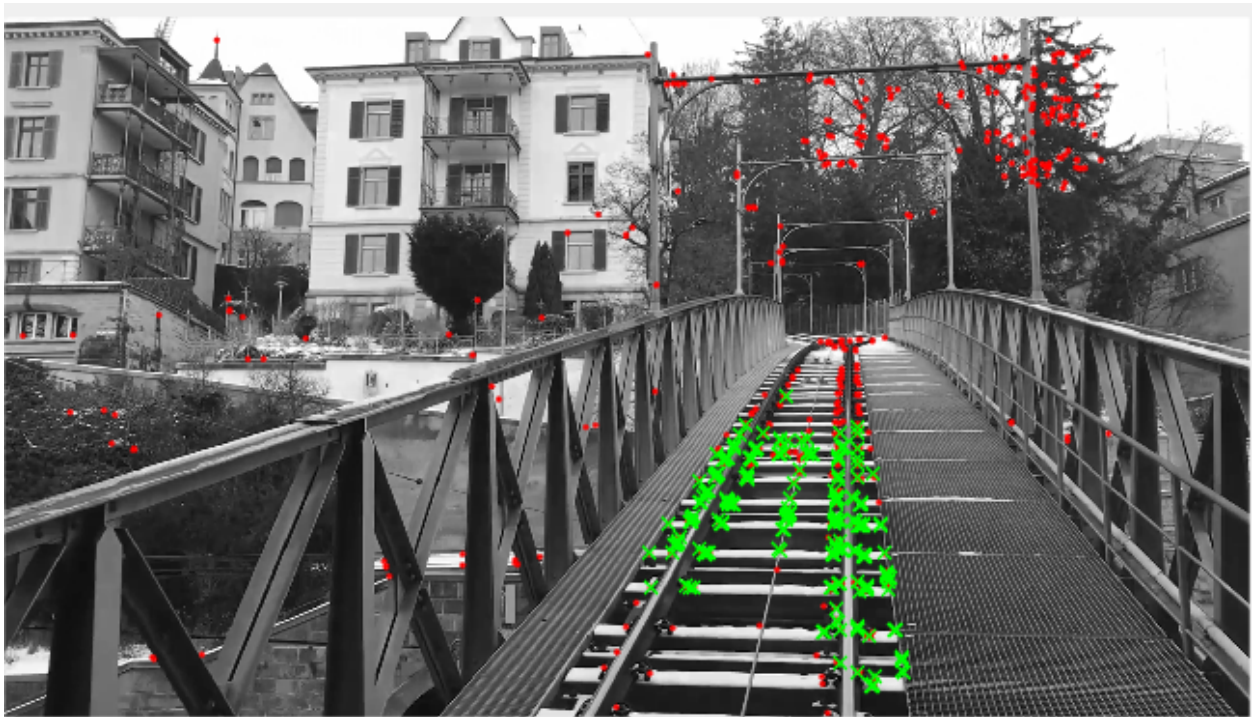


Visual Odometry Pipeline

Pascal Buholzer, Fabio Dubois, Milan Schilling, Miro Voellmy

January 8, 2017



Contents

1	Introduction	3
1.1	Coordinate Frames	3
1.2	Conventions	3
2	Implementation	4
2.1	Pipeline overview	4
2.2	Options and parameters	4
2.2.1	Camera calibration	5
2.3	Initialization	6
2.3.1	getBootstrapFrames()	6
2.3.2	findCorrespondences()	7
2.3.3	eightPointRansac()	7
2.3.4	linearTriangulation()	7
2.3.5	bundleAdjust()	7
2.3.6	applySphericalFilter()	7
2.4	Continuous Operation	8
2.4.1	p3pRansac()	8
2.4.2	findCorrespondences.cont()	8
2.4.3	updateKpTracks()	9
2.4.4	triangulateNewLandmarks()	9
2.4.5	bundleAdjustment()	9
2.4.6	Reinitialization	9
3	Results	11
3.1	Bootstrapping methods	11
3.2	Overall performance	11
3.2.1	Key parameters	11
3.2.2	Run time characteristics summary	12
3.2.3	Comparison to ground truth	12
3.2.4	Speed & Real-Timeness	14
4	Discussion	15
4.1	Bundle Adjustment	15
4.2	Pose estimation algorithm	16
4.3	Reinitialization	16
5	Conclusion	17
5.1	Future work	17

1 Introduction

During this mini project a monocular visual odometry pipeline was developed. This pipeline takes the consecutive gray-scale images of a single digital camera as input. The output of the pipeline is the position of the camera in relation to its initial position for each frame. The pipeline is programmed in such a way that the Markov assumption is valid. This means that the current computation step is only dependent on the previous step to reduce the required computation effort.

1.1 Coordinate Frames

In this mini project the coordinate frames were defined as shown in Fig. 1. The camera coordinates are in a way oriented, that the x-y plane lies parallel to the image plane, while the z-axis is pointing towards the scenery. The world frame however is oriented in such a way that the x-y plane is parallel to the ground and the z-axis is pointing upwards. The origin of the world frame is at the same location as the origin of the first bootstrap image.

Transformation between frames are described by homogenous transformation matrices. T_{AB} maps points from frame B to frame A .

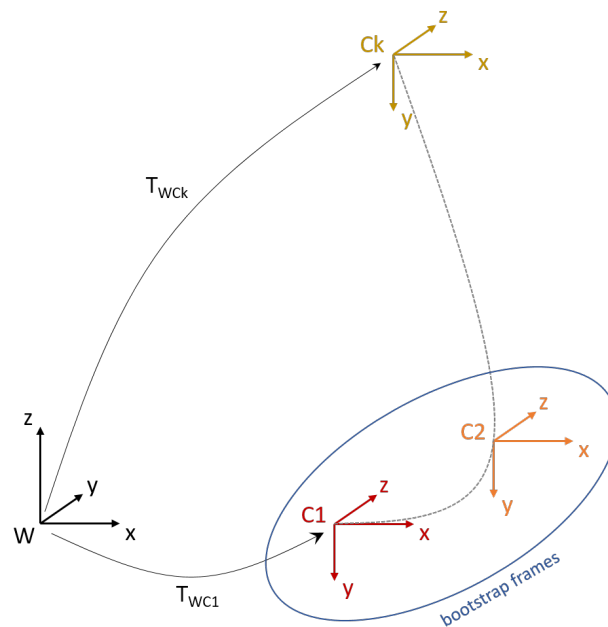


Figure 1: Coordinate Frames

1.2 Conventions

- Index of previous frame: i
- Index of current frame: j
- Index of frame for newly added candidate keypoint: *first*
- Pose difference between previous to current frame: $T_{C_i C_j}$
- $[u/v]$: Pixel coordinates
- Query keypoints: Keypoints newly generated in frame j
- Candidate keypoint: A keypoint without associated landmark
- Harris Matcher: Descriptor matching keypoint tracker (based on Harris features) developed during the lecture.

2 Implementation

This pipeline was developed in MATLAB 2016b. Since the group consisted of four students, a Git repository was used to be able to work on different files simultaneously, and to enable version control.

2.1 Pipeline overview

As shown in Fig. 2 the pipeline consists of three parts:

1. Bootstrapping
2. Initialization
3. Continuous operation

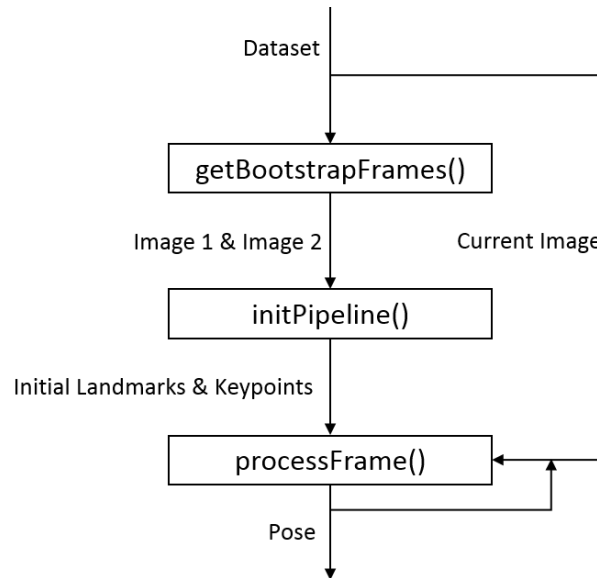


Figure 2: Overview flow chart

2.2 Options and parameters

The pipeline was designed in a modular way. Key algorithms are abstracted into self-contained functions, as described in the pipeline overview. Next to this 'functional programming' approach all the tuning variables (e.g. number of keypoints, bearing angle thresholds, etc.) were centrally aggregated in a parameter struct. For further insight please consult the function `loadParameters.m`.

In order to run the visual odometry two launch procedures were implemented:

- **Debug Mode:** For development and debug mode the `main.m` with default parameters can be executed. Numerous individual plots are displayed with insightful information about matching, inlier rejection and triangulation.
- **Simple GUI:** Out of performance reasons a more compact and user-friendly display of the pipeline output was created with a GUI designed with the Matlab GUIDE application, see figure Fig. 3. Only the most crucial entities, like number of landmarks, are visualized for intuitive understanding.

How to run the GUI Please follow the steps below to run the visual odometry through the GUI environment:

1. adapt dataset paths in `loadParameters.m`
2. type into MATLAB command window: `gui_simple`
3. in *Parameters* panel select dataset to run on and toggle respective radio buttons
4. hit *Run* to trigger the visual odometry

Advanced parameter tuning can be achieved by changing the default parameters in `loadParameters.m`.

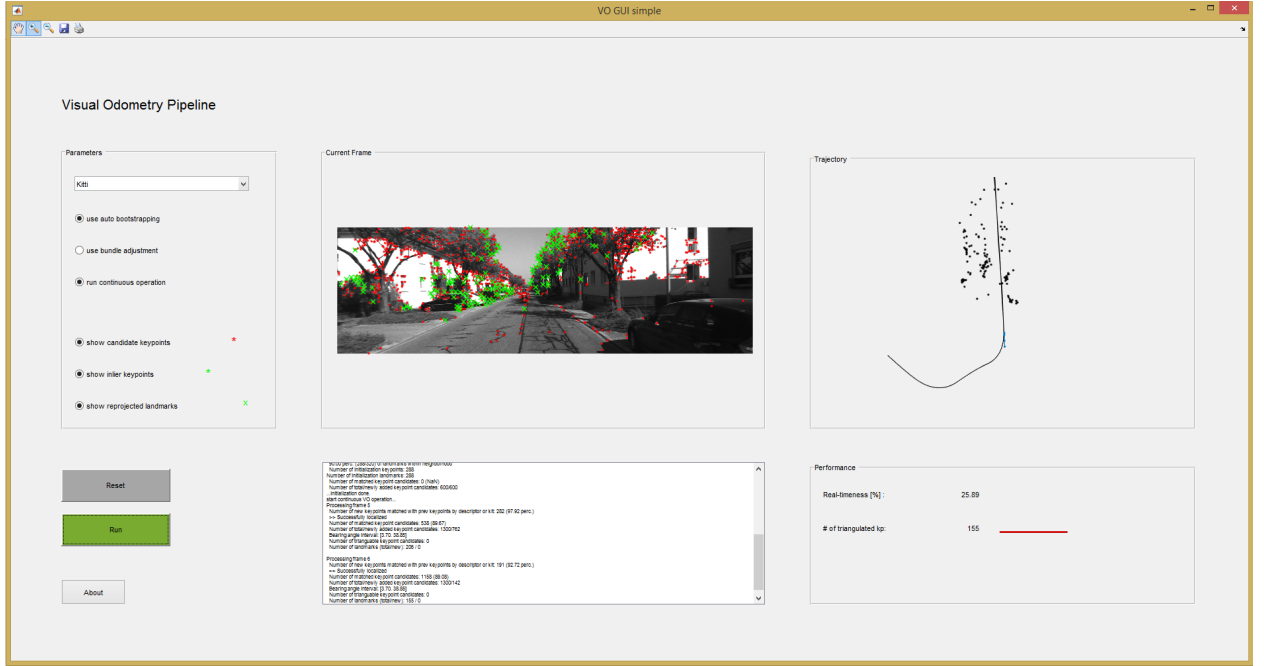


Figure 3: Graphical user interface
candidate keypoints (red ●), inlier keypoints (green ●), reprojected landmarks (green ×)

2.2.1 Camera calibration

Using the Camera Calibration Toolbox for Matlab¹ from Jean-Yves Bouguet (Caltech) the camera of the an iPhone 7 Plus was calibrated. A set of 17 calibration images were extracted from a calibration sequence filming a checker-board pattern from various angles.

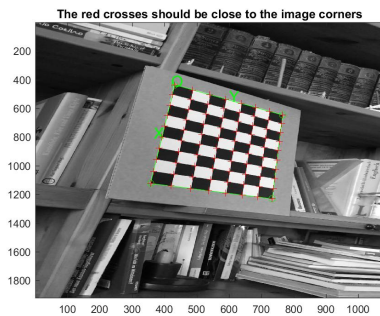


Figure 4: Calibration pattern coordinate system aligned after corner selection

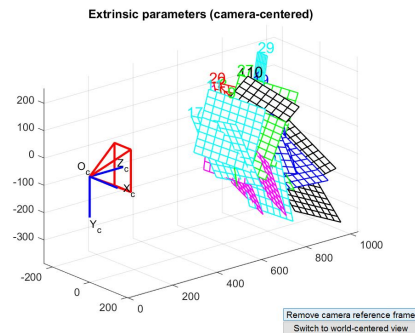


Figure 5: Relative to camera pattern poses

The camera intrinsics parameters together with a fourth-order polynomial approximation of the radial-tangential distortion were calculated based on the assumption of the skew being zero. The exact values and error margins are to be found on the Github repository.

Given these calibration parameters the VO pipeline was also applied on self-generated datasets called *Poly-Up* and *Poly-Down*. Note, that no undistortion was performed, since testing showed it to be negligible.

¹www.vision.caltech.edu/bouguetj/calib_doc/

2.3 Initialization

The first step of the initialization is the bootstrapping which outputs an image pair for further processing. This image pair varies depending on the dataset and is required since the baseline between consecutive images is often too small for accurate landmark triangulation and pose estimation.

Features are then matched across the two bootstrap images and the pose of the second camera is estimated using an 8-point-Ransac. Landmarks are generated using the pose and the matched features. These landmarks and the pose are then refined using bundle adjustment. To end the initialization the landmarks which are unrealistic are discarded.

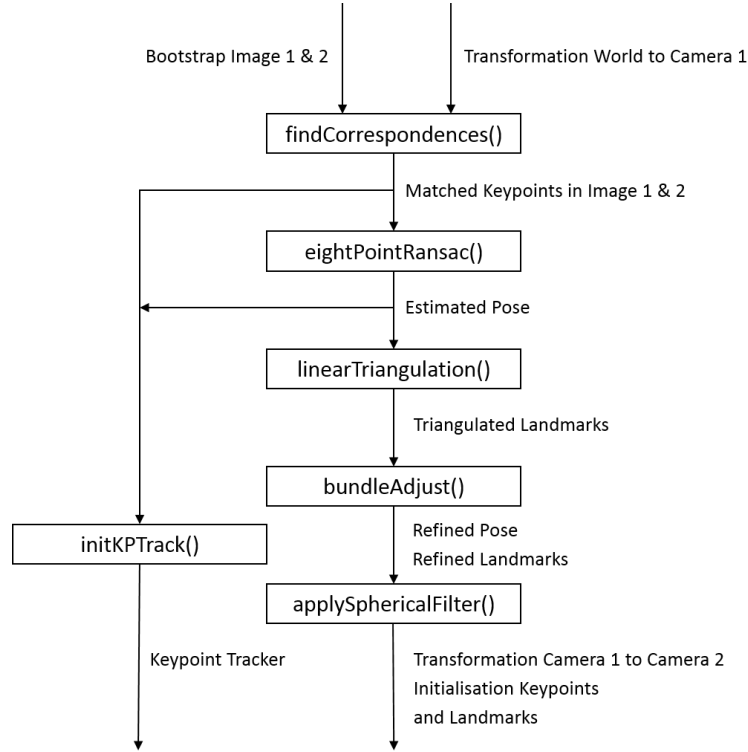


Figure 6: Initialization flow chart

2.3.1 getBootstrapFrames()

In order to pick reasonable image pairs for initialization an automatic procedure was implemented, as a additional degree of 'autonomy' compared to hard-coded index pairs.

Since a good initialization relies heavily on the number of inlier keypoints a hard constrain on their minimum number of `min_num_inlier_kp = 600` was set.

Two approaches were investigated:

- *SSD matching + Baseline/Depth*:
correspondence search via sum of squared differences (SSD) on Harris features & 8-point-Ransac for outlier rejection & baseline/depth-ratio (as proposed in the lecture)

$$\frac{C1_baseline}{C1_av_depth} \geq min_b2dratio = 0.1 \quad (1)$$

- *KLT + bearing angle*:
correspondence search and inlier selection via Kanade-Lukas-Tomasi (KLT) tracking on query Harris corners & minimum average bearing angle

$$av_bearing_angle_deg \geq min_av_angle_deg = 10deg \quad (2)$$

As elaborated in Section 3.1 the second approach outperformed the other and is selected as auto-bootstrapping method.

2.3.2 findCorrespondences()

The same approaches discussed in Section 2.3.1 are used to find the final matching between the bootstrap frames.

2.3.3 eightPointRansac()

Normalized 8-point-Ransac is used to estimate the essential matrix and filter the outliers of the keypoint matching. The error function consists of the shortest distance from the query keypoint to the epipolar line. All keypoints which have an error smaller than a certain threshold are regarded as inliers. The final inlier-set is the one with the most inliers. The essential matrix obtained by the 8-point-Ransac is then decomposed with respect to the obtained translation such that we get the final transformation.

2.3.4 linearTriangulation()

3D-landmarks are generated using linear triangulation of the previously matched keypoints and transformation.

2.3.5 bundleAdjust()

The landmarks and transformation between the first two frames are optimized using bundle adjustment (non-linear least-squares minimization of the reprojection errors). This improved the pose of the initialization as shown in Fig. 7.

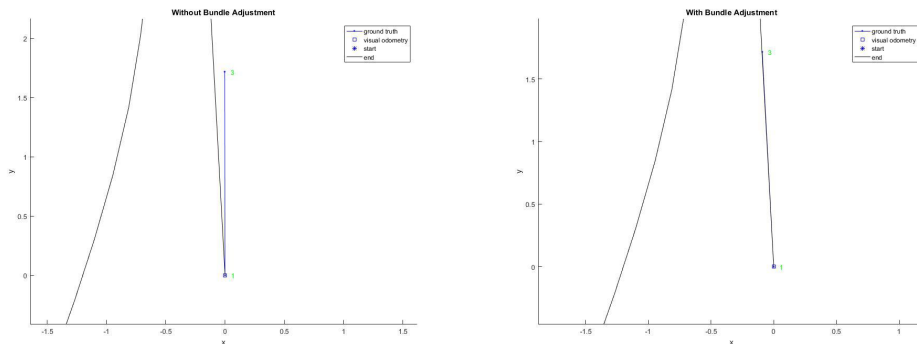


Figure 7: Initialization trajectory before (left) and after (right) bundle-adjustment

2.3.6 applySphericalFilter()

Discards all landmarks which are not within a half-sphere in front of the camera. This removes on one hand landmarks triangulated behind the camera (negative z-components) and also landmarks further away than the filter cutoff radius. With this filtering only landmarks in the neighborhood are kept, which are more reliably triangulated (smaller uncertainty cone). Note, that this absolute thresholding relies on an accurate scale estimation.

2.4 Continuous Operation

Continuous operation of the VO pipeline is implemented in the 'process frame' function. It tracks keypoints with corresponding landmarks over several frames while estimating the pose difference between successive frames. Further, a keypoint tracker finds new candidate keypoints which will become new landmarks if a candidate keypoint was tracked far enough and achieved 'good' trianguability. This ensures to never run out of landmarks and keypoints if the image changes over time. Optionally, a bundle adjustment can be triggered. The routines of the continuous operation shown in Fig. 8 are described below.

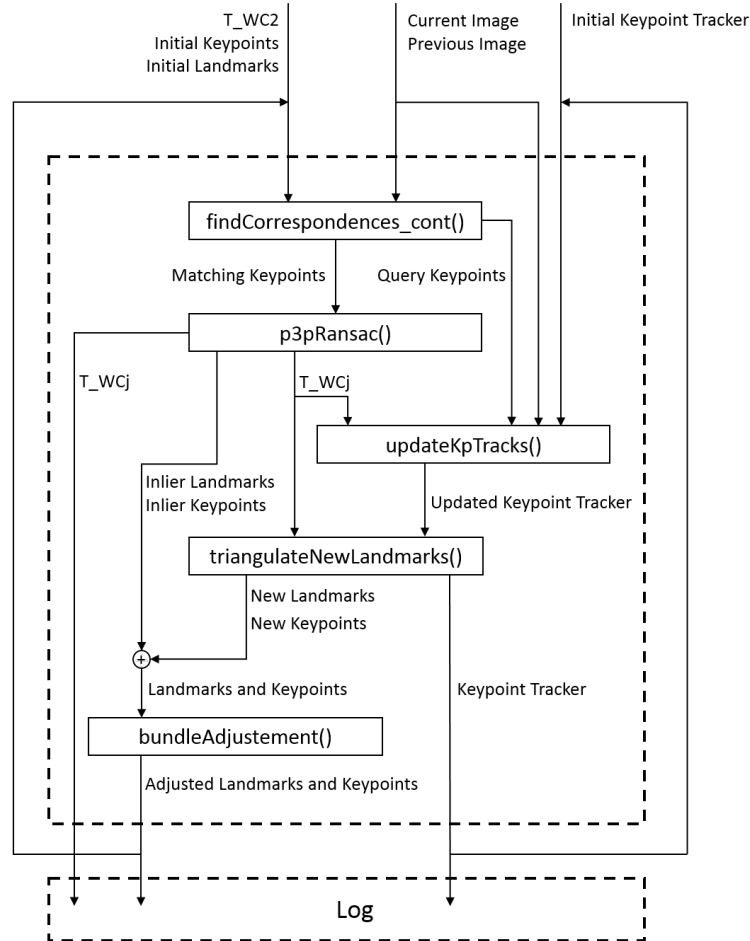


Figure 8: Continuous operation flow chart

2.4.1 p3pRansac()

To estimate the pose difference $T_{C_i C_j}$ from frame i to j we use the p3p-RANSAC algorithm also used in exercise 5. If wished the RANSAC can also use DLT pose estimation. Using these RANSAC algorithm ensures to remove outliers from our landmarks. We don't use DLT refinement after the p3p RANSAC since the best guess from p3p often gave better results.

2.4.2 findCorrespondences_cont()

Tracking keypoints with existing landmarks from frame i to frame j is achieved by the function 'findCorrespondences_cont'. The user can choose whether to use a KLT or Harris matcher. As a by-product, the generated query keypoints are saved to be used by the candidate keypoint tracker in a successive step so they don't have to be generated twice which saves computation time. The number of generated query keypoints is adjustable by a parameter. In case the KLT tracker is active (which does not return query keypoints by default) new query keypoints are generated using Harris features. The amount of newly generated keypoints is the difference of remaining and wished candidate keypoints.

2.4.3 updateKpTracks()

In every frame, candidate keypoints from previous frames are tracked to j -frame. Every candidate keypoint track consists of the following entries: $\{[u/v]_j, [u/v]_{first}, T_{WC_{first}}, nr_trackings\}$. After successive tracking, $[u/v]_j$ as well as the number of successful successive trackings are updated. The user can choose whether to use a KLT or a Harris matcher. In case a candidate keypoint could not be tracked its whole track gets removed from the tracker. If there are less candidate keypoints in the tracker than desired, newly generated keypoints (generated in `find_correspondences_cont`) are added to the tracker. Every newly added keypoint is stored together with the current pose $T_{WC_{first}}$.

2.4.4 triangulateNewLandmarks()

In order to localize correctly (see Section 2.4.1), a sufficient number of landmarks is required. If the number of landmarks in the current frame drops below a certain threshold, the 'triangulateNewLandmarks' function is called to generate new landmarks from the candidate keypoint tracks. To check, whether a candidate keypoint is ready for triangulation the bearing angle between its first and its current observation is calculated. This is done for every keypoint candidate. If the bearing angle of a candidate keypoint exceeds a certain threshold, it is removed from the keypoint tracker and a new landmark gets triangulated with the algorithm developed in exercise 4.

Adaptive bearing angle threshold: The bearing angle threshold is adaptive to the number of remaining landmarks. The higher the number of remaining landmarks the higher the threshold. This ensures generation of new landmarks in case the filter starts to run out of landmarks but improves triangulation results once enough landmarks are in the pipeline. See also parameter 'increase bearing angle threshold' in Table 3.

Filters: Two filters are implemented to discard outliers within the newly generated landmarks (e.g. from image noise):

1. Spherical filter: Already described in Section 2.3.6.
2. Reprojected error filter: Landmarks passed the spherical filter get reprojected into the current frame. If the reprojected point is too far apart from the candidate keypoint coordinates it is discarded.

Landmarks and their corresponding keypoints which passed these two filters are appended to the existing landmarks and keypoints vectors and will be used during the next localization iteration.

2.4.5 bundleAdjustment()

After the new landmarks have been generated bundle adjustment occurs. The bundle adjustment has been implemented in a sliding window fashion. The bundle adjustment occurs at a set frequency over the frames within the sliding window. As depicted in Fig. 9 can both the size of the sliding window and the frequency of the bundle adjustment be adjusted according to the dataset. The pose of the first frame of the sliding window is fixed in order to keep the trajectory continuous.

The only exception regarding the keyframes occurs if the pipeline has crashed. The bundle adjustment is executed directly after the automatic reinitialization and the keyframes will be adjusted accordingly.

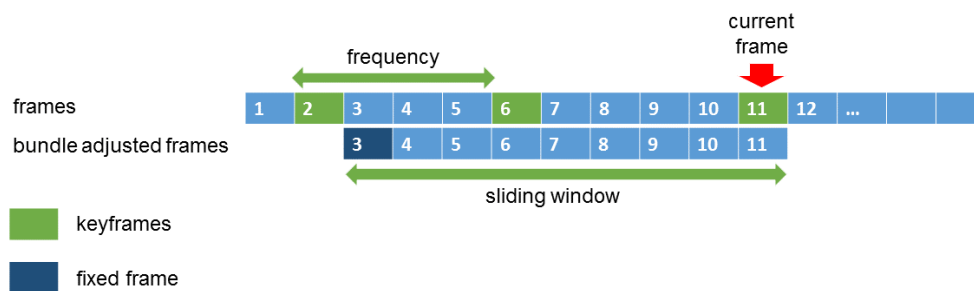


Figure 9: The bundle adjust window moves along with the frames and optimizes all the poses within the window each time the front of the window reaches a keyframe.

2.4.6 Reinitialization

The quality of the transformation matrix estimated by the p3p-RANSAC depends on the number of inlier landmarks. If there are enough inlier landmarks, the pose found by the p3p-RANSAC is a good estimation.

However, if there are too few inlier landmarks, the estimated pose can be considered as too inaccurate for a further use. The fact that too few landmarks were found is a sign that the last few poses weren't good enough. It is therefore a good idea to discard some of the previous poses. The reinitialization discards the last n poses and reinitializes the pipeline with the current image j and the image $j - (n + 1)$. Then the pipeline continues with these newly generated landmarks from pose $T_{WC_{j-n(1)}}$. The keypoint tracker is reinitialized as well and all old (and probably inaccurate) candidate keypoints are discarded.

3 Results

3.1 Bootstrapping methods

The comparison showed a significantly higher performance of the KLT tracker approach, this being orders of magnitude faster and yielding many more feature correspondences also over large distances. Table 1 depicts the index tuples retrieved through the second approach with `min_num_inlier_kp = 600` being required.

approach		Kitti	Malaga	Parking
<i>KLT + Bearing angle</i>	first idx	1	1	1
	second idx	4	6	5
	# keypoints	600	600	600
	angle [deg]	1.37	4.59	4.07

Table 1: Bootstrapping pair indices for different datasets with minimum keypoint inliers

When relaxing the minimum number of inlier keypoints allowed to reach the desired baseline/depth ratio of 10 % for the first and a bearing angle of 10 degrees for the second bootstrapping approach results are shown in Table 2.

approach		Kitti	Malaga	Parking
<i>SSD Harris desc. + baseline/depth ratio</i>	first idx	1	1	1
	second idx	6	10	7
	# keypoints	34	39	54
	ratio [%]	10.1	10.9	11.5
<i>KLT + bearing angle</i>	first idx	-	1	1
	second idx	-	16	17
	# keypoints	-	99	74
	angle [deg]	-	10.4	10.55

Table 2: Bootstrapping pair indices for different datasets and no keypoint number constraint

We observed that the first approach yields far too few keypoints, due to the `matchDescriptor()` method. Furthermore, the pure lateral camera motion of dataset Parking suits well this baseline/depth bootstrapping approach, which confirms our intuition given the 'triangulability condition' described in Section 2.3.1.

Deploying the second approach on the Kitti dataset evidently fails due to the straight camera motion, inhibiting angles ($[0.5, 3.7]$ deg) to reach the desired 10 degrees. As soon as there is a distinctive rotational movement involved, e.g. in Malaga and Parking, a useful bootstrapping is performed.

3.2 Overall performance

This chapter describes the performance of the tuned VO pipeline with different datasets and evaluates some of the key characteristics.

3.2.1 Key parameters

Table 3 shows the parameters that worked best for our implementation and led to the results presented. Tuning the datasets was always a tradeoff between speed, accuracy and robustness. The main parameter, that improved trajectory error a lot was the pixel tolerance of the RANSAC. It had to be set to a small value in order to not improve the trajectory and robustness.

Parameter / Dataset	Kitti	Malaga	Parking	Poly-up	Poly-down
min # of landmarks for reinit	65	90	80	50	50
min bearing angle for triangulation [deg]	1.5	3	3.7	2	2
increase bearing angle threshold [# landmarks]	230	250	230	250	250
RANSAC tolerance [pixel]	2	8	2	3	3
# candidate keypoints in tracker	1500	1300	1300	3500	3500

Table 3: Key parameters chosen for every dataset

3.2.2 Run time characteristics summary

Table 4 summarizes the main run time characteristics of the different datasets with our pipeline.

The **parking dataset** performed fastest and had the best trajectory. Since we could not perform bundle adjustment, there is a slight drift of the estimated trajectory. However, the scene can also be nicely reconstructed in 3D from the point cloud generated (see Fig. 12).

The **Kitti dataset** starts well but suffers a lot from scale drift also due to missing bundle adjustment and thus summing up errors. However, the first few curves are tackled relatively well.

The **Malaga dataset** also suffers from scale drift and reinitialization is needed from time to time. Additionally, the two large curves are not captured very well, the trajectory stays at one point while turning.

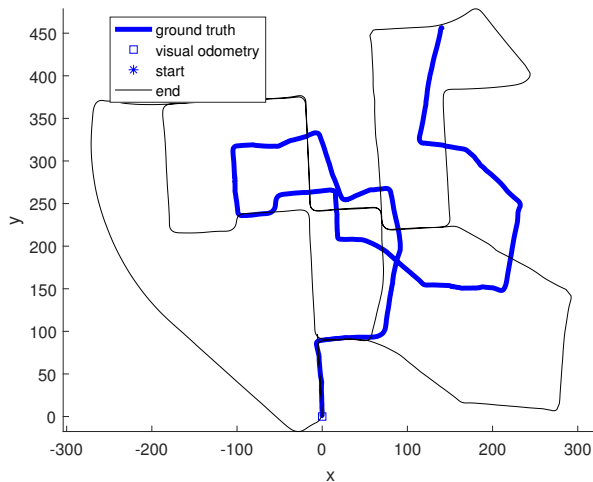
Both **Poly datasets** are quite short but proved to be tricky regarding feature selection. Therefor a high number of candidate keypoints were tracked. Some reinitializations were needed to overcome 'feature-sparse' situation, e.g. midway on the bridge or moments of illumination changes. Still the trajectories look quite smooth and the landmarks represent the ramp and surrounding environment (fence and vegetation) quite well (given an 'arbitrary' initial scale estimate).

Characteristics / Dataset	Kitti	Malaga	Parking	Poly up	Poly down
# reinit required	23	46	2	4	3
# average landmarks	~350	~350	~ 400	~ 250	~ 250
Frame rate [Hz]	1.4	1.6	2.6	1	1
Trajectory quality	+	-	++	+	+

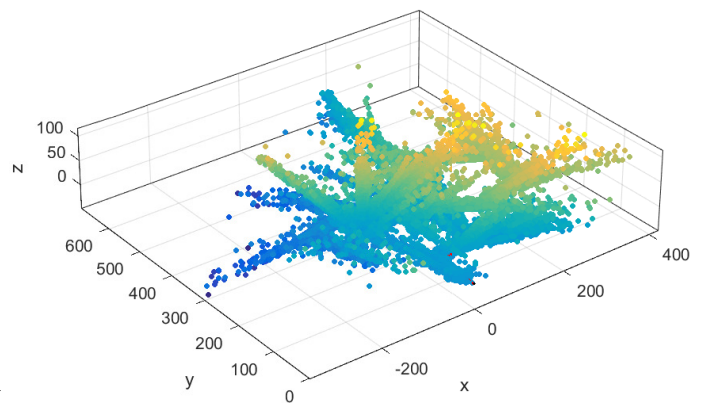
Table 4: Runtime characteristics

3.2.3 Comparison to ground truth

Kitti

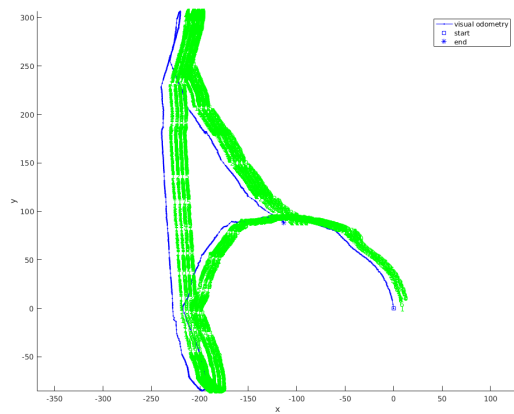


(a) Trajectory vs. ground truth

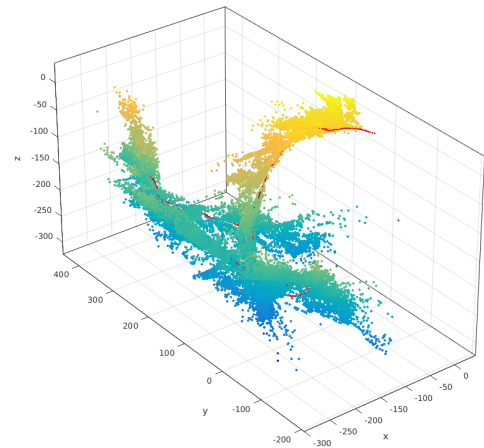


(b) 3D landmarks

Figure 10: Kitti Dataset Results

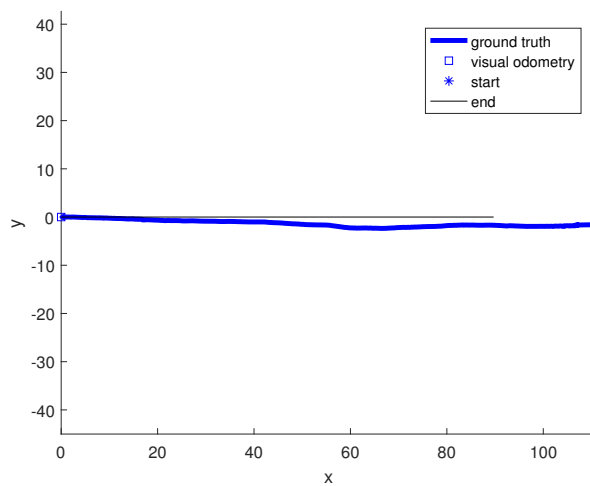
Malaga

(a) Trajectory vs. ground truth

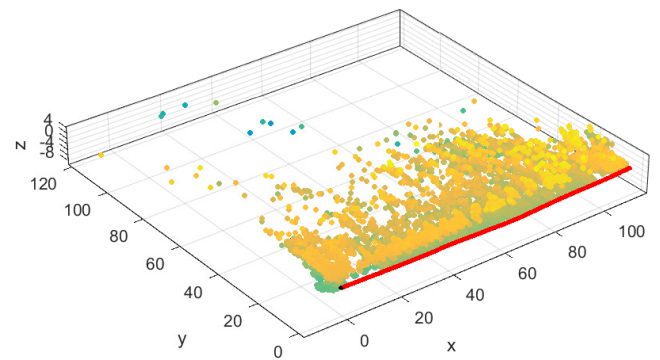


(b) 3D landmarks

Figure 11: Malaga dataset results

Parking

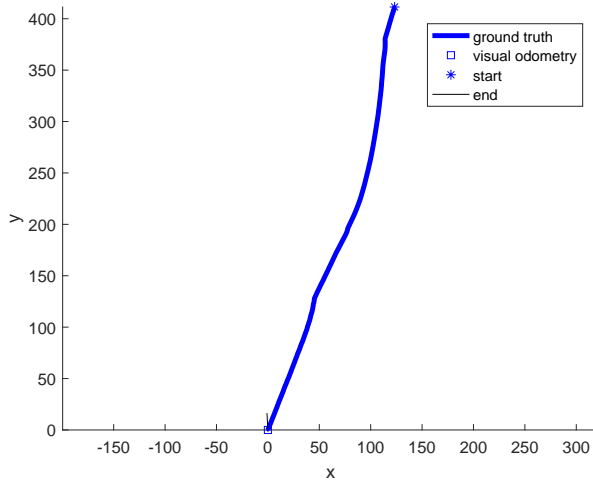
(a) Trajectory vs. ground truth



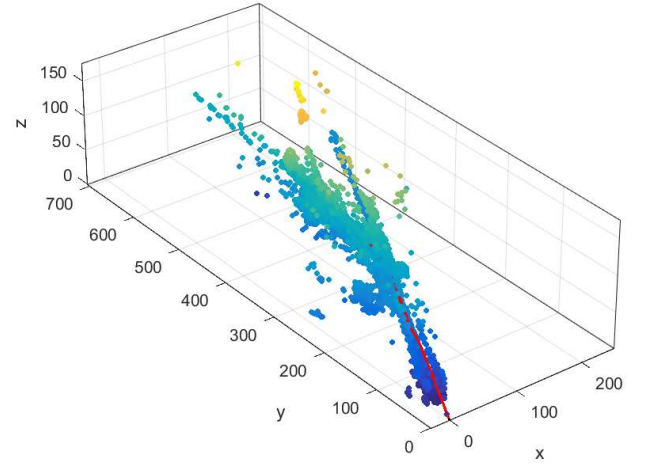
(b) 3D landmarks

Figure 12: Parking dataset results

Poly-Up



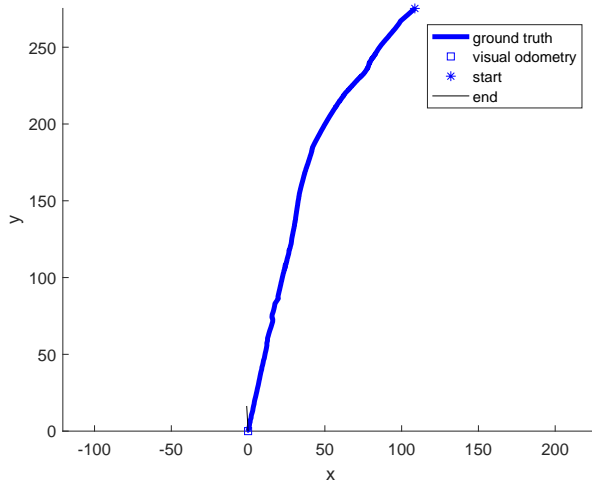
(a) Trajectory



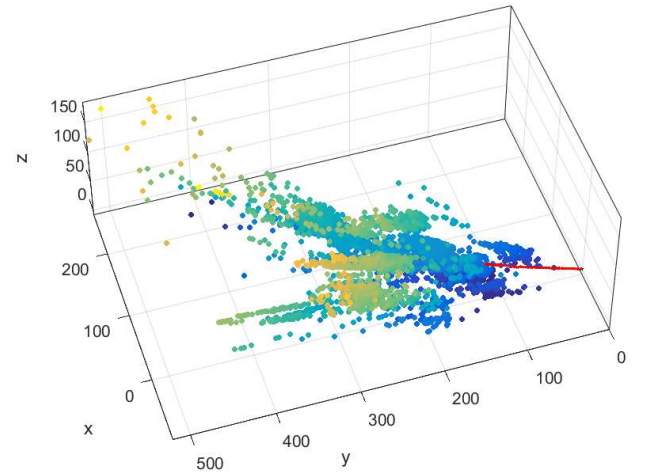
(b) 3D landmarks

Figure 13: Poly-Up dataset results

Poly-Down



(a) Trajectory



(b) 3D landmarks

Figure 14: Poly-Down dataset results

3.2.4 Speed & Real-Timeness

A lot of effort was invested into efficient programming (e.g. avoiding for-loops in Matlab). The parameter with the biggest impact on the frame rate was the number of iterations for the P3P-RANSAC (a classical for-loop). The frame rate depends almost linearly on this parameter. Further the more landmarks and candidate keypoints tracked the slower the pipeline. However, we preferred to generate more landmarks over having a faster pipeline in order to increase robustness and accuracy.

A profiling run with the Matlab profiler tool showed that most of the computational resources is spent in the following three functions in descendant order (self-time): `p3p()`, `selectKeypoints()` and `harris()`. The frame rates were obtained with a Core-i7 (3.5GHz) processor.

4 Discussion

4.1 Bundle Adjustment

While the approach described in Section 2.4.5 was implemented the results did not meet our expectations. Even after thoroughly evaluating and improving the code for dozens of hours no satisfying results were obtained with bundle adjustment enabled.

The plots in Fig. 15 show, that the reprojected landmarks and keypoints before and after bundle adjustment match correctly. Further testing showed that the keypoints are also tracked correctly across the frames. The various parameters such as the number of fixed frames, sliding window and the bundle adjustment frequency were tested across the board. The results were always worse than without bundle adjustment however. One example of strangely optimized landmarks is shown in Section 4.1. The landmarks all correlated into a line such that the trajectory fails.

Since this bug messed up the bundle adjustment, all datasets and results are with bundle adjustment deactivated.

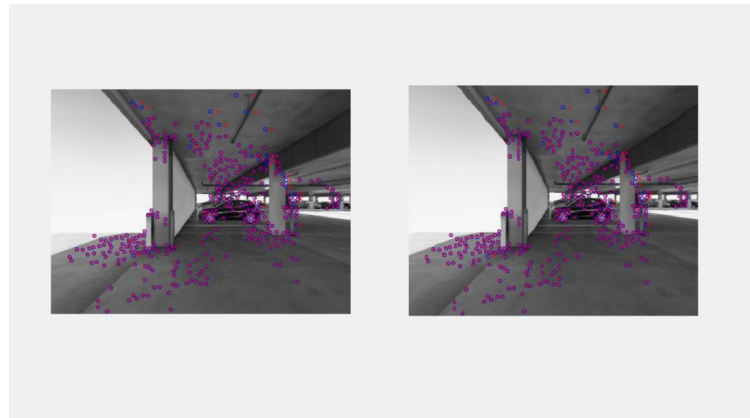


Figure 15: Reprojected landmarks and corresponding keypoints before and after bundle adjustment.

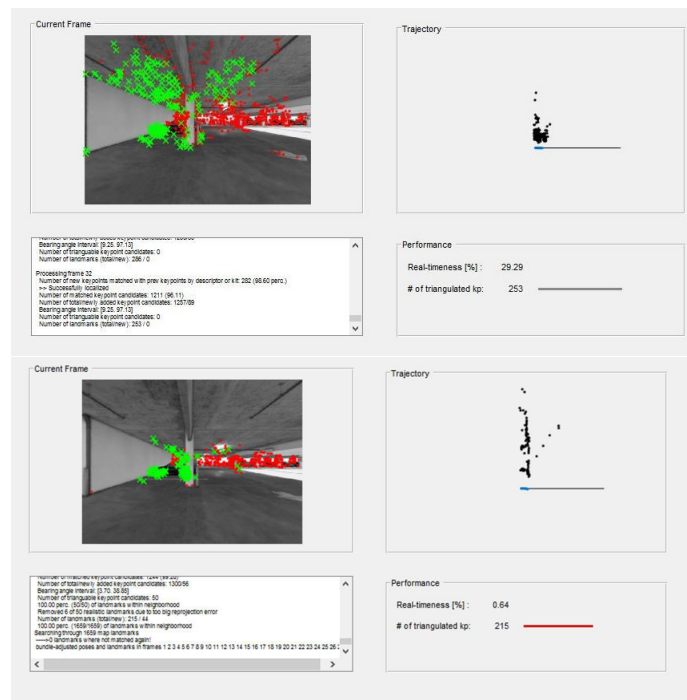
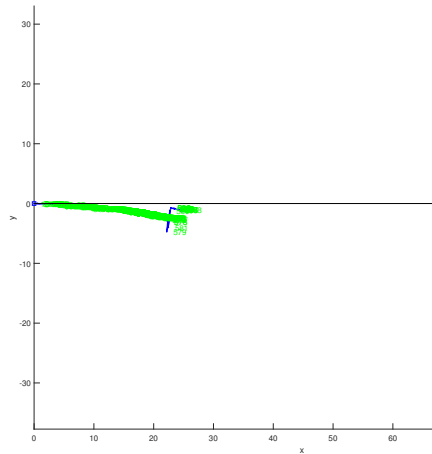


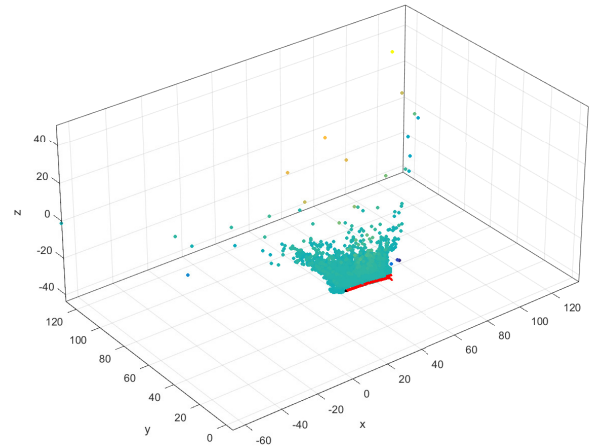
Figure 16: Bundle adjusted landmarks correlate in a line.

4.2 Pose estimation algorithm

We investigated pose estimation with DLT refinement after the P3P-Guess and without refinement (using the last P3P guess from RANSAC). It turned out, that very often the P3P-guess was a lot better and especially more robust than the DLT refinement. Fig. 17 shows the estimated trajectory with DLT refinement. It is clearly worse than the P3P estimate as shown in Fig. 12. That's why we used the last P3P guess in our implementation.



(a) Trajectory vs. ground truth



(b) 3D landmarks

Figure 17: Parking dataset results with DLT refinement

4.3 Reinitialization

The Reinitialization prevents the pipeline from dying of a lack of landmarks. If there are too few landmarks left, the reinitialisation interrupts the keypoint tracker and generates new landmarks. How often this happens depends on the parameters mentioned in Section 3.2.1 as well as on the dataset. The pipeline loses keypoints when areas with a high keypoint density are moving out of the image quickly. We could observe a small buckle of the trajectory every time the reinitialization was done. However, most of the times trajectory improved significantly after reinitialisation.

5 Conclusion

A monocular visual odometry pipeline was developed. Satisfying results were obtained on simple datasets (parking dataset) as well as on more difficult ones (Kitti). Due to the lack of an absolute scale measurement, scale drift can be observed. The pipeline was also successfully tested on a self-recorded and calibrated image series. Various tracking, filtering and recovery approaches including a KLT-tracker, a reinitialisation algorithm, scale normalization and automatic bootstrapping are implemented and can be selected using either a simple GUI or an advanced parameter file to evaluate it's impact. A bundle adjustment was also implemented but due to time restrictions it only worked well for initialization. We expect it would have improved localization even more. The pipeline runs on 1-3 Hz depending on the chosen parameters.

Even tough this mini project proved to be very time intensive we learned a lot about computer vision algorithm and improved our Matlab knowledge. However, there would still be many ideas on how to further improve the pipeline. See also Section 5.1.

5.1 Future work

A list of potential improvements

Possible future work, improvements (loop closure, ...)

- do bootstrapping using both information about the rotation, through difference of bearing angles, and information about translatory motion, from baseline/depth ratio.
- Include bundle adjustment
- Optimize runtime e.g. use a parallel for loop for P3P-RANSAC
- Reinitialize depending on $T_{C_i C_j}$ instead of number of landmarks remaining. An unusually big translation or rotation $T_{C_i C_j}$ is a sign that the pipeline start diverging.
- Loop closure and place recognition (e.g. with bag of words approach).