

Feature-based Visual Odometry Prior for Real-time Semi-dense Stereo SLAM

Nicola Krombach, David Droeschel, Sebastian Houben, Sven Behnke

*Autonomous Intelligent Systems Group, Computer Science Institute VI
University of Bonn, Friedrich-Ebert-Allee 144, 53113 Bonn, Germany*

Abstract

Robust and fast motion estimation and mapping is a key prerequisite for autonomous operation of mobile robots. The goal of performing this task solely on a stereo pair of video cameras is highly demanding and bears conflicting objectives: on one hand, the motion has to be tracked fast and reliably, on the other hand, high-level functions like navigation and obstacle avoidance depend crucially on a complete and accurate environment representation. In this work, we propose a two-layer approach for visual odometry and SLAM with stereo cameras that runs in real-time and combines feature-based matching with semi-dense direct image alignment. Our method initializes semi-dense depth estimation, which is computationally expensive, from motion that is tracked by a fast but robust keypoint-based method. Experiments on public benchmark and proprietary datasets show that our approach is faster than state-of-the-art methods without losing accuracy and yields comparable map building capabilities. Moreover, our approach is shown to handle large inter-frame motion and illumination changes much more robustly than its direct counterparts.

Keywords:

visual simultaneous localization and mapping, visual odometry,
feature-based SLAM, semi-dense SLAM

Email addresses: krombach@ais.uni-bonn.de (Nicola Krombach),
droeschel@ais.uni-bonn.de (David Droeschel), houben@ais.uni-bonn.de (Sebastian Houben), behnke@ais.uni-bonn.de (Sven Behnke)

1. Introduction

A key feature of nearly all mobile robots is a reliable, robust, and fast state estimation that is essential for most high-level operations like autonomous navigation or exploration. Many mobile robots rely on cameras since they are inexpensive and lightweight and can be used for a variety of tasks including visual obstacle detection, 3D scene reconstruction, visual odometry, and even visual simultaneous localization and mapping (SLAM).

Visual odometry (VO) means estimating the egomotion solely from images captured by a monocular or stereo camera system. There are a large variety of VO methods that can be classified into feature-based and direct methods. SLAM broadens this task by also requiring to compute a representation of the robot’s surrounding referred to as map. Most VO and SLAM methods are feature-based and work by detecting keypoints and matching them between frames. In contrast, direct methods estimate the camera motion by minimizing the photometric error over all pixels. Since this minimization consists of aggregating the matching cost over all image pixels, it is computationally more demanding than determining the reprojection error of sparse set of feature points. Hence, direct methods are often computationally more demanding, yet more accurate, than their feature-based counterparts. In this work, we propose a novel approach that combines direct image alignment with sparse feature matching for stereo cameras. By combining both paradigms, we are able to process images with high frame rate and to also track large inter-frame motion while maintaining the accuracy and quality of a direct method. Due to the distinctiveness of the tracked features, our method performs well on datasets with low frame rates, which is often a problem for direct methods as they need sufficient image overlap.

We extend monocular LSD-SLAM [1] to work with a stereo setup and restrict semi-dense matching to key frames for achieving a higher frame rate. In order to estimate the motion between key frames, we employ a feature-based VO method and use the estimated motion as initialization for the direct image alignment. Thus, we restrict the search space for direct image alignment and gain real-time performance. This paper builds upon our recent work [2] where we introduced a VO algorithm deploying both feature-based and semi-direct matching techniques. Here, we expand this approach to a fully-fledged SLAM system.

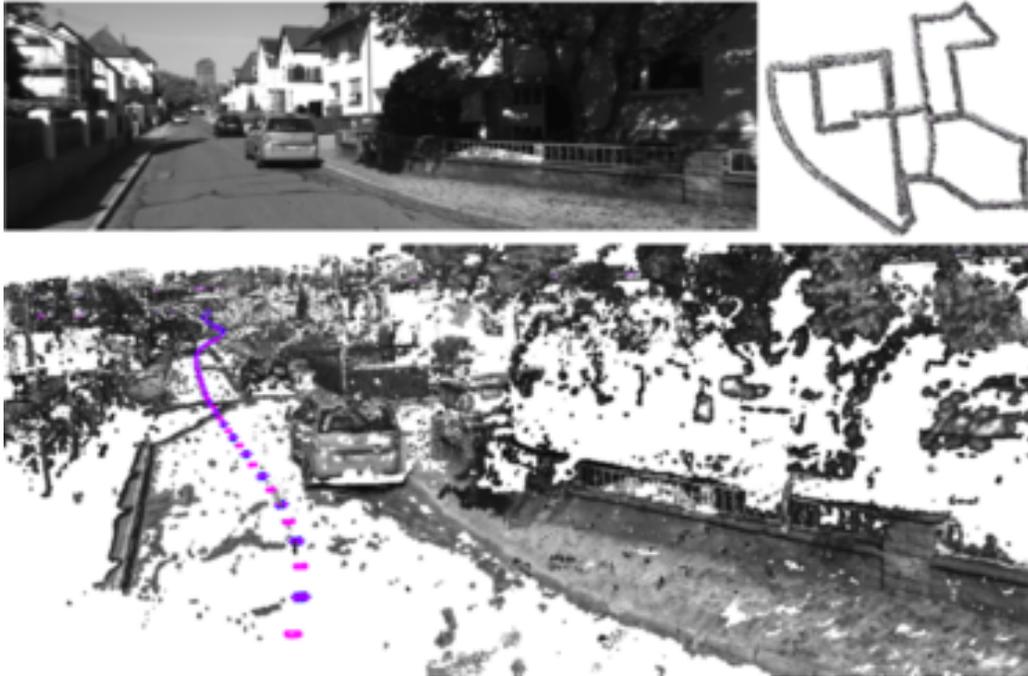


Figure 1: Semi-dense 3D reconstruction of KITTI 00: Top left: Camera image. Bottom: Semi-dense 3D reconstruction with colored camera trajectory (key frames blue, feature-based tracked frames pink). Top right: Bird's eye view of the complete reconstructed scene.

2. Related Work

Visual odometry (VO) and visual SLAM are both vivid areas of research and have seen rapid progress in the past years. Currently, feature-based and direct approaches present two of the dominant paradigms.

2.1. Feature-based Methods

The general pipeline for feature-based methods can be summarized as follows: Image features are detected and either matched between frames or tracked over time. Based on these feature correspondences, the relative motion between two frames is computed. To compensate for drift, many methods make use of pose-graph optimization.

Popular feature-based methods are MonoSLAM [3] and Parallel Tracking and Mapping (PTAM) [4]. PTAM is a widely used feature-based monocular SLAM method, which allows robust state estimation in real-time and has

been successfully used on MAVs with monocular cameras [5]. ORB-SLAM [6] has been proposed as a monocular visual SLAM method that tracks ORB features in real-time and furthermore uses them for local and global bundle adjustment, and candidate retrieval for loop closing. Most monocular methods have a dedicated initialization stage where both the map and the camera movement are estimated at the same time. To this end, the initial movement must bear a certain amount of parallax and can, thus, not be completely arbitrary. On a second note, when using monocular methods, additional sensors are needed to estimate the absolute scale of a scene. In contrast, stereo, depth camera or multi-camera methods [7, 8] have a constant measurement of scale and, hence, do neither suffer from scale drift nor do they need a particular initialization stage. A later extension of ORB-SLAM [9] incorporates stereo and depth cameras by using a dense depth estimation in order to initialize new ORB map points. On the other hand S-PTAM [10, 11] matches sparse features directly and does not rely on a stereo depth images that would need to be computed in a preprocessing step. Likewise, a multi-camera version of ORB-SLAM [12] has been presented that does not rely on dense depth images but triangulates sparse features if stereo pairs are given. Due to their complementary nature, feature-based methods also incorporate readings from an inertial measurement unit (IMU) as high-frequency short-term estimates between frames. Straightforwardly VO and inertial readings are fused in a filter-based approach [13, 14] which is termed *loose coupling*. In particular the work by Forster [15] has allowed for *tight-coupling*, i.e., integrating both IMU readings and visual odometry in a single non-linear cost function. This technique has since then found its way into another variant of ORB-SLAM [16].

In our work, we rely on a well-established and efficient feature-based library for stereo visual odometry [7] which provides a good trade-off between accuracy and runtime.

2.2. Direct Methods

In contrast to feature-based methods, which abstract images into a sparse set of feature points, direct methods use the entire image information in order to minimize the photometric error. In an early work [17], Comport et al. formulate pixel-wise quadrifocal constraints for sparse corresponding stereo matches in two subsequent pairs of images from a stereo setup. If extended to the entire image data, these methods are computationally more intensive

than feature-based methods. First introduced for monocular cameras by Engel et al. [18, 1], LSD-SLAM (Large-scale Semi-Dense SLAM) estimates a pixel-wise inverse depth for a reference frame by means of successive small-baseline stereo estimations. The inverse depth and the according pixel-wise variance is propagated to a new keyframe as soon as the stereo baseline becomes too large. As the pose and point-wise depth estimation is done by minimizing a cost function on image data via gradient descent, the motion must be small or a good initial pose estimate must be given in order to not converge to a local minimum. This is the reason that large inter-frame motion is problematic. The relative poses between the keyframes are asynchronously optimized in a pose graph approach in which two keyframes are connected by a rigid transform with an additional scaling factor. Direct approaches have been extended to stereo and RGB-D cameras. Engel et al. [19] use both fixed-baseline stereo and temporal stereo (as in monocular LSD-SLAM) to refine the depth estimate of the current reference keyframe. Since static stereo is performed initially for every new stereo keyframe, a more reliable depth estimate is available right from the beginning. Hence, the pose and depth refinement become more robust and can deal with larger inter-frame motion. With RGB-D cameras, direct methods are more straightforward as a point-wise depth estimate with constant variance is given in every frame. Steckler and Behnke [20] transform the depth image into a coarser representation, named a surfel map, for aggregation and track the camera motion with ORB features. They later utilize this approach for dense image registration and combine it with a sparse feature matcher in order to compute visual odometry [21]. Dense direct methods often need to use GPUs to achieve real-time performance [22, 23]. By using only pixels with sufficient gradient, LSD-SLAM [1] reduces the computational demand and real-time semi-dense SLAM becomes possible with a strong CPU. The extension to stereo cameras [19] uses both fixed-baseline stereo depth and temporal multi-view stereo in order to estimate a semi-dense environment representation. Recently, a method for directly optimizing the depth of sparse feature points for visual odometry, DSO (Direct Sparse Odometry), has been proposed by Engel et al. [24]. Building upon the same optimization scheme like LSD-SLAM, they optimize for all parameters (including the depth of numerous sparsely chosen image points) for a sliding window of a few keyframes. In order to increase robustness, the cameras must be calibrated photometrically and exposure times have to be taken into account. Schöps et al. [25] use a visual-inertial odometry approach to compute a short time horizon of camera poses and

obtain a dense stereo estimation by plane sweeping multiple images[26]. The latter method can be scheduled in parallel and achieves real-time applicability by use of a GPU.

2.3. Hybrid Methods

Regarding the reconstruction of the environment, direct methods have the advantage of estimating a dense map while feature-based methods can only rely on the sparse features that have been tracked. Dense direct methods are computationally demanding and are often executed as a final step for estimating a globally consistent dense map after pose tracking with sparse interest-points succeeded. To speed up global optimization, already tracked sparse feature-points can be used as initialization for dense mapping [27]. The semi-direct method by Forster et al. [28] uses direct motion estimation for initial feature extraction and continues by using only these features. A novel release includes edgelets as features, encompasses IMU readings, and yields a significant speedup [29]. A recent paper by Piazza et al. [30] presents a real-time capable algorithm to compute and update a 3D manifold mesh on a CPU. This allows for deriving a dense 3D map from a set of sparse points as provided by any of the above SLAM systems. A combination of a feature-based and a direct method has been presented by Younes et al. [31] as feature-assisted direct monocular odometry. They present a VO method that is based on DSO[24] but uses feature-based tracking when optimization yields little relative improvement. In contrast, we always continuously combine feature-based and semi-dense direct tracking over time, taking advantage of the fast tracking from the feature-based method and the accurate alignment of image gradients from direct methods. The feature-based tracking result is immediately fed to the direct tracking at runtime as an initial guess.

3. Method

Our method is mainly based on the monocular version of LSD-SLAM that we extended to work with stereo cameras. By using stereo cameras instead of a single monocular camera, the absolute scale of the scene becomes observable, eliminating scale ambiguity and the need for additional sensors, e.g., inertial measurement units.

To ensure a high frame rate, we restrict the semi-dense direct alignment to key frames only and estimate the motion for all other frames by the feature-based method LIBVISO2 [7]. This motion estimate is used as initial estimate

for direct alignment of key frames. The semi-dense environment mapping runs in a parallel thread.

3.1. Notation

We follow the notation of Engel et al. [1]. The monochrome stereo images captured at time i are denoted with $I_i^{l/r} : \Omega \subset \mathbb{R}^2 \rightarrow \mathbb{R}$, with image domain Ω . Each key frame $KF_i = \{I_i^l, I_i^r, D_i, V_i\}$ consists of the left and right stereo images $I_i^{l/r}$, the semi-dense inverse depth map $D_i : \Omega_{D_i} \rightarrow \mathbb{R}^+$, and the corresponding pixel-wise variance map $V_i : \Omega_{D_i} \rightarrow \mathbb{R}^+$. The inverse of the depth z of a pixel is denoted as $d = z^{-1}$. Camera motions are represented as twist coordinates $\boldsymbol{\xi} \in \mathfrak{se}(3)$ with corresponding transformation $\mathbf{T}_{\boldsymbol{\xi}} \in SE(3)$. A 3D point $\mathbf{p} = (p_x, p_y, p_z)^T$ is projected into image coordinates $\mathbf{u} = (u_x, u_y, 1)^T$ by the projection function $\pi(\mathbf{p}) := \mathbf{K} (p_x/p_z, p_y/p_z, 1)^T$ with intrinsic camera matrix \mathbf{K} . Thus, the inverse projection function $\pi^{-1}(\mathbf{u}, d)$ maps a pixel with corresponding inverse depth to a 3D point $\mathbf{p} = \pi^{-1}(\mathbf{u}, d) := (d^{-1}\mathbf{K}^{-1}\mathbf{u})^T$.

3.2. LSD-SLAM

The processing pipeline of LSD-SLAM [1] consists of the three main components: Tracking, depth map estimation, and global map optimization.

Tracking, i.e., frame-wise relative pose estimation, is based on maximizing photo-consistency and thus minimizing the photometric error between the current frame and the most recent key frame using Gauss-Newton optimization:

$$E(\boldsymbol{\xi}) := I_{KF}(\pi(\mathbf{p})) - I(\pi(\mathbf{T}_{\boldsymbol{\xi}} \mathbf{p})) \quad , \forall p \in \Omega \quad (1)$$

where \mathbf{p} is warped from I to I_{KF} by $\boldsymbol{\xi}$. New frames are tracked towards a key frame and the rigid body motion of the camera $\boldsymbol{\xi} \in \mathfrak{se}(3)$ is estimated.

In the depth map estimation, tracked frames are then used to refine the existing depth map of the key frame by many short-baseline stereo comparisons. Given the transformation between a tracked frame and the key frame, that has been estimated prior in the tracking, the epipolar lines are calculated. Afterwards, for each pixel with sufficient gradient its depth hypothesis is updated with stereo measurements. The depth is calculated by finding the best matching point along the epipolar line, that is the point which minimizes the SAD error measured over five equidistant points along the epipolar line. Given the estimated depth the depth map of the most recent key frame is then refined by either creating new depth hypotheses or improving existing ones. New key frames are created when the distance exceeds a certain

threshold and are initialized by propagating depth of the previous key frame towards the new frame. Once a key frame is replaced, it is added to the pose-graph for further refinement and loop closing.

3.3. LIBVISO2

LIBVISO2 [7] is a fast feature-based VO library for monocular and stereo cameras. Similar to other feature-based methods, it consists of feature matching over subsequent frames and egomotion estimation by minimizing the reprojection error. Features are extracted by filtering the images with a corner and blob mask and performing non-maximum and non-minimum suppression on the filtered images. Starting from all feature detections in the current left image, candidates are matched in a circular fashion over the previous left image, the previous right image, the current right image, and back to the current left image. If the first and last features of such a circle match differ, the match is rejected. Based on all found matches, the egomotion is then estimated by minimizing the reprojection error using Gauss-Newton and outliers are removed using RANSAC.

3.4. Semi-dense Alignment of Stereo Key Frames

We build upon the open source release of monocular LSD-SLAM and extend it with stereo functionality. In contrast to monocular visual odometry, stereo allows to compute absolute depth maps and, thus, does not suffer from scale drift. By extending LSD-SLAM to stereo, we combine the existing depth map computation over time with instant stereo depth from the current image pair. While monocular LSD-SLAM uses a random initialization and has to bootstrap over the first frames, we take advantage of using stereo cameras and initialize our method with absolute depth values. We use ELAS [32] to compute the depth map of the initial key frame. The following key frames are registered with their previous key frame by minimizing the photometric error of their left reference frames as well as the depth error. While in the monocular case, absolute depth is not observable; with stereo cameras absolute depth is observable for every incoming stereo image pair. This allows us to minimize the depth error in addition to the photometric error. Hence, for direct tracking with stereo, we extend the minimization of the photometric residual r_p to take the depth residual r_d into account:

$$\begin{aligned} r_p(\mathbf{p}, \boldsymbol{\xi}) &= \|I_{KF_i}(\pi(\mathbf{p})) - I_j(\pi(\mathbf{T}_\xi \mathbf{p}))\|, \\ r_d(\mathbf{p}, \boldsymbol{\xi}) &= \|D_{KF_i}(\pi(\mathbf{p})) - D_{stereo_j}(\pi(\mathbf{T}_\xi \mathbf{p}))\|, \end{aligned} \tag{2}$$

where $\boldsymbol{\xi}$ is the camera motion from the i -th key frame to the new j -th frame and D_{stereo_j} is the initial instant stereo depth map of the j -th frame. The minimization is performed using a weighted least squares formulation and solved with the Gauss-Newton method. The residual is formulated as stacked residual \mathbf{r} and is weighted with a 2×2 weight matrix \mathbf{W} :

$$\begin{aligned} \mathbf{r}(\boldsymbol{\xi}) &= \mathbf{W}(\boldsymbol{\xi}) \sum_{\mathbf{p} \in \Omega_{D_i}} \begin{pmatrix} h(r_p(\mathbf{p}, \boldsymbol{\xi})) \\ h(r_d(\mathbf{p}, \boldsymbol{\xi})) \end{pmatrix}; \\ \mathbf{W}(\boldsymbol{\xi}) &= \sum_{\mathbf{p} \in \Omega_{D_i}} \begin{pmatrix} w_p(r_p(\mathbf{p}, \boldsymbol{\xi})) & 0 \\ 0 & w_d(r_d(\mathbf{p}, \boldsymbol{\xi})) \end{pmatrix}, \end{aligned} \quad (3)$$

where both residuals are weighted with the Huber norm denoted as $h(\cdot)$.

3.5. Hybrid Odometry Estimation

Our idea is to take advantage of the different strengths of both approaches and, thereby, combine fast feature matching with precise semi-dense image alignment for efficient and reliable state estimation. The modular structure of our approach is illustrated in Figure 2.

We initialize the first key frame with a dense depth map computed by ELAS (Figure 2 (1)). Subsequent frames are then tracked towards the key frame incrementally using feature-based LIBVISO2 (Figure 2 (2)). The relative poses of the tracked frames are concatenated and form the relative pose of the camera to the key frame:

$$\boldsymbol{\xi}_{feat} = \boldsymbol{\xi}_{i_n} \circ \boldsymbol{\xi}_{i_{n-1}} \circ \cdots \circ \boldsymbol{\xi}_{i_0} . \quad (4)$$

The current absolute pose of the camera at step j and key frame i can be retrieved by:

$$\boldsymbol{\xi}_{ij} = \boldsymbol{\xi}_{KF_i} \circ \boldsymbol{\xi}_{i_{j-1}} . \quad (5)$$

We perform feature-based odometry as long as the motion is sufficiently small. As soon as the motion exceeds the motion threshold

$$\epsilon_{motion} = \frac{1}{n} \sum_{k=1}^n \sqrt{(\mathbf{u}_i^k - \mathbf{u}_{i-1}^k)^2}, \quad (6)$$

we perform direct registration again (Figure 2 (3)) and the previous key frame is replaced with the new frame, where n is the number of matched feature

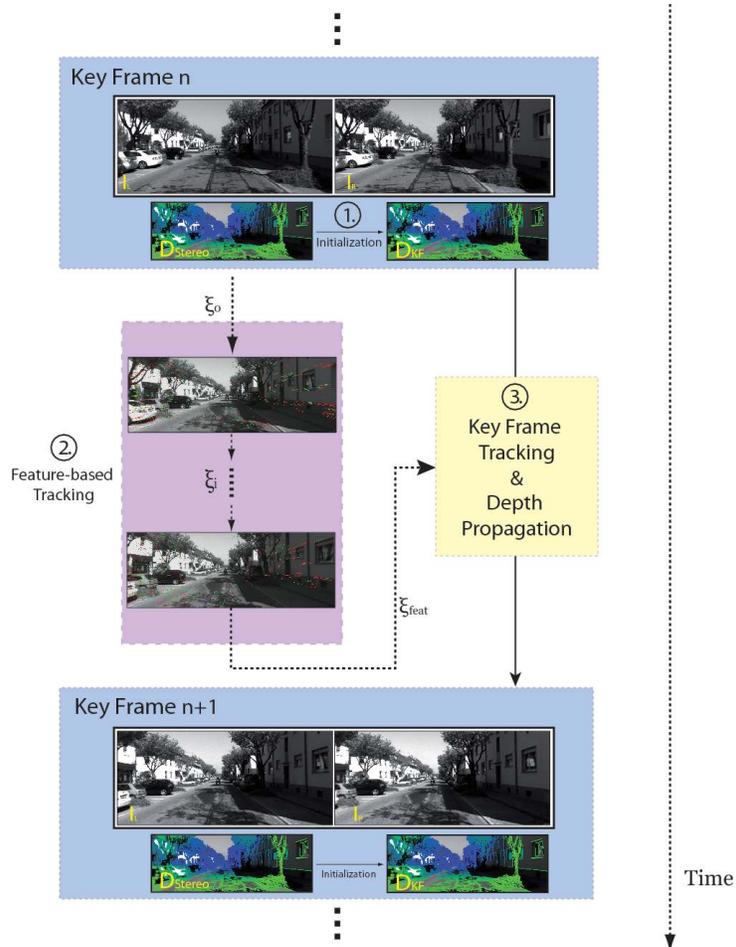


Figure 2: Overview of our combined semi-direct approach. While direct tracking is only performed on key frames, feature-based tracking is performed for frames in between. The output of the feature-based odometry serves as prior for the direct tracking.

points and (\mathbf{u}_i^k) and (\mathbf{u}_{i-1}^k) are corresponding feature matches between the current and the previous image.

The motion ξ_{feat} serves as initial estimate for the direct registration of the new frame towards the key frame:

$$\xi_{KF_{i+1}} = \xi_{KF_i} \circ \xi_{feat} . \quad (7)$$

This allows us to track larger motions faster and more robustly. The depth map of a new key frame is initialized by instant stereo correspondences and

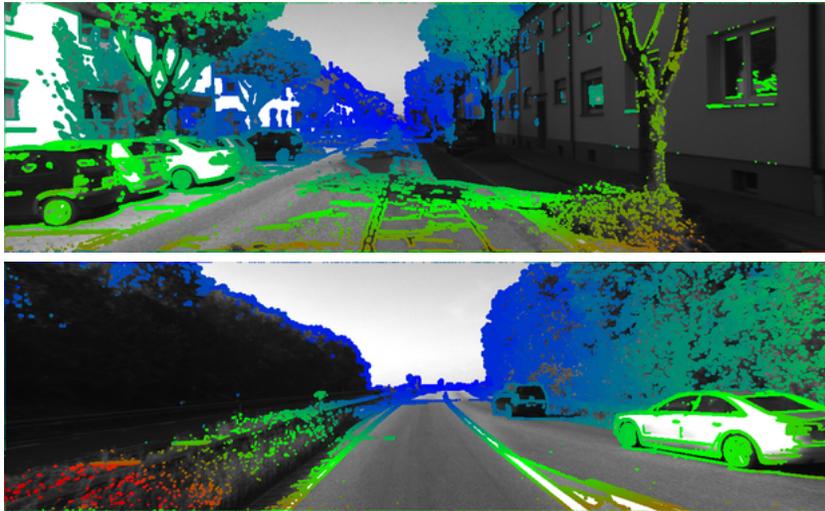


Figure 3: Computed semi-dense depth maps for KITTI datasets (sequences 00 and 01). Color depicts distance to the sensor.

then fused with the previous depth map by propagation as described in the next section. Once a new key frame is initialized, we start feature-based matching again.

3.6. Map Update

The depth map of each key frame is updated with instant stereo measurements as well as with propagated depth from the previous key frame. If a new key frame is created, the depth map is computed by instant stereo from the left and right images. For reasons of runtime, we use a simple but fast block matching along epipolar lines instead of ELAS [32] which is used for the first keyframe only. Corresponding pixels are found by minimizing the sum of absolute distances (SAD) error over a 15×15 pixel window. After initializing the depth map with stereo measurements, the depth estimates are refined by propagating depth hypotheses of the previous depth map to the new frame:

$$\mathbf{p}_{new}(\mathbf{p}) = \mathbf{R}_{C,KF} \mathbf{p} + \mathbf{t}_{C,KF}, \quad (8)$$

where \mathbf{p} is the 3D point in the previous key frame. The rotation $\mathbf{R}_{C,KF}$ and translation $\mathbf{t}_{C,KF}$ describe the coordinate transformation from the key frame coordinate system KF to the candidate coordinate system C . If the residual between the instant and propagated depth is high, the depth value with

smaller variance is chosen. Otherwise both estimates— d_{stereo} and d_{prop} —are fused to a new depth estimate d_{new} as a variance-weighted sum:

$$d_{new} = (1 - \omega) d_{stereo} + \omega d_{prop} . \quad (9)$$

The variance ω for each depth hypothesis is determined as described by Engel et al. [33].

For fish eye lenses we extend this weighting scheme: as fish eye lenses suffer from distortion at the image borders, we further increase the variance of a pixels depth hypothesis depending on the distance $r(u, v)$ to the optical center (c_x, c_y) :

$$r(u, v) = \sqrt{(u - c_x)^2 + (v - c_y)^2} . \quad (10)$$

Figure 3 shows the resulting semi-dense depth maps for two KITTI sequences.

3.7. Global Map Optimization

So far, we presented an approach performing incremental visual odometry by directly tracking incoming stereo images in combination with semi-dense depth reconstruction. This gradual pose estimation technique accumulates errors over time.

In order to alleviate this caveat we use G²O [34] for global pose graph optimization. The pose graph is constructed from the key frame poses as vertices and their relative transformations as edges. Instead of optimizing $SIM(3)$ constraints, i.e., assuming that two separate camera frames are related via a rigid-body motion with an additional unknown scaling factor, as in the original proposal by Engel et al. [33], we set constraints between key frames as their $SE(3)$ rigid-body motion as due to our stereo setup we have no scalar ambiguity. Once a key frame is created, its pose is added to the key frame graph as a vertex V_i . Subsequently, we search the existing vertices in the graph for additional constraints that can refine the pose graph. To this end, the closest n key frames, that have sufficient scene overlap in terms of parametrizable euclidean distance as well as parametrizable angular overlap, are matched against the newly created key frame: We estimate the transformation between both frames both ways, by registering the constraint candidate against the key frame and vice versa using semi-dense direct matching as described in Equation (1). Only if the matching succeeds for both directions and the resulting transformations agree, they are added

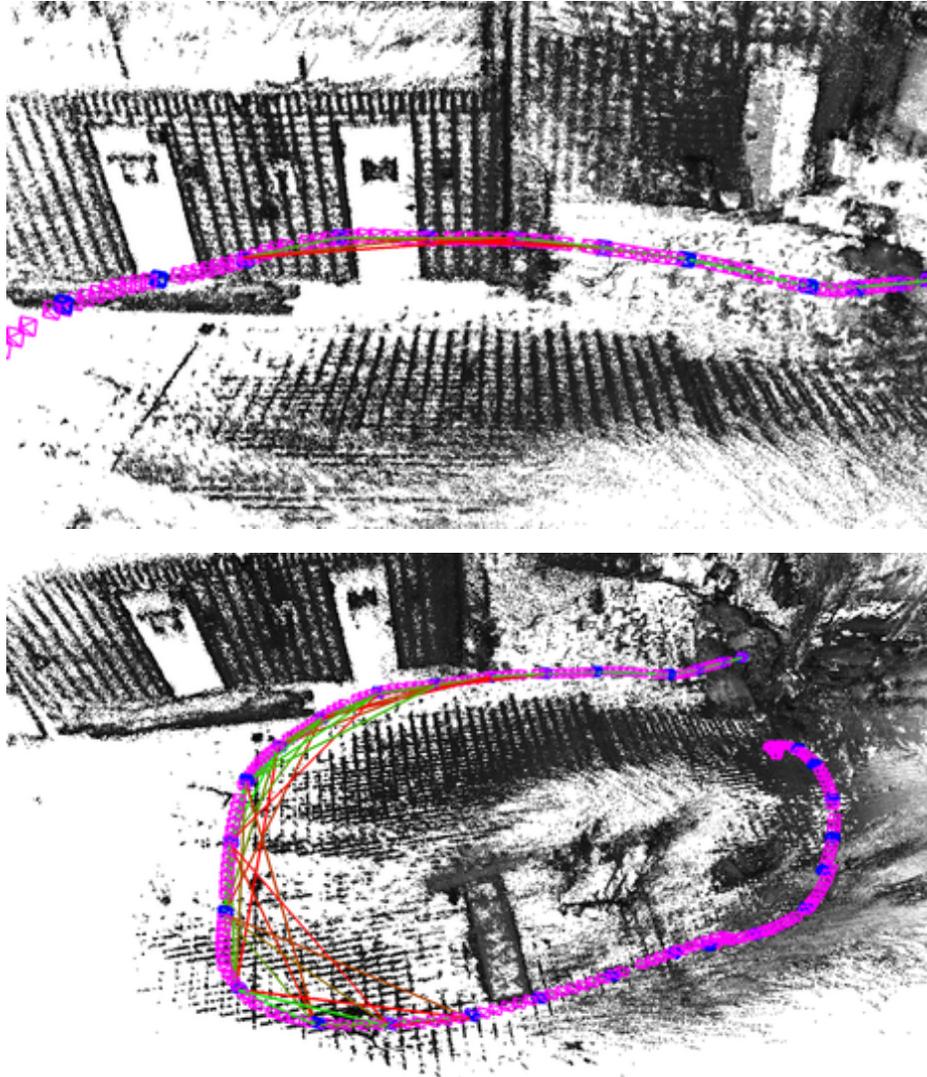


Figure 4: Trajectory of a camera and reconstructed map. Pink frusta indicate camera poses for every frame taken, blue frusta denote keyframes. Edges between keyframes are color-coded using the reprojection error between them. The coding ranges from green (low reprojection error) to red (high reprojection error). In the global map optimization edges between key frames are added to the pose-graph. Nearby key frames usually match better than distant key frames (visualized as line color: green for a good match, red for a poor match). The top picture shows the graph before optimization with many red edges indicating poor consistency. Below: After graph optimization the key frame poses have been refined yielding a lower overall error.

to the graph as an additional constraint. All edge constraints E_{ij} between the vertices V_i, V_j define a cost function C that is optimized using G²O:

$$C(V_0, \dots, V_l) = \sum_{E_{ij}} \|T_V(V_i, V_j) - T_E(E_{ij})\|^2. \quad (11)$$

where T_V denotes the relative pose between two key frames corresponding to V_i and V_j and T_E yields the relative pose held by the edge E_{ij} . Figure 4 shows an exemplifying result of the global optimization stage.

4. Evaluation

For the evaluation of our semi-direct approach we perform experiments on three challenging stereo datasets: the well-known KITTI-dataset [35], the EuRoC dataset [36] and a proprietary dataset recorded with our high-performance MAV presented in Section 4.3. The datasets differ in terms of frame rate, apparent motion, stereo baseline, and base platform. All experiments have been conducted on an Intel Core i7-4702MQ running at 2.2 GHz with 8 GB RAM.

We compare the quality of our combined approach in terms of accuracy and runtime to LSD-SLAM [1] and LIBVISO2 [7], as well as to state-of-the-art methods like S-PTAM [10] and ORB-SLAM [6]. The execution of the referred methods has been obtained using the provided default parameters.

As ground truth for all sequences is available, we employ the evaluation metrics by Sturm [37] and measure the absolute trajectory error (ATE) by computing the root mean squared error (RMSE) over the whole trajectory. In addition, we also provide the median error for better insight, because single outliers can greatly affect the final result. The ATE is a popular measure for the evaluation of visual SLAM systems as it measures the Euclidean distance between ground truth poses and estimated poses at corresponding timestamps, and thereby allows to evaluate the global consistency of SLAM systems. In a first step the trajectories are rigidly aligned because they stem from different coordinate systems. Moreover, a similarity alignment is performed for the monocular systems to estimate the absolute scale of the estimated trajectory. For an intuitively accessible visualization, trajectories are always shown in bird’s eye perspective neglecting height differences in the trajectory. In the following sections we first present detailed results for each dataset, individually. Moreover, we evaluate the performance of visual

KITTI Sequence	Absolute Trajectory Error RMSE (Median) in m			
	Ours	LIBVISO2	ORB-SLAM	S-PTAM
00	5.79 (4.54)	29.71 (18.49)	8.30 (6.04)	7.83 (6.30)
01	61.55 (54.57)	66.54 (60.46)	335.52 (303.79)	204.65 (157.10)
02	18.99 (14.38)	34.26 (27.36)	18.66 (15.03)	20.78 (17.28)
03	0.63 (0.52)	1.67 (1.54)	11.91 (9.19)	10.53 (10.41)
04	0.67 (0.46)	0.80 (0.66)	2.15 (1.73)	0.98 (0.88)
05	5.47 (4.14)	22.14 (19.07)	4.93 (4.73)	2.80 (2.24)
06	2.06 (1.80)	11.54 (10.26)	16.01 (15.56)	4.00 (4.01)
07	2.34 (1.67)	4.41 (4.37)	4.30 (3.65)	1.80 (1.53)
08	8.42 (7.04)	47.67 (34.84)	38.80 (18.12)	5.13 (4.26)
09	5.46 (3.33)	89.83 (77.57)	7.46 (6.91)	7.27 (4.61)
10	1.68 (1.37)	49.35 (36.00)	8.35 (7.55)	2.08 (1.70)
mean	10.28 (8.53)	32.54 (26.42)	41.49 (35.66)	25.74 (20.26)
w/o S 01	5.15 (3.93)	29.14 (23.02)	12.09 (8.85)	7.85 (6.57)

Table 1: ATE Results on KITTI Dataset

SLAM compared to pure visual odometry and provide quantitative result. Afterwards, we shortly summarize the obtained average results for accuracy and runtime and conclude with qualitative results of our 3D reconstruction. Finally, the following videos summarize our visual SLAM approach on the EuRoc dataset ¹ and visual odometry on the KITTI dataset ².

4.1. KITTI

The KITTI dataset [35] is a very popular dataset for the evaluation of visual and laser-based odometry or SLAM methods. It contains 22 stereo sequences accompanied by laser scans, and ground truth from a localization unit consisting of a GPS and an IMU. The stereo camera rig and the laser scanner are mounted on top of a standard station wagon—the autonomous driving platform Anniway [38]. The stereo rig has a baseline of approximately 54 cm.

Rectified images are provided with 10 Hz at a resolution of 1240×376 pixels. The sequences are recorded in real-world driving situations along

¹<https://www.youtube.com/watch?v=7NkHf6syRIo>

²<https://www.youtube.com/watch?v=PRYgnIBDVGI>

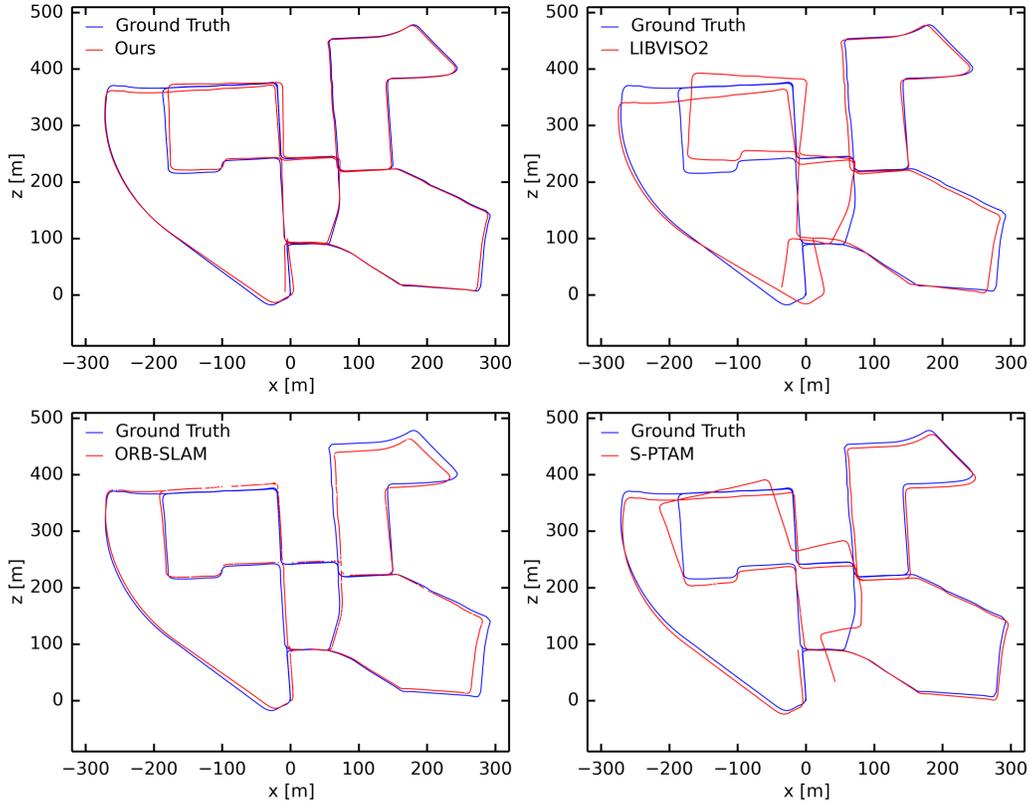


Figure 5: Results for KITTI Sequence 00. Comparison of our method to LIBVISO2 (top row), ORB-SLAM and S-PTAM (bottom row). Our methods achieves the lowest ATE (5.79 m)

urban, residual and countryside roads. The distance traveled ranges from a few 100 meters up to 5 kilometers with driving speeds up to 80 km/h.

The dataset is very challenging because the low frame rate in combination with fast driving speed leads to large inter-frame motions of up to 2.8 m per frame. This strongly restricts the number of possible feature correspondences. Moreover, moving obstacles in form of passing vehicles, bicycles, or pedestrians that have great impact on the performance of visual odometry systems, are included frequently.

We compare the performance of our semi-direct method with four state-of-the-art methods for visual odometry and SLAM.

We selected LIBVISO2 and LSD-SLAM for reference as our method is built upon them. Moreover, we chose two established feature-based SLAM

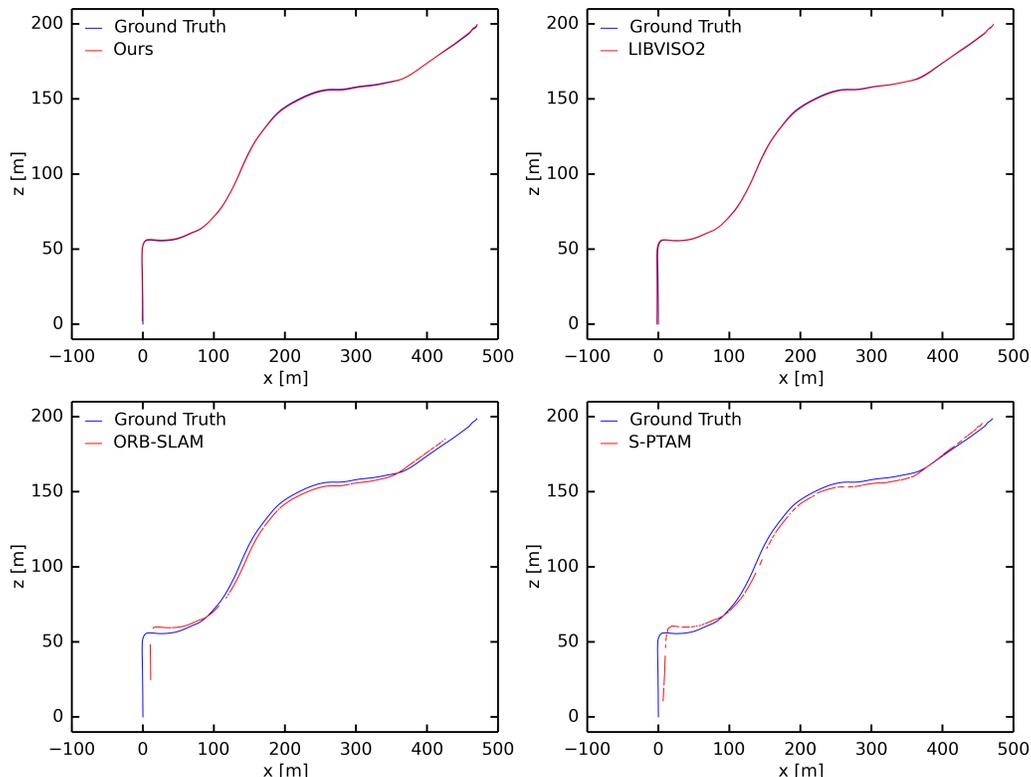


Figure 6: Results for KITTI Sequence 03. Comparison of our method to LIBVISO2 (top row), ORB-SLAM and S-PTAM (bottom row). Our method and LIBVISO2 show accurate trajectories.

algorithms that presented promising results: ORB-SLAM as a monocular and S-PTAM as a stereo method. All processing is done on the original image resolution of the rectified images of 1240×376 .

Both error measures—RMSE and Median—for the training sequences 00 to 10 of the KITTI dataset are listed in Table 1.

Unfortunately, LSD-SLAM fails on all sequences of the KITTI dataset. This is probably caused by too large inter-frame motion for a pure monocular direct method, as sufficient scene overlap is important for successful tracking. Moreover, it can be seen, that all SLAM methods lack performance on sequence 01, resulting in a very high ATE. Sequence 01 contains images from driving on a highway, thus it is hard to find re-occurring feature points in subsequent frames.

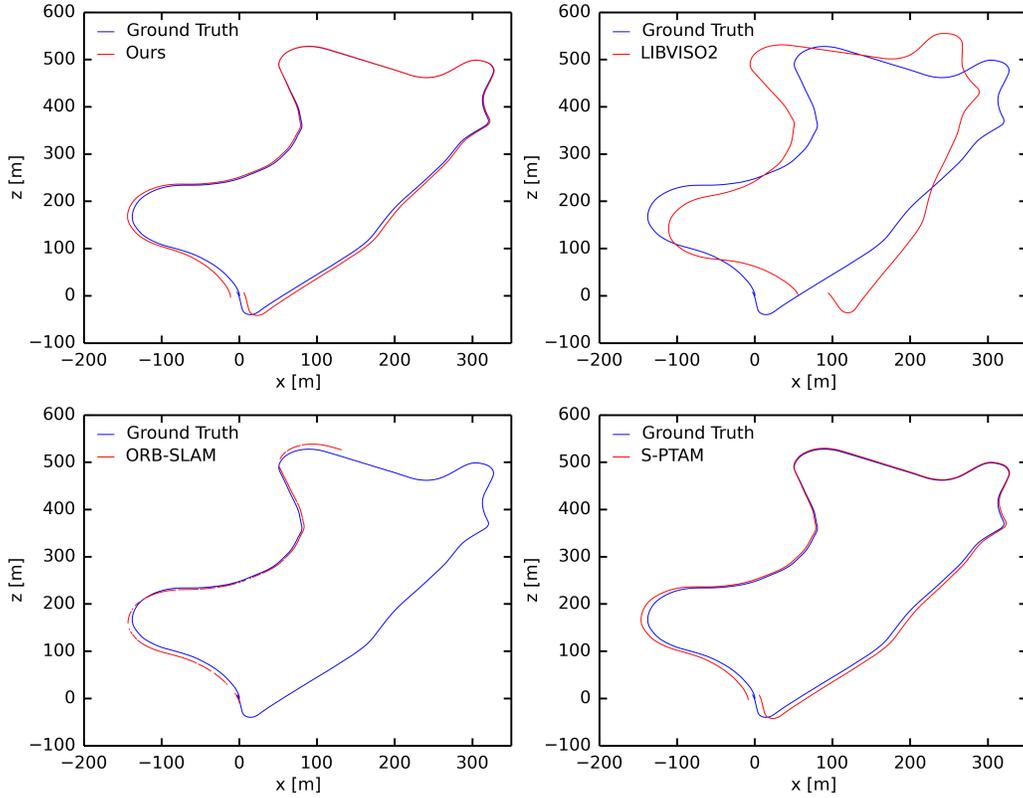


Figure 7: Results for KITTI Sequence 09. Comparison of our method to LIBVISO2 (top row), ORB-SLAM and S-PTAM (bottom row). While ORB-SLAM loses track and LIBVISO2 accumulates drift, S-PTAM and our method stay close to the ground truth.

Overall our method is equally good and in seven of eleven cases even better than state-of-the-art methods. Especially sequences 03 and 04 show very accurate results below 1 m ATE. In three of the cases S-PTAM and in one case (sequence 02) ORB-SLAM performs better. As LIBVISO2 is a pure odometry method, it performs significantly worse than the SLAM methods on all datasets.

When averaging over the eleven training sequences our method ranks first, followed by S-PTAM, ORB-SLAM and LIBVISO2. However, the bad results from sequence 01 greatly affect the final average computation as all methods accumulate high ATEs in sequence 01. One could argue that such high ATE values count as outlier or failure. Therefore, we also show resulting means when omitting sequence 01 for all methods. It follows that these results show

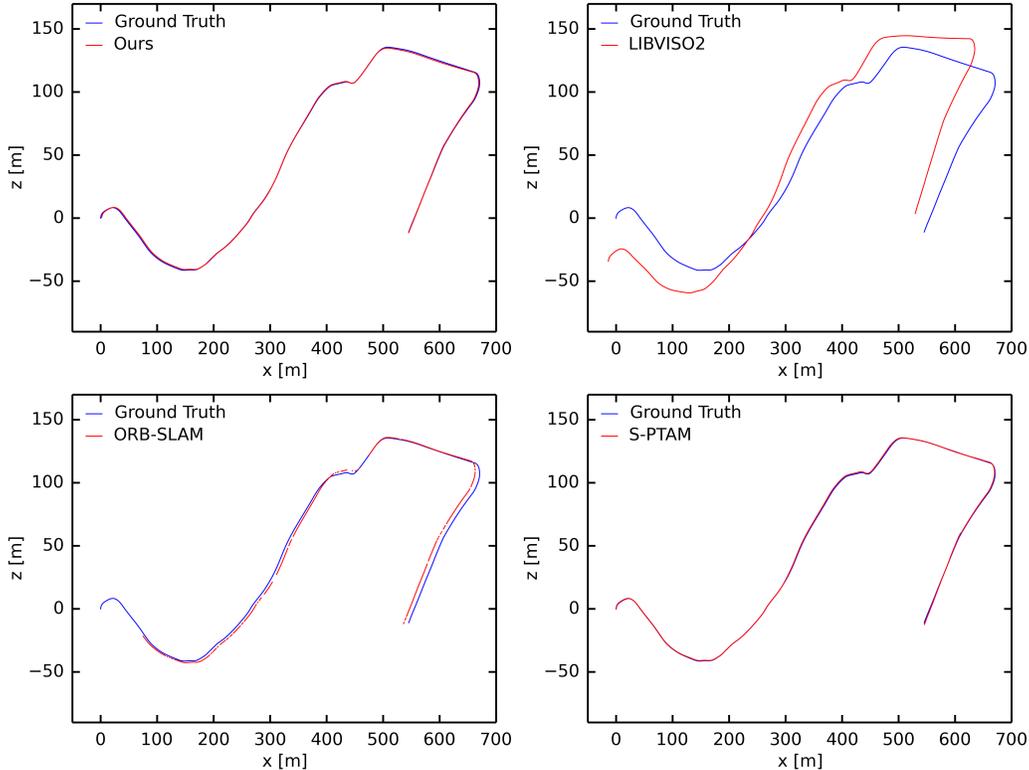


Figure 8: Results for KITTI Sequence 10. Comparison of our method to LIBVISO2 (top row), ORB-SLAM and S-PTAM (bottom row) on a longer trajectory without loop closures.

significantly lower ATEs. When omitting sequence 01 our method achieves a mean (median) ATE of 5.15 m (3.99 m) compared to the S-PTAM result of 7.85 m (6.57 m).

For a better visualization exemplary trajectories are shown in birds-eye perspective for the sequences 00, 03, 09 and 10.

Sequence 00 is shown in Figure 5. It can be seen, that our approach performs best, followed by ORB-SLAM, S-PTAM and LIBVISO2. Moreover, limitations of the approaches become visible: as LIBVISO2 is a pure odometry method, it accumulates more drift over time, and as ORB-SLAM is a monocular method, scale is not always estimated correctly. Rotations are challenging for all methods. In this sequence S-PTAM fails to track rotations frequently and exhibits drift for the last part of the trajectory.

Figure 6 shows the results for sequence 03, a trajectory without full loop

closure. We choose this sequence to compare the drift over time, when no full loop can be closed. All estimated trajectories are close to the ground truth. However, our method is—with 0.63 m ATE—distinctively more accurate than ORB-SLAM (11.91 m) and S-PTAM (10.53 m). Additionally, LIBVISO2 also shows accurate results with an ATE of 1.67 m and does not accumulate much drift for this trajectory.

In sequence 09 a full loop closure appears at the very end of the trajectory, that is not always detected from the SLAM methods before the sequence ends. This behavior is shown in Figure 7. Again, LIBVISO2 suffers from drift over time while the results from our Semi-Direct SLAM (5.46 m) are more accurate than the results from S-PTAM (7.27 m). Moreover, it can be seen, that ORB-SLAM lost track at some point and failed to relocalize. Thus, more than half of the trajectory remains uncovered. This is not reflected in the error measure, because the ATE is only computed over existing measurements.

Sequence 10 is similar to sequence 03 as it contains no full loop but covers a longer path and performs more rotations. Results for this sequence are visualized in Figure 8. They show that our method performs well even if the path of LIBVISO2 drifts over time. ORB-SLAM fails to initialize directly from the beginning but retrieves a trajectory consistent with the ground truth later on, even with little offset. S-PTAM again shows similar results to Semi-Direct SLAM although Semi-Direct SLAM performs slightly better (1.68 m to 2.08 m respectively).

However, as LSD-SLAM fails on the KITTI sequences, we compare our semi-direct approach to its fully direct version without feature-based initial estimates. In particular, we compare the combined semi-direct approach to its building blocks—LIBVISO2 and direct stereo tracking—separately. As LIBVISO2 is a pure odometry method, we evaluate it against results from our semi-direct odometry without closing loops.

In Figure 9 the results from three different datasets (00, 02 and 06) are shown in birds-eye perspective. The left column shows the resulting path LIBVISO2 computed and the right column the path from the semi-direct odometry. It can be seen that LIBVISO2 accumulates more drift over time than the semi-direct approach, while at the same time the semi-direct approach remains closer to the ground truth trajectory.

When comparing our semi-direct approach to its fully direct version without feature-based odometry as initial estimate, we noticed that a fully direct version has problems with strong turns in the dataset. Moreover, the

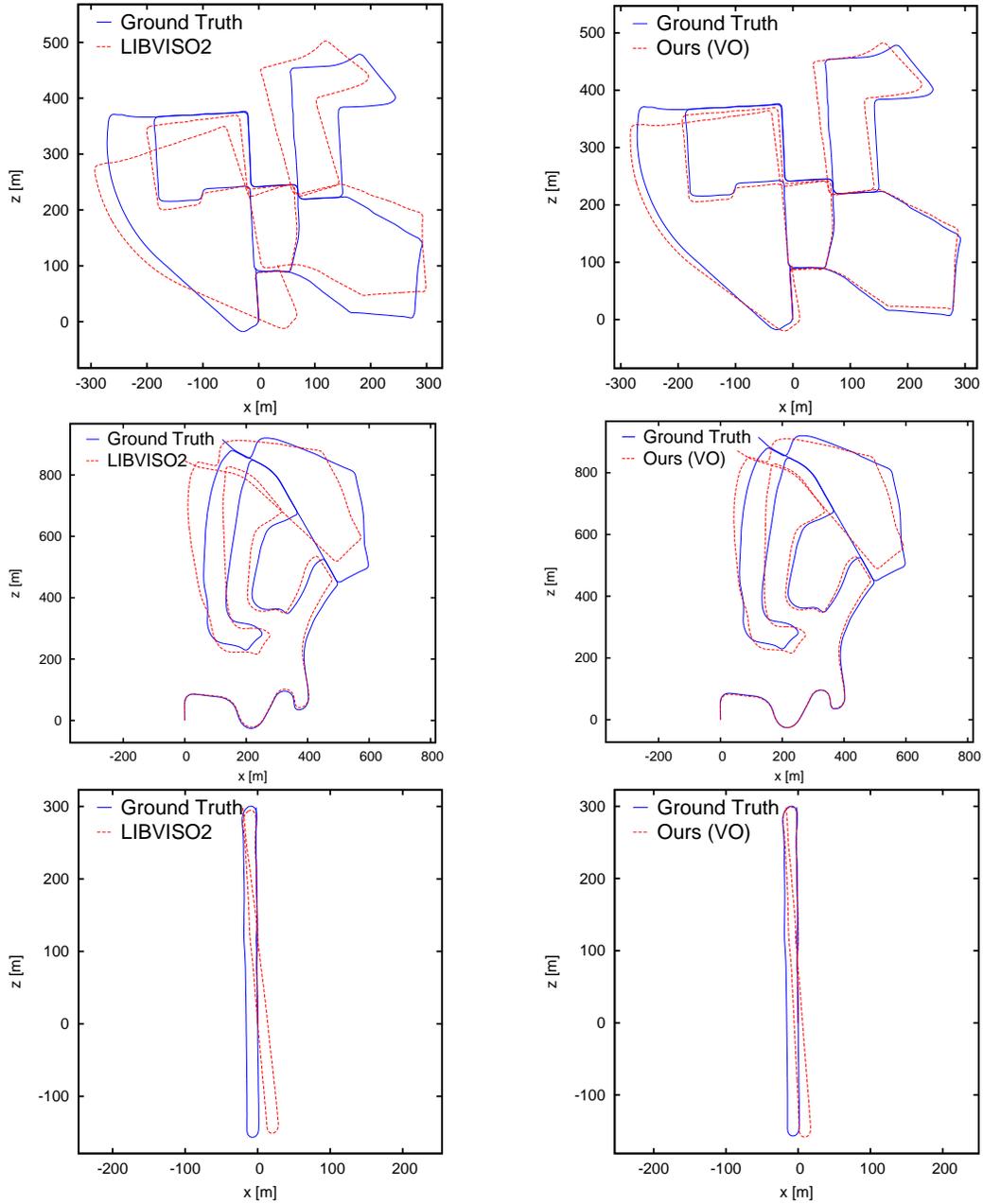


Figure 9: Comparison of the results from LIBVISO2 (left) to our semi-direct odometry (right). Top Row: KITTI Sequence 00, Middle Row: KITTI Sequence 02, Bottom Row: KITTI Sequence 06. In direct comparison to LIBVISO2 our method accumulates less drift.

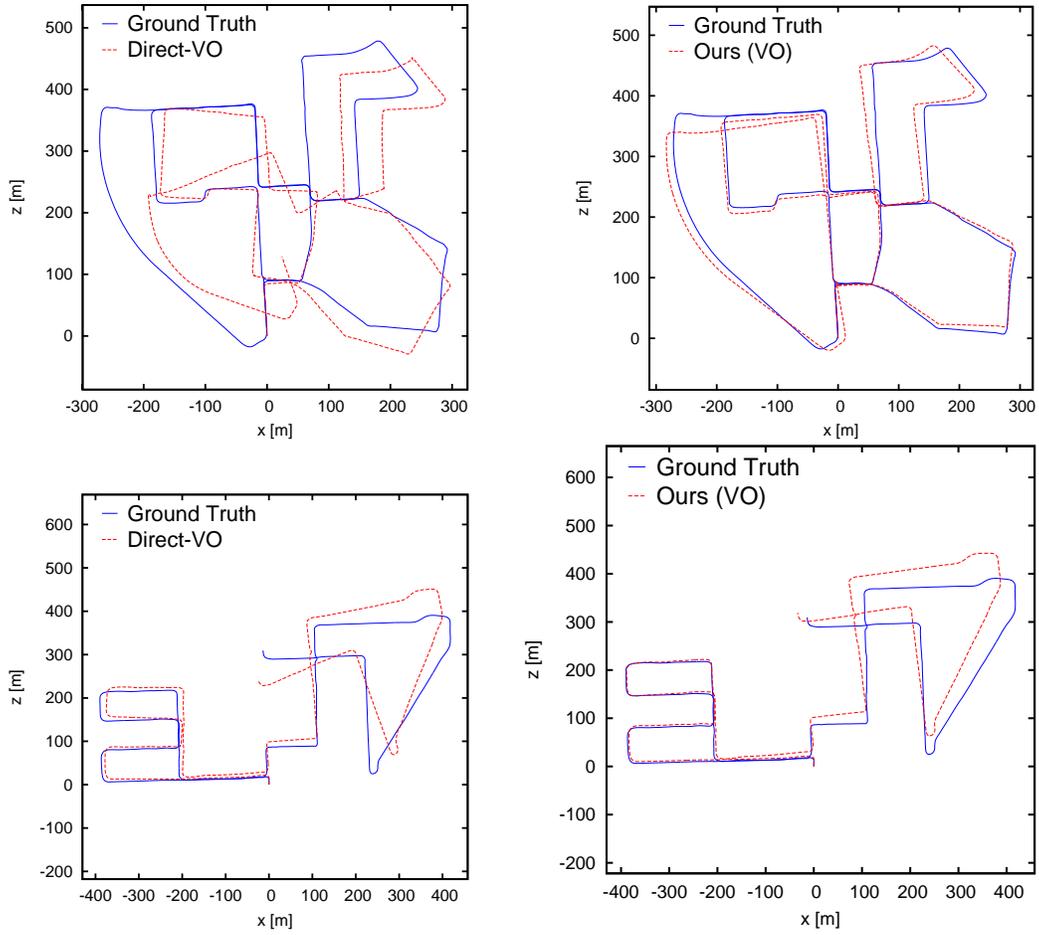


Figure 10: Comparison of the results from direct odometry (left) to our semi-direct odometry (right). Top Row: KITTI Sequence 00, Bottom Row: KITTI Sequence 08. In direct comparison to the direct visual odometry our method is clearly more robust to fast rotations and to large motions.

dataset is very challenging for a fully direct method as it contains large inter-frame motions and difficult lighting changes. Large inter-frame motions are problematic for direct methods because they assume small pixel displacements [39]. Without a good initial estimate they often fail to retrieve large displacements. Difficult lighting changes, induced by auto-exposure and changing sunlight, are tough, as they violate the brightness constancy assumption. Thereby, it can be seen in Figure 10 that the fully direct odometry accumulates more drift over time than our semi-direct version. Again, while the semi-direct approach is shown in the right column, the fully direct approach is visualized in the left column for sequences 00 and 06. Fully direct tracking tends to fail especially at strong turns and at street crossings where lighting changes increase because the car leaves shadowed street canyons. In contrast, our approach is more robust to strong rotations and illumination changes.

The semi-direct method performs better than its isolated building blocks. The direct tracking is in principle more accurate but has problems with large motions. However, when a good initial estimate is available, as in our case from LIBVISO2, direct tracking succeeds even at large motions and with a low frame rate.

Generally speaking, a combined semi-direct odometry performs better than both—feature-based and direct—odometries alone. Overall, our approach shows promising results on the KITTI dataset when compared to other state-of-the-art methods.

4.2. *EuRoC*

In addition to the evaluation on the KITTI dataset, we perform further experiments on the well-known visual-inertial EuRoC MAV dataset that contains stereo images and synchronized IMU readings from the on-board computer of an Asctec Firefly hex-rotor helicopter. We choose six trajectories with different difficulties from the two Vicon datasets V0 and V1. The data has been collected from flights in a room that is equipped with a Vicon motion capture system offering 6D ground truth poses.

The MAV carries a visual-inertial sensor [40] that captures stereo images of WVGA resolution at 20 Hz and synchronized IMU measurements at 200 Hz.

Both datasets contain three trajectories with increasing difficulty named as: easy (.01), medium (.02) and difficult (.03). The easy trajectories have good illumination, are feature rich, and show no motion blur and only low

EuRoC Dataset	Absolute Trajectory Error RMSE (Median) in m				
	Ours	Libviso2	LSD-SLAM	ORB-SLAM	S-PTAM
V1_01	0.12 (0.11)	0.31 (0.31)	0.19 (0.10)	0.79 (0.62)	0.28 (0.19)
V1_02	0.11 (0.10)	0.29 (0.27)	0.98 (0.92)	0.98 (0.87)	0.50 (0.35)
V1_03	0.75 (0.45)	0.87 (0.64)	X	2.12 (1.38)	1.36 (1.09)
V2_01	0.18 (0.12)	0.40 (0.31)	0.45 (0.41)	0.50 (0.42)	2.38 (1.78)
V2_02	0.27 (0.22)	1.29 (1.08)	0.51 (0.48)	1.76 (1.39)	4.58 (4.18)
V2_03	0.87 (0.66)	1.99 (1.66)	X	X	X
mean	0.38 (0.28)	0.85 (0.71)	0.53 (0.48)	1.23 (0.94)	1.82 (1.52)

Table 2: ATE Results on EuRoC Dataset

optical flow and low varying scene depth. They capture a static scene. The difficulty increases in the medium trajectories by adding difficult lighting conditions, high optical flow and medium varying scene depth. However, they still show a static scene and a feature rich environment without motion blur. In contrast, the difficult scenes contain areas with only few visual features and more repetitive structures. Moreover, they add motion blur and highly unstable lighting conditions. The MAV performs very aggressive flight maneuvers resulting in high optical flow and highly varying scene depth in a non-static scene.

The dataset is known to have different issues that make a reliable state-estimation more challenging: for example, the stereo images were captured using an automatic exposure control that is independent for both cameras. Therefore, shutter times are different, which results in different image brightnesses, making stereo matching and feature tracking more challenging. This is especially important, as direct methods minimize the photometric error.

Moreover, as the ground truth is recorded from a different physical device than the images, the accuracy depends on the synchronization scheme used [36].

The resulting ATEs are listed in Table 2. As the difficult datasets V1_03 and V2_03 contain very dynamic movements and fast rotations with an MAV, LSD-SLAM often loses track after a few seconds and is then unable to re-localize for the rest of the trajectory. In Table 2 this is denoted as failure (X). Similarly, S-PTAM and ORB-SLAM lose track for the difficult trajectory V2_03. This dataset shows very challenging conditions with strong motion blur and fast aggressive maneuvers. Moreover, the absence of sufficient visual

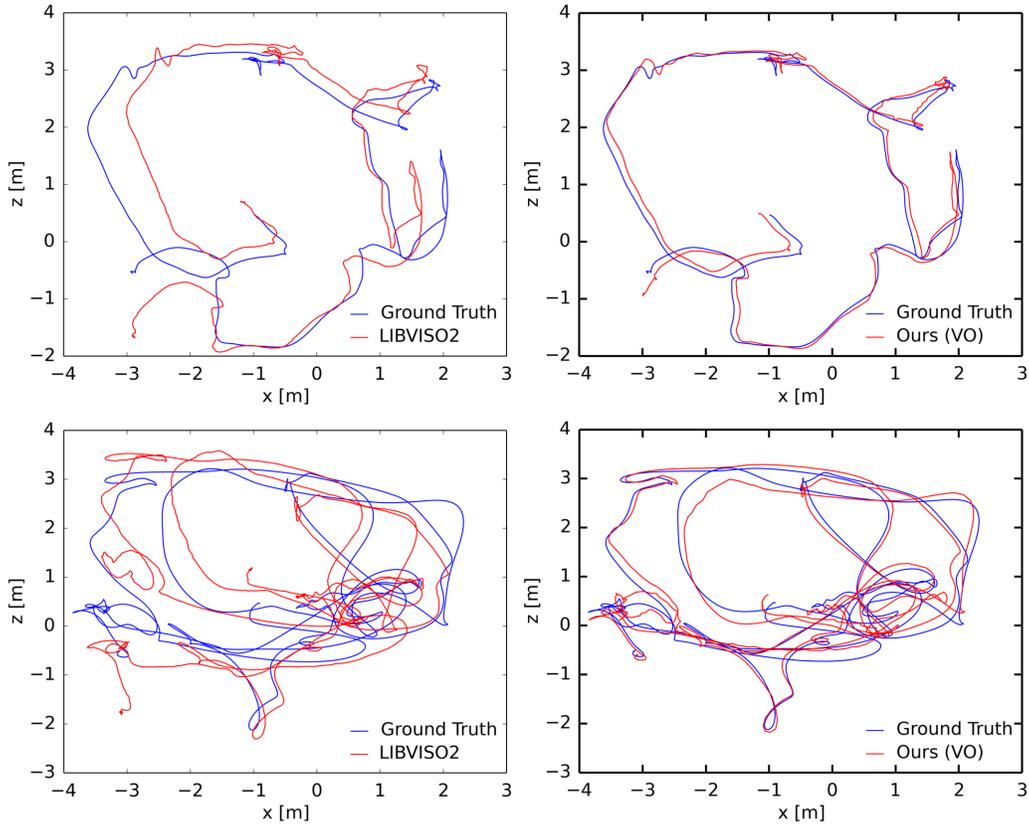


Figure 11: Comparison of the results from LIBVISO2 (left) to our semi-direct odometry (right) on datasets V2.01 and V2.02 with ground truth from a Vicon motion capture system. Again our method is much closer to the ground truth even without SLAM.

features makes it hard for feature-based methods to succeed.

Table 2 also shows, that our approach outperforms the other methods and reliably recovers the motion for all test sequences. Additionally, it can be seen, that the results of LIBVISO2 are improved on every trajectory.

On average semi-direct SLAM achieves a higher accuracy, with 0.38 m ATE, than LSD-SLAM, with 0.53 m, ORB-SLAM, with 1.23 m ATE, and S-PTAM with 1.82 m ATE. LSD-SLAM, ORB-SLAM and S-PTAM often perform poorly at fast motions in combination with rotations, and then tend to lose track temporarily.

Additionally, we again directly compare results from Semi-Direct Visual Odometry to LIBVISO2 and to Direct Odometry from LSD-SLAM. Figure 11

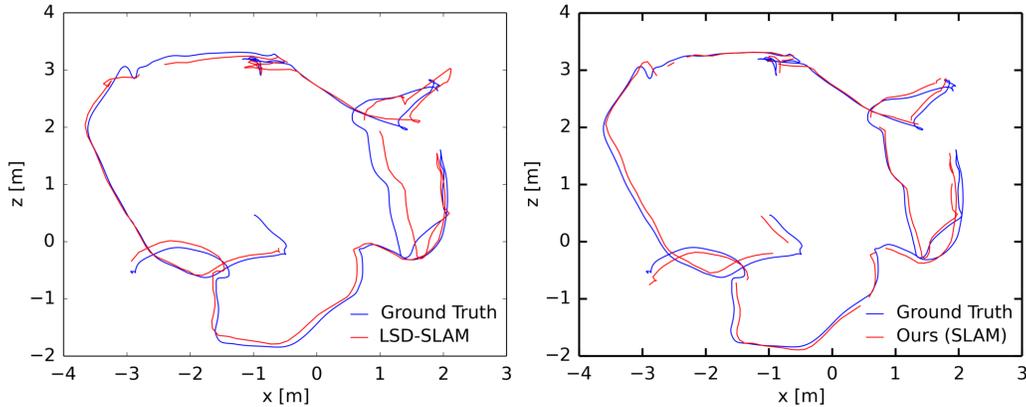


Figure 12: Comparison of the results from LSD-SLAM (left) to our semi-direct SLAM (right) on dataset V2_01 with ground truth from a Vicon motion capture system. While LSD-SLAM shows an ATE of 0.45 m our methods performs better with an ATE of 0.18 m

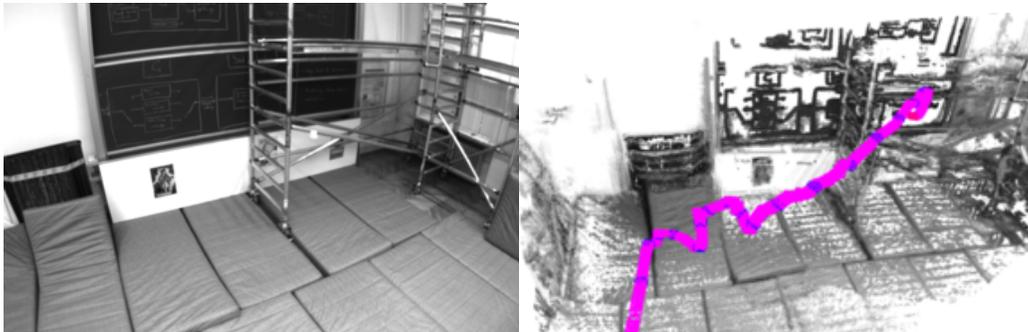


Figure 13: Exemplary results of flight Vicon_m: While the left image shows a capture of the recorded scene, the right image shows the retrieved camera trajectory and reconstructed semi-dense depth. Key frames are shown in blue, while feature-based tracked frames are shown in pink.

shows the resulting trajectories for datasets V2_01 and V2_02 of LIBVISO2 and Semi-Direct Visual Odometry. Both methods were performed without loop closures, and, thus, accumulate small errors in the estimates over time. It can clearly be seen that the Semi-Direct Odometry is closer to the ground truth from the Vicon system than LIBVISO2. Even though the datasets contain rapid rotations, our method stays close to the ground truth path.

In comparison to LSD-SLAM, our approach is more robust to fast rotations in the trajectory as can be seen in Figure 12. While LSD-SLAM computes wrong estimates at strong turns, our method follows the path more

precisely.

In addition to the official datasets we performed one manual flight in the Vicon room, named Vicon_m, where we evaluated the mapping abilities of our approach. As an example sequence for our mapping abilities, Figure 13 shows a sequence captured on a manual flight: the MAV captures a corner of the Vicon room and is able to reconstruct a semi-dense 3D representation of the recorded scene. As can be seen, details of the scaffold are retrieved as well as the ground plane and drawings on the blackboard. The recovered camera trajectory is shown as well. Key frames are colored in blue while frames that were tracked feature-based with LIBVISO2 are shown in pink.

In total, we showed that our method is more robust to dynamic motions than the other evaluated methods and achieves a lower ATE on all evaluated datasets.

4.3. MAV

In the previous section, we demonstrated that our semi-direct approach is capable of accurate pose estimation with standard stereo cameras. We furthermore evaluate the performance of our approach on different datasets, that have been acquired with our MAV, shown in Figure 14. In contrast to the setups in previous datasets, our MAV is equipped with fish eye lenses and a wide baseline.

Our MAV is built as high-performance platform with a multimodal omnidirectional sensor setup [41]. As MAVs have very limited payload, we use only lightweight components and are capable of navigating indoor and outdoor. Especially for (fully) autonomous navigation in unknown and dynamic environments, a multimodal and omnidirectional sensor setup is of great advantage. The strengths of the different sensors can be combined and their measurements can be fused in an occupancy grid map.

The MAV is built as hexarotor with six 14" propellers each connected to a MK3644/24 motor. For better stability and collision protection the MAV is surrounded with a non-rigid milled frame that does not only protect the rotors, but also serves as mount for various sensors. For on-board computation in real-time, the MAV is equipped with a mini-ITX board, namely a Gigabyte GB-BXi7-4770R with an Intel Core i7-4770R quad-core CPU, 16 GB DDR-3 memory and a 480 GB SSD to process all sensor outputs.

We employ a multimodal sensor setup consisting of IMU, laser scanners and cameras. Moreover, our system is equipped with two laser scanners and six cameras for high-level autonomous operation and navigation. In

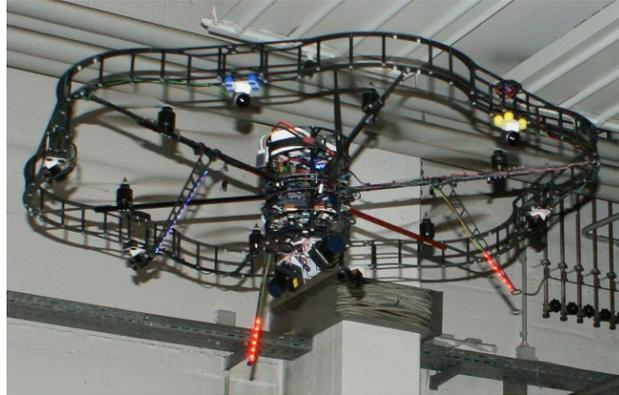


Figure 14: High performance MAV during flight. The omnidirectional sensor setup includes three fish eye stereo pairs covering a wide field of view for autonomous navigation.

particular, we use two rotating Hokuyo UST-20LX laser scanners, each with a scan range of 20 m and 270° apex angle. Together they can perform a full 3D scan of the environment with 4 Hz. They are used for obstacle perception and SLAM-based 6DOF localization [42].

For visual obstacle detection and visual SLAM, the MAV is equipped with an omnidirectional camera setup. The cameras are mounted to the non-rigid body frame using dampers to filter out vibrations induced by the six propellers. The camera mounting can easily be switched from a fully omnidirectional setup with independent optical axes to a stereo setup with three stereo camera pairs, as can be seen in Figure 15. The multi-camera setup allows omnidirectional perception of the environment and allows robust state estimation due to redundant information sources, i.e., even if one stereo pair faces a homogeneous wall with no texture the other two pairs still allow for robust localization. We use XIMEA MQ013MG-E2 global-shutter monochrome USB 3.0 cameras with 1.3 MP resolution in combination with Lensagon BF2M2020S23 fish-eye lenses for a wide field of view. By making use of the available independent USB controllers of the on-board system, we distribute the USB traffic and thus can achieve high camera frame rates at full resolution. Each stereo pair is connected to a USB 3.0 HUB, which is connected to a dedicated on-board USB 3.0 port that offers full USB 3.0 speed. Through this setup we ensure that for each camera enough bandwidth is available. Assuming that each HUB offers 2400 Mbit/s (300 MB/s), each camera may use up to 1200 Mbit/s (150 MB/s). Theoretically, each camera

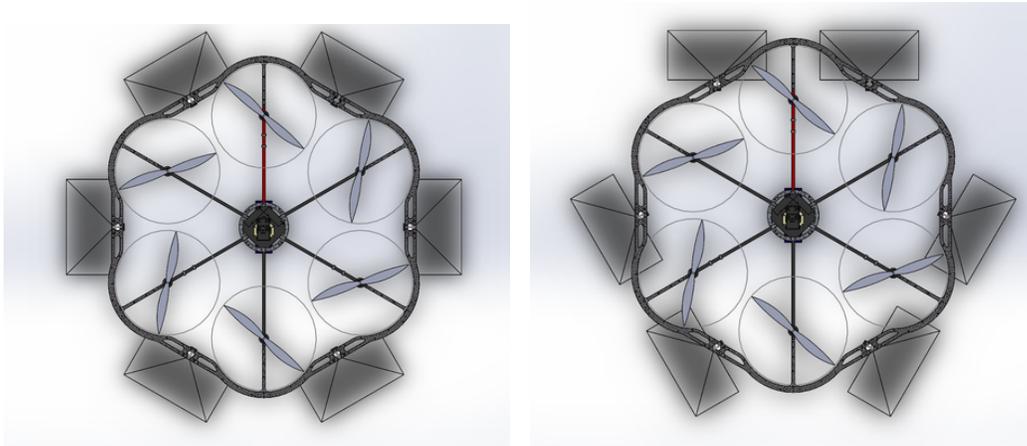


Figure 15: Top down view of possible camera configurations: the left image shows a fully omnidirectional setup with independent optical axis, while on the right a stereo setup consisting of three independent stereo pairs is shown.

can achieve the best possible frame rate of 60 Hz in 8-bit mode and 57 Hz in 16-bit mode.

However, the real data rate is limited by additional system and protocol overhead when reading and writing from the connected devices. Under real lighting conditions and depending on exposure times we achieve up to 50 Hz for each camera in 16-bit mode. Our camera driver not only ensures that the images are published synchronously, but also offers advanced functionality like downsampling, gamma correction or rectification.

We use laser-based SLAM [42] as ground truth and again compare the results with those of state-of-the-art SLAM methods. In total we captured four flights in a decommissioned car service station with challenging lighting conditions. While on the first two flights, named *rect1* and *rect2*, the MAV covers a rectangular path without many loop closures, the other two flights, *loop1* and *loop2*, include three to four full loops.

A general prerequisite for stereo computation is to rectify the images. To allow different models for calibration we build a general rectification nodelet in ROS, that rectifies the images given respective look-up tables as input. The look-up tables can be either calculated offline beforehand or online using, e.g., `cv::StereoRectify`, the computer vision library OpenCV. The rectification nodelet publishes rectified images together with camera info messages that contain the necessary calibration parameters from intrinsic and extrinsic calibration.

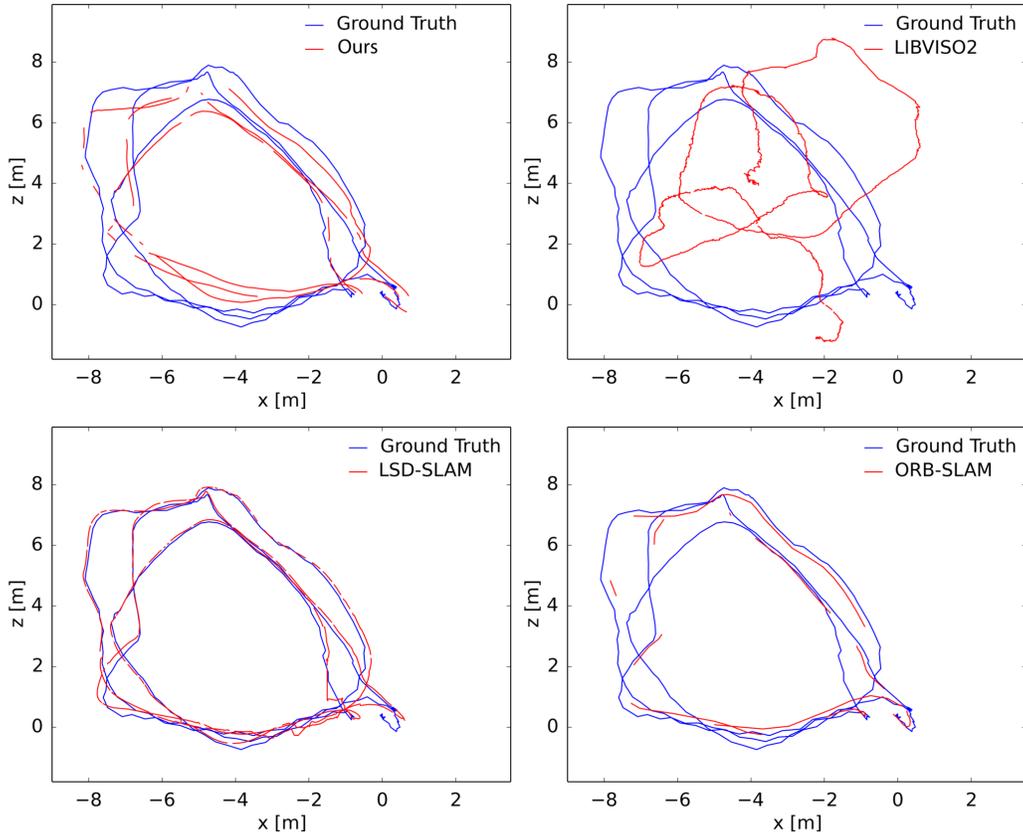


Figure 16: Results for MAV loop 1. Comparison of our method to LIBVISIO2 (top row), LSD-SLAM and ORB-SLAM (bottom row) on a challenging dataset that contains large loop closures. As can be seen the monocular methods perform best while LIBVISIO2 accumulates strong drift. Still, our method is able to reconstruct the trajectory with an ATE of 0.63 m while S-PTAM fails completely.

Moreover, we added functionality to down-sample the rectified images by a factor c for further run time enhancement. The images are captured with full resolution of 1280×1024 in 16 bit-encoding and are down-sampled to half the resolution and 8 bit in the rectification step. Figure 17 shows the result from the rectification step on an image from the recorded dataset.

The rectification of the images runs in parallel for all six cameras and takes 1 ms for a single image when downsampling to half the original resolution. For an even smaller resolution of 320×256 the rectification takes 0.7 ms and for the full resolution 4 ms.

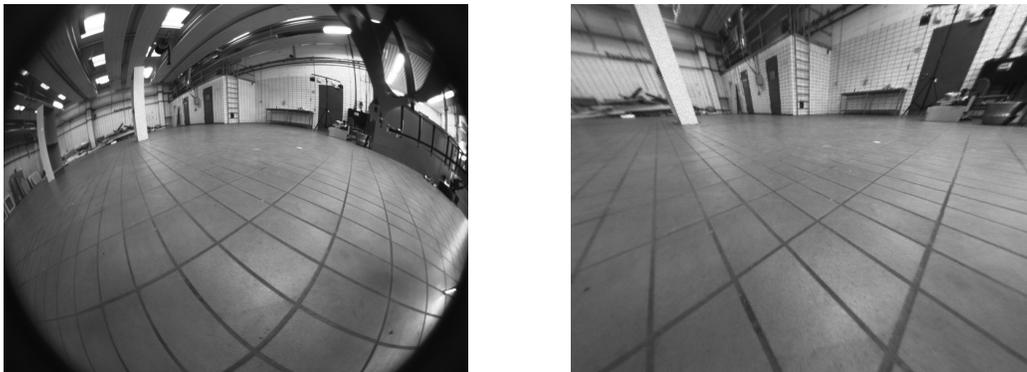


Figure 17: Stereo Rectification: Raw fish eye images (left) are rectified onto a plane with half the resolution (left).

Experiments on the four flights show that retrieving the correct camera motion is more challenging than on the previous datasets. Especially, the stereo methods do poorly on these datasets. S-PTAM fails to initialize correspondences on all datasets and, thus, cannot be taken into account for comparison.

Exemplary results are shown for a trajectory with repeating loop closures in Figure 16. It can be seen that the stereo methods Semi-Direct SLAM and LIBVISO2 show a higher offset to the ground truth trajectories than the monocular methods. Especially LIBVISO2 accumulates high errors at this circular trajectory and the result is not as accurate as before, thereby hindering the semi-direct approach. As LIBVISO2 performs no loop closure detection, errors in the absolute trajectory cannot be resolved which leads to a globally inconsistent trajectory. Semi-Direct SLAM uses only the relative motion estimates of LIBVISO with regards to the current key frame. Thereby, Semi-Direct SLAM is still able to reconstruct a path close to the ground truth with an ATE of 0.63 m. Contrarily, LSD-SLAM and ORB-SLAM achieve an ATE below 0.31 m.

The fact that monocular methods seem to perform better than stereo methods, suggests that the underlying lens distortion model was chosen to allow for a rectification mapping and, hence, semi dense stereo matching, but does not fit the used fish eye lenses very accurately. Additionally, the non-rigid mounting of the stereo cameras introduces difficult conditions for stereo correspondence search along fixed epipolar lines. We assume that the wide non-rigid baseline of 53.37 cm in combination with the perspective rec-

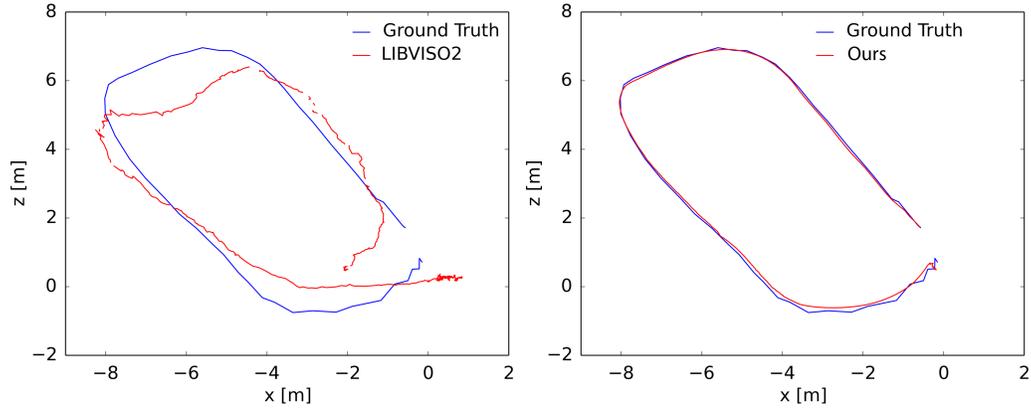


Figure 18: Comparison of LIBVISO2 (left) and our method (right) on the MAV dataset rect1. Even though, the dataset contains no loop closure, our method shows accurate results.

MAV Dataset	Absolute Trajectory Error RMSE (Median) in m			
	Ours	Libviso2	LSD-SLAM	ORB-SLAM
rect1	0.13 (0.11)	1.24 (0.49)	0.30 (0.29)	0.98 (0.24)
rect2	0.84 (0.81)	1.61 (1.59)	0.38 (0.37)	0.59 (0.25)
loop1	0.63 (0.57)	1.66 (0.99)	0.31 (0.28)	0.25 (0.21)
loop2	1.58 (0.71)	2.61 (1.90)	0.54 (0.42)	1.19 (0.78)
mean	0.80 (0.55)	1.78 (1.24)	0.38 (0.34)	0.75 (0.37)

Table 3: ATE Results on MAV Dataset

tification onto a plane impedes the stereo correspondence search. It would generally be more appropriate to model the fish eye lenses as a projection onto a sphere. As described in Section 3.6, we use an additional weighting scheme that downweights the influence of inaccurate depth measurements close to the image borders in order to cope with strong distortions. Moreover, we repeatedly estimate the extrinsic transformation of the cameras online to handle the non-rigidity. Therefore, we are able to retrieve stereo correspondences and estimate the trajectory on this challenging dataset, unlike S-PTAM which fails to initialize any correspondences.

Figure 18 shows trajectories for the sequence `rect1` computed by LIBVISO2 and Semi-Direct SLAM. The output of LIBVISO2 shows very noisy estimates and leads to a comparably high ATE of 1.24 m. In contrast Semi-Direct SLAM produces a smoother trajectory with an ATE of 0.13 m. However, in general we achieve a higher ATE than the monocular methods. Table 3 summarizes the resulting ATE on all datasets.

In terms of accuracy the monocular methods perform better than all stereo methods. This time S-PTAM is unable to track features on all datasets and fails in recovering any motion. It is remarkable that monocular methods perform better than stereo methods on these datasets which supports our assumption that the rectification of the fish eye images onto a plane in combination with non-rigidly mounted stereo cameras is very challenging for stereo computations. Moreover, the wide baseline is problematic as the image overlap between both stereo images is reduced.

On average, we achieve an ATE of 0.8 m while monocular LSD-SLAM achieves an average ATE of 0.38 m.

4.4. *Odometry versus SLAM*

In this section, we will compare the quantitative results of visual odometry to visual SLAM. As visual odometry tends to drift over time, global optimization methods such as bundle adjustment or pose graph optimization help to reduce the drift.

In Semi-Direct SLAM loop closures are detected between key frames and are added as additional constraints to the global pose graph (see Section 3.7).

The trajectories of the EuRoC dataset contain many possible loop closures. Therefore, we show comparative results between visual odometry and SLAM exemplary on this dataset. Qualitative results are listed in Table 4. In addition to the ATE as error measure, we also state the percentage im-

EuRoC Dataset	Absolute Trajectory Error		
	RMSE (Median) in m and Improvement in %		
	Our Semi-Direct VO [m]	Our Semi-Direct SLAM [m]	Improvement [%]
V1_01	0.26 (0.18)	0.12 (0.11)	53.85 (38.89)
V1_02	0.59 (0.59)	0.11 (0.10)	81.36 (83.05)
V1_03	0.81 (0.76)	0.75 (0.44)	7.41 (42.11)
V2_01	0.22 (0.13)	0.18 (0.12)	18.18 (7.69)
V2_02	0.31 (0.25)	0.27 (0.22)	12.90 (12.00)
V2_03	1.13 (0.97)	0.87 (0.66)	23.01 (31.96)
mean	0.55 (0.48)	0.38 (0.28)	30.72 (42.71)

Table 4: Odometry compared to SLAM on EuRoC

provement gained by SLAM. We measure the improvement as

$$Improvement = \frac{VO - SLAM}{VO}. \quad (12)$$

The average improvement for all seven trajectories lies at 30.72% denoting an absolute improvement of 0.17m on average. It can clearly be seen that for each trajectory the odometry result is further improved by SLAM. The improvements range from 7.41% up to 81.36%. The maximum improvement reached an absolute enhancement of 0.48m. As the trajectories V1.01 and V1.02 show significant improvements of 53.85% and 81.36% respectively, both of the results are visualized in Figure 19. The advantages of SLAM are visible in either example. In comparison to the pure odometry, SLAM retrieves trajectories closer to the ground truth. The bottom row of Figure 19 highlights the improvement of 81.36% on dataset V1.02. This dataset is of medium difficulty and contains very dynamic translational and rotational movements. It can be seen, that the odometry might be locally accurate but exhibits accumulated drift. In the global graph SLAM the drift is corrected by loop closures resulting in a better aligned trajectory.

In contrast to the EuRoC dataset the KITTI dataset shows notably less loop closure possibilities. However, when loop closures are found, the global consistency of the map is re-established. Sequence 06 contains a distinct loop. While visual odometry produces an ATE of 4.37m on Sequence 06, the result is corrected after closing the loop and the ATE decreases to 2.06m, showing an improvement of 52.9%. Figure 20 illustrates this phenomenon:

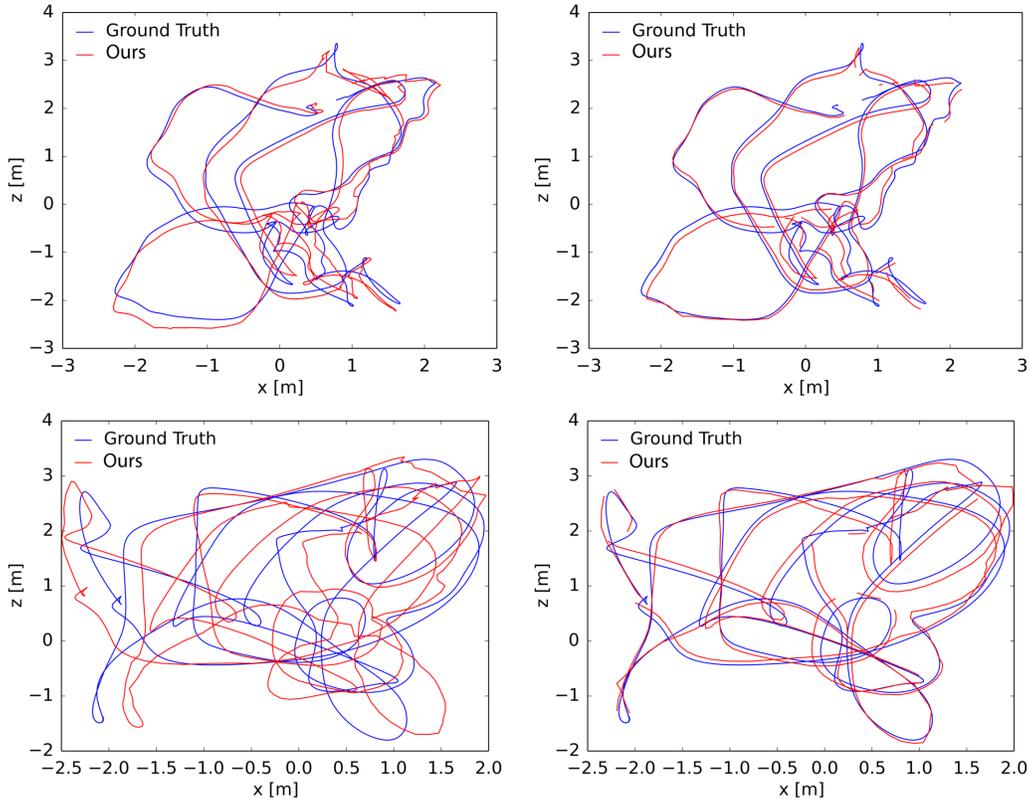


Figure 19: Comparison of Semi-Direct Odometry (left) and Semi-Direct SLAM (right) on the EuRoC dataset V1_01 and V1_02. With SLAM loop closures are found and accumulated drift is corrected yielding percentage improvements of 53.85% and 81.36% respectively

while the odometry drifts over time and does not retrieve the circular path, the SLAM extension closes the loop and continues the trajectory on the previously driven path.

Additionally, we also evaluate the performance of SLAM in comparison to pure odometry on our MAV. Similarly to above results, loop closures greatly help to reduce the drift on the datasets loop1 and loop2. While on dataset loop1 the odometry yields an estimate with 1.1 m ATE, the visual SLAM recovers the camera motion with 0.63 m. On dataset loop2 the odometry result improves from 2.16 m to 1.58 m when performing SLAM. The relative improvements on these datasets are 42.7% and 26.9% respectively.

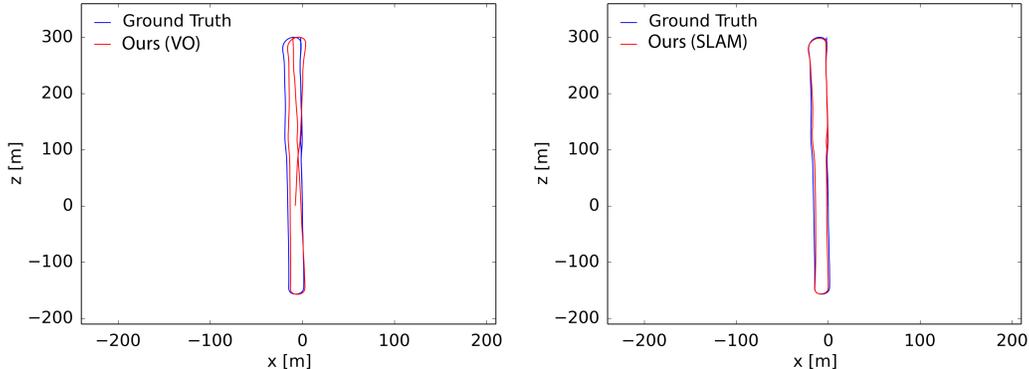


Figure 20: Comparison of Semi-Direct Odometry (left) and Semi-Direct SLAM (right) on KITTI Sequence 06. Example of a full loop closure found by our SLAM method while pure odometry drifts. SLAM achieves an improvement of 52.9%.

Dataset	Absolute Trajectory Error RMSE (Median) in m				
	Ours	Libviso2	LSD-SLAM	ORB-SLAM	S-PTAM
KITTI	10.28 (8.53)	32.54 (26.42)	X	41.49 (35.66)	25.74 (20.26)
EuRoC	0.38 (0.28)	0.85 (0.71)	0.53 (0.48)	1.23 (0.94)	1.82 (1.52)
MAV	0.80 (0.55)	1.78 (1.24)	0.38 (0.34)	0.75 (0.37)	X

Table 5: Average ATE Results on the different evaluated datasets

5. Accuracy

We have shown on different challenging datasets that in terms of accuracy we achieve similar results as current state-of-the-art stereo methods. The mean results for all datasets are summarized in Table 5. As can be seen in the table, our method achieves a lower ATE than the other evaluated methods on the KITTI and EuRoC datasets. On our MAV, monocular methods outperform the stereo methods. However, in comparison to the other stereo methods, our approach performs better and more robustly. Moreover, we

RPE	Ours (VO)	Libviso2	Direct VO
Translation Error (%)	0.8061	0.8449	0.8168
Rotation Error ($^{\circ} \text{m}^{-1}$)	0.0051	0.0052	0.0053

Table 6: Relative pose errors of the odometry methods. Translational drift is measured in percentage and rotational drift in $^{\circ} \text{m}^{-1}$.

Dataset	Method	Tracking	Mapping	Constraint Search	Optimization	Total (VO)	Total (SLAM)
KITTI	Ours	26.5 ms	36.6 ms	253.5 ms	564.6 ms	63.1 ms	881.2 ms
	LSD-SLAM	-	-	-	-	-	-
	ORB-SLAM	30.7 ms	254.0 ms	7.8 ms	1315.6 ms	284.6 ms	1608.0 ms
	S-PTAM	71.1 ms	5.7 ms	-	2036.9 ms	77.4 ms	2114.3 ms
	LIBVISO2	33.8 ms	-	-	-	33.8 ms	-
EuRoC	Ours	22.6 ms	39.6 ms	153.5 ms	684.2 ms	62.2 ms	899.9 ms
	LSD-SLAM	27.6 ms	85.6 ms	158.1 ms	207.3 ms	113.2 ms	478.5 ms
	ORB-SLAM	17.9 ms	159.2 ms	3.7 ms	535.6 ms	177.1 ms	716.4 ms
	S-PTAM	47.3 ms	1.5 ms	-	976.9 ms	48.8 ms	1025.7 ms
	LIBVISO2	24.8 ms	-	-	-	24.8 ms	-
MAV	Ours	17.5 ms	25.8 ms	140.0 ms	130.7 ms	43.3 ms	313.3 ms
	LSD-SLAM	28.7 ms	67.3 ms	314.0 ms	637.3 ms	79.0 ms	951.3 ms
	ORB-SLAM	24.3 ms	221.2 ms	11.0 ms	353.8 ms	245.5 ms	610.3 ms
	S-PTAM	-	-	-	-	-	-
	LIBVISO2	25.3 ms	-	-	-	25.3 ms	-

Table 7: Average runtimes of all evaluated methods

measure relative pose errors as proposed by Geiger et al. [35] to measure the performance and drift of pure odometry over large-scale sequences as in the KITTI dataset. Table 6 summarizes the results of our Semi-Direct Odometry in comparison to LIBVISO2 and Direct Odometry. Translational and rotational errors are measured separately. Results show that our method shows less translational and rotational drift over time. Moreover, as already seen above in the exemplary trajectory plots, the fully direct odometry has a higher rotational error than the other methods as direct alignment of frames becomes harder during large rotations.

In summary, our semi-direct approach shows accurate results for all datasets. Even on challenging fish eye stereo the whole trajectory can be retrieved and loop closures are found while S-PTAM fails to find any correspondences.

5.1. Runtime

For state-estimation with visual odometry or SLAM on mobile robots, real-time capability is an important factor. We thereby measure the efficiency of our method in terms of average runtime in ms.

We measure the average runtime as well as the runtime of the different blocks because it is oftentimes sufficient if tracking can be done with high frequency since global optimization usually does not run in real-time. The runtimes are broken down to the individual blocks: tracking, mapping, search for constraints and pose graph optimization. Timings for all datasets are listed in Table 7. Missing values are denoted with '-', e.g., S-PTAM does

not perform a constraint search as the other methods, and LIBVISO only performs tracking. The table clearly highlights that the SLAM parts, consisting of the constraint search and pose graph optimization, are the bottleneck for all systems.

In general, it can be seen that our approach is able to track incoming frames with 30 Hz. The mapping thread also runs in parallel to tracking at approximately 30 Hz. However, global optimization is still very costly for all methods. Especially in large-scale sequences the runtime rises.

5.2. Qualitative Results

A major advantage of our semi-direct approach is that 3D point clouds are estimated at runtime yielding an accurate semi-dense reconstruction of the environment. Thus, we are not only able to estimate the current pose of the camera but also maintain a 3D map of the environment which can be used for additional tasks like obstacle avoidance.

Exemplary qualitative results are shown for sequence 00 of the KITTI dataset. As can be seen in Figure 21, an accurate and consistent 3D reconstruction is achieved by Semi-Direct SLAM. For direct comparison to feature-based SLAM methods, the resulting sparse map built by ORB-SLAM is shown in Figure 22. While the reconstruction of ORB-SLAM only contains sparse points, our reconstruction allows detailed inference to existing objects in the scene. Most objects, that are visible in the camera image, can be recovered in our semi-dense map. For example, one can clearly distinguish between individual trees and cars. Contrarily, in the sparse map of ORB-SLAM one can only guess vaguely where the street is located.

Figure 23 shows the estimated pose graph of the camera trajectory and reconstructed map of the medium difficult EuRoC dataset V1_02. The images demonstrate that our estimated 3D reconstruction is globally consistent. The objects shown in the exemplary given camera image can easily be retrieved in the reconstructed map.

In conclusion, we state that our method builds globally consistent semi-dense 3D maps of the environment. It is well suited for large-scale sequences as in the KITTI dataset, as well as for smaller indoor sequences like the EuRoC dataset. We hence believe that the semi-dense 3D reconstruction yields a great benefit for autonomous visual navigation.

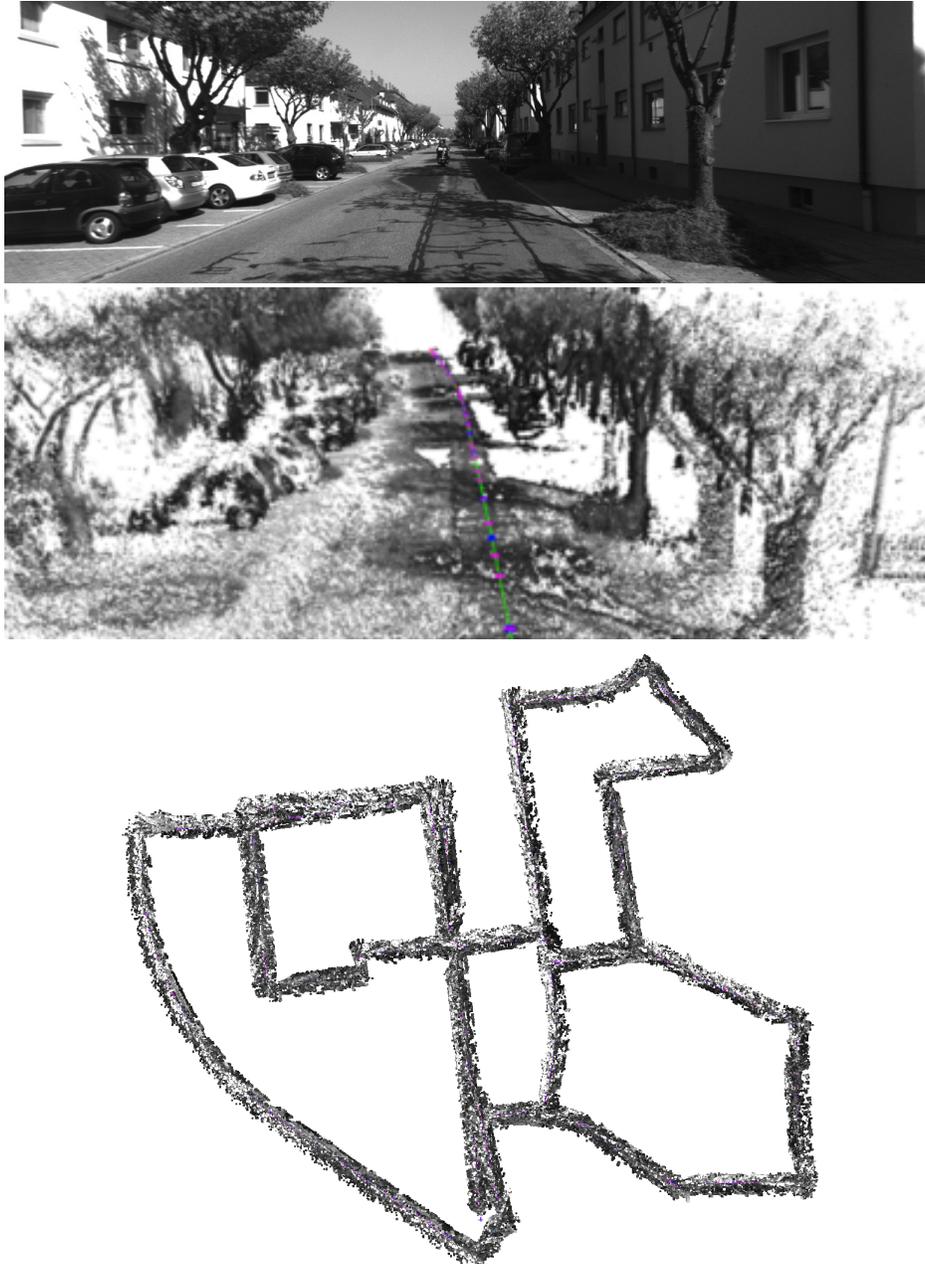


Figure 21: Semi-dense 3D Reconstruction of KITTI 00: The top image shows the reconstructed scene as captured by the camera. Below the semi-dense 3D reconstruction of this scene and the complete reconstruction of this dataset is shown.

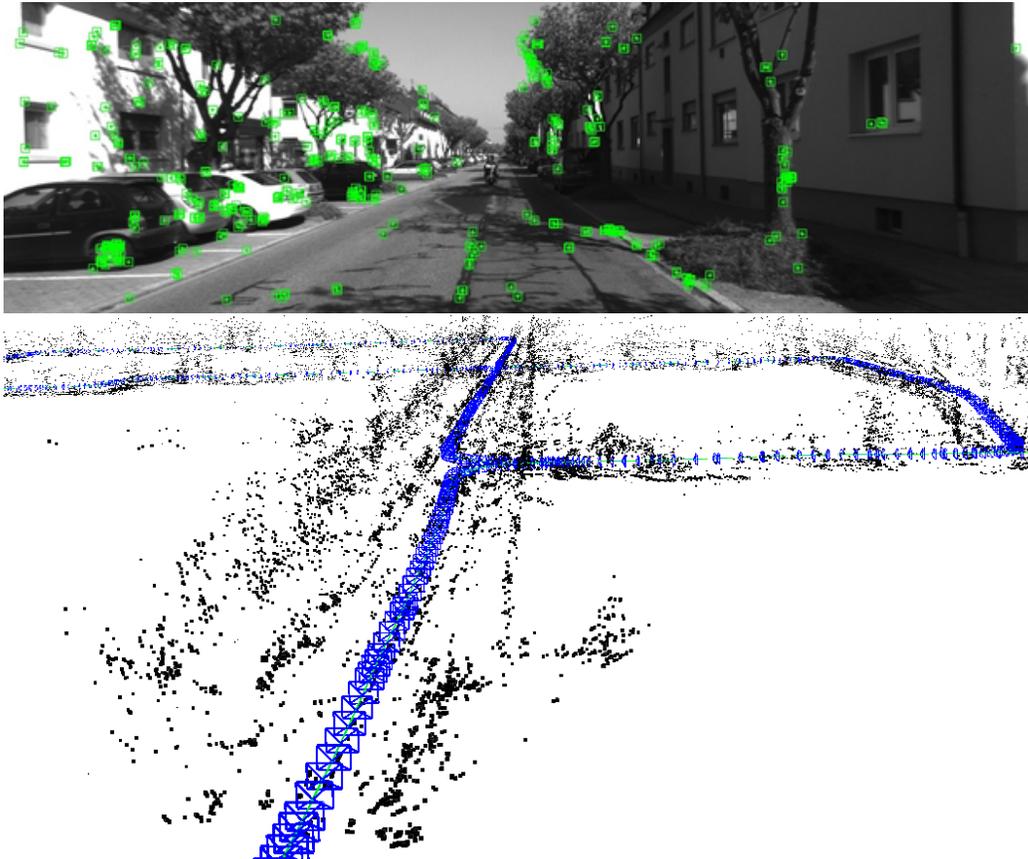


Figure 22: Sparse feature-based 3D Reconstruction of KITTI 00 by ORB-SLAM. The top view shows an exemplary scene where ORB features are tracked. The lower image shows the sparse map that is obtained by tracking ORB features.

6. Conclusions

In this paper, we proposed a novel hybrid visual odometry and SLAM method that combines feature-based tracking with semi-dense direct image alignment. Our method fuses depth estimates from motion between key frames with instantaneous stereo depth estimates.

The performance of our method has been evaluated in terms of accuracy, runtime, and scene reconstruction on three challenging datasets. Our experiments show that for tracking egomotion between image frames, we achieve accuracy similar to the state-of-the-art at high frame rate without the necessity to reduce the image resolution. Due to the feature-based tracking as

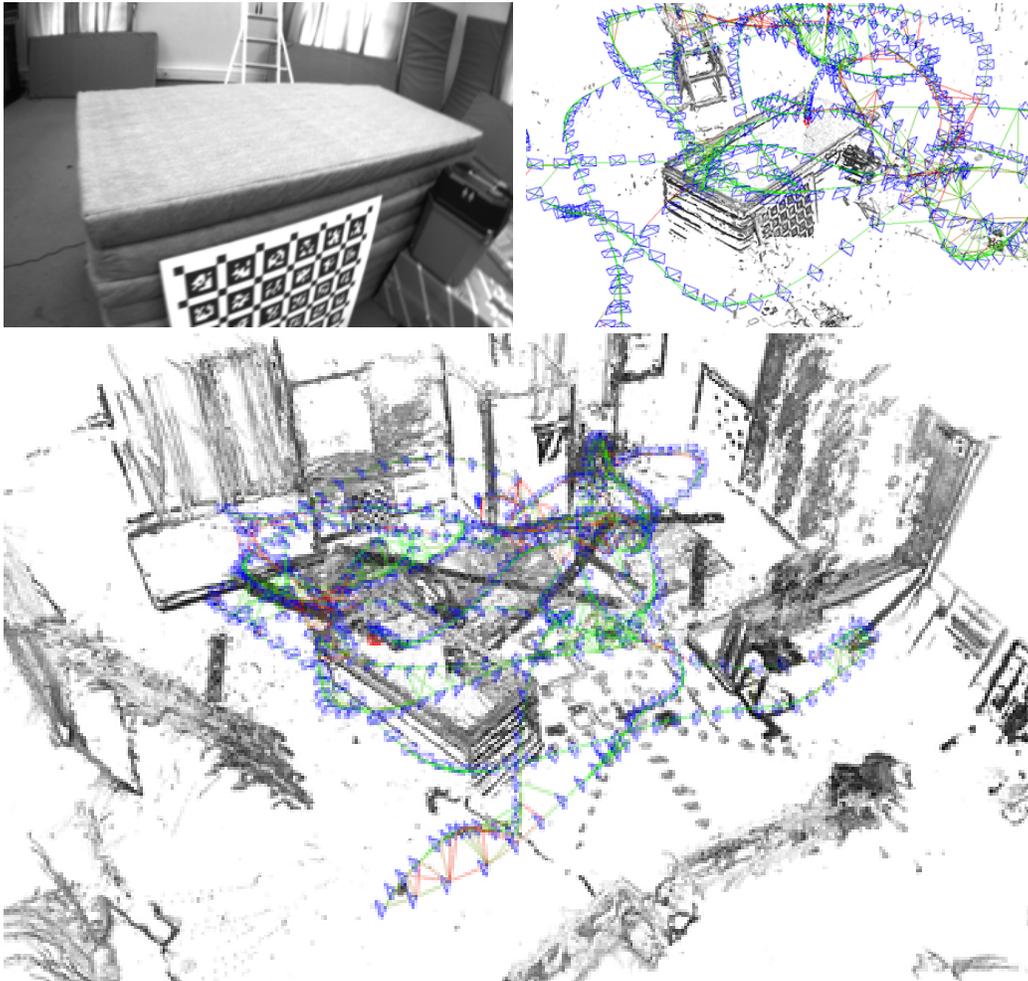


Figure 23: Semi-dense 3D Reconstruction of the EuRoC Dataset V1.02 with medium difficulty. Results show a globally consistent semi-dense map. The depicted key frame graph visualizes the trajectory. Key frames are shown in blue, while edge-constraints are shown in green and red, depending on their confidence.

prior for semi-dense direct alignment, our method is computationally less expensive and can estimate the relative camera motion in real-time. In future work, we plan to incorporate high frequency IMU readings and to evaluate other feature-based tracking priors, e.g. ORB features.

Acknowledgement

This work has been supported by the German Federal Ministry for Economic Affairs and Energy (BMWi) in the Autonomics for Industry 4.0 project InventAIRy.

References

- [1] J. Engel, T. Schöps, D. Cremers, LSD-SLAM: Large-scale direct monocular SLAM, in: European Conf. on Computer Vision (ECCV), 2014, pp. 834–849.
- [2] N. Krombach, D. Droschel, S. Behnke, Combining feature-based and direct methods for semi-dense real-time stereo visual odometry, in: Int. Conf. on Intelligent Autonomous Systems (IAS), 2016, pp. 855–868.
- [3] A. Davison, I. Reid, N. Molton, O. Stasse, Monoslam: Real-time single camera SLAM, *Pattern Analysis and Machine Intelligence* 29 (6) (2007) 1052–1067.
- [4] G. Klein, D. Murray, Parallel tracking and mapping for small AR workspaces, in: Int. Symposium on Mixed and Augmented Reality (ISMAR), 2007, pp. 225–234.
- [5] S. Weiss, D. Scaramuzza, R. Siegwart, Monocular-SLAM-based navigation for autonomous micro helicopters in GPS-denied environments, *Journal of Field Robotics* 28 (6) (2011) 854–874.
- [6] R. Mur-Artal, J. Montiel, J. D. Tardós, ORB-SLAM: A versatile and accurate monocular SLAM system, *Trans. on Robotics* 31 (5) (2015) 1147–1163.
- [7] A. Geiger, J. Ziegler, C. Stiller, Stereoscan: Dense 3D reconstruction in real-time, in: Intelligent Vehicles Symposium (IV), 2011, pp. 963–968.

- [8] M. Nieuwenhuisen, D. Droeschel, J. Schneider, D. Holz, T. Labe, S. Behnke, Multimodal obstacle detection and collision avoidance for micro aerial vehicles, in: European Conf. on Mobile Robots (ECMR), 2013, pp. 7–12.
- [9] R. Mur-Artal, J. D. Tardos, Orb-slam2: An open-source slam system for monocular, stereo, and rgb-d cameras, *IEEE Transactions on Robotics* 33 (5) (2017) 1255–1262.
- [10] T. Pire, T. Fischer, J. Civera, P. D. Cristoforis, J. J. Berlles, Stereo Parallel Tracking and Mapping for robot localization, in: Int. Conf. on Intelligent Robots and Systems (IROS), 2015, pp. 1373–1378.
- [11] T. Pire, T. Fischer, G. Castro, P. De Cristoforis, J. Civera, J. J. Berlles, S-ptam: Stereo parallel tracking and mapping, *Robotics and Autonomous Systems* 93 (2017) 27–42.
- [12] S. Houben, J. Quenzel, S. Behnke, Efficient multi-camera visual-inertial SLAM for micro aerial vehicles, in: Int. Conf. on Intelligent Robots and Systems (IROS), 2016, pp. 1616–1622.
- [13] M. Achtelik, M. Achtelik, S. Weiss, R. Siegwart, Onboard IMU and monocular vision based control for MAVs in unknown in- and outdoor environments, in: Int. Conf. on Robotics and Automation (ICRA), 2011, pp. 3056–3063.
- [14] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, P. Furgale, Keyframe-based visual-inertial odometry using nonlinear optimization, *Int. Journal of Robotics Research*.
- [15] C. Forster, L. Carlone, F. Dellaert, D. Scaramuzza, On-manifold preintegration for real-time visual-inertial odometry, *IEEE Transactions on Robotics* 33 (1) (2017) 1–21.
- [16] R. Mur-Artal, J. D. Tardos, Visual-inertial monocular slam with map reuse, *IEEE Robotics and Automation Letters* 2 (2) (2017) 796–803.
- [17] A. Comport, E. Malis, P. Rives, Accurate quadrifocal tracking for robust 3D visual odometry, in: Int. Conf. on Robotics and Automation (ICRA), 2007, pp. 40–45.

- [18] J. Engel, J. Sturm, D. Cremers, Semi-dense visual odometry for a monocular camera, in: *Computer Vision (ICCV), 2013 IEEE International Conference on*, IEEE, 2013, pp. 1449–1456.
- [19] J. Engel, J. Stückler, D. Cremers, Large-scale direct SLAM with stereo cameras, in: *Int. Conf. on Intelligent Robots and Systems (IROS)*, 2015, pp. 1935–1942.
- [20] J. Stückler, S. Behnke, Multi-resolution surfel maps for efficient dense 3d modeling and tracking, *Journal of Visual Communication and Image Representation* 25 (1) (2014) 137–147.
- [21] J. Stückler, A. Gutt, S. Behnke, Combining the strengths of sparse interest point and dense image registration for rgb-d odometry, in: *Int. Symposium on Robotics (ISR) and 8th German Conf. on Robotics (ROBOTIK)*, 2014, pp. 1–6.
- [22] R. A. Newcombe, A. Davison, Live dense reconstruction with a single moving camera, in: *Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2010, pp. 1498–1505.
- [23] M. Pizzoli, C. Forster, D. Scaramuzza, REMODE: Probabilistic, monocular dense reconstruction in real time, in: *Int. Conf. on Robotics and Automation (ICRA)*, 2014, pp. 2609–2616.
- [24] J. Engel, V. Koltun, D. Cremers, Direct sparse odometry, *Pattern Analysis and Machine Intelligence* 40 (3) (2018) 611–625.
- [25] T. Schöps, T. Sattler, C. Häne, M. Pollefeys, 3d modeling on the go: Interactive 3d reconstruction of large-scale scenes on mobile devices, in: *3D Vision (3DV), 2015 International Conference on*, IEEE, 2015, pp. 291–299.
- [26] D. Gallup, J.-M. Frahm, P. Mordohai, Q. Yang, M. Pollefeys, Real-time plane-sweeping stereo with multiple sweeping directions, in: *Conf. on Computer Vision and Pattern Recognition (CVPR)*, IEEE, 2007, pp. 1–8.
- [27] R. Mur-Artal, J. D. Tardós, Probabilistic semi-dense mapping from highly accurate feature-based monocular SLAM, in: *Robotics: Science and Systems*, 2015.

- [28] C. Forster, M. Pizzoli, D. Scaramuzza, SVO: Fast semi-direct monocular visual odometry, in: *Int. Conf. on Robotics and Automation (ICRA)*, 2014, pp. 15–22.
- [29] C. Forster, Z. Zhang, M. Gassner, M. Werlberger, D. Scaramuzza, Svo: Semidirect visual odometry for monocular and multicamera systems, *IEEE Transactions on Robotics* 33 (2) (2017) 249–265.
- [30] E. Piazza, A. Romanoni, M. Matteucci, Real-time cpu-based large-scale 3d mesh reconstruction, *arXiv preprint arXiv:1801.05230*.
- [31] G. Younes, D. Asmar, J. Zelek, Fdmo: Feature assisted direct monocular odometry, *arXiv preprint arXiv:1804.05422*.
- [32] A. Geiger, M. Roser, R. Urtasun, Efficient large-scale stereo matching, in: *Asian Conf. on Computer Vision (ACCV)*, 2010, pp. 25–38.
- [33] J. Engel, J. Sturm, D. Cremers, Semi-dense visual odometry for a monocular camera, in: *Int. Conf. on Computer Vision (ICCV)*, 2013, pp. 1449–1456.
- [34] R. Kümmerle, G. Grisetti, H. Strasdat, K. Konolige, W. Burgard, G2o: A general framework for graph optimization, in: *ICRA*, 2011, pp. 3607–3613.
- [35] A. Geiger, P. Lenz, R. Urtasun, Are we ready for autonomous driving? the KITTI vision benchmark suite, in: *Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2012, pp. 3354–3361.
- [36] M. Burri, J. Nikolic, P. Gohl, T. Schneider, J. Rehder, S. Omari, M. W. Achtelik, R. Siegwart, The EuRoC micro aerial vehicle datasets, *The Int. Journal of Robotics Research*.
- [37] J. Sturm, N. Engelhard, F. Endres, W. Burgard, D. Cremers, A benchmark for the evaluation of RGB-D SLAM systems, in: *Int. Conf. on Intelligent Robots and Systems (IROS)*, 2012, pp. 573–580.
- [38] S. Kammel, J. Ziegler, B. Pitzer, M. Werling, T. Gindele, D. Jagzent, J. Schrder, M. Thuy, M. Goebel, F. v. Hundelshausen, O. Pink, C. Frese, C. Stiller, Team annieway’s autonomous system for the 2007 darpa urban challenge, *Journal of Field Robotics* 25 (9) (2008) 615–639.

- [39] M. Irani, P. Anandan, About direct methods, in: Int. Workshop on Vision Algorithms: Theory and Practice, Int. Conf. on Computer Vision (ICCV), Springer-Verlag, London, UK, UK, 2000, pp. 267–277.
- [40] J. Nikolic, J. Rehder, M. Burri, P. Gohl, S. Leutenegger, P. T. Furgale, R. Siegwart, A synchronized visual-inertial sensor system with FPGA pre-processing for accurate real-time SLAM, in: Int. Conf. on Robotics and Automation (ICRA), 2014, pp. 431–437.
- [41] M. Beul, N. Krombach, M. Nieuwenhuisen, D. Droeschel, S. Behnke, Autonomous Navigation in a Warehouse with a Cognitive Micro Aerial Vehicle, Springer International Publishing, Cham, 2017, pp. 487–524.
doi:10.1007/978-3-319-54927-9_15.
URL https://doi.org/10.1007/978-3-319-54927-9_15
- [42] D. Droeschel, M. Nieuwenhuisen, M. Beul, D. Holz, J. Stückler, S. Behnke, Multilayered mapping and navigation for autonomous micro aerial vehicles, *Journal of Field Robotics* 33 (4) (2016) 451–475.