# Probabilistic RGB-D Odometry based on Points, Lines and Planes Under Depth Uncertainty

Pedro F. Proença*, Yang Gao

*Surrey Space Centre, Faculty of Engineering and Physical Sciences, University of Surrey, GU2 7XH Guildford, U.K*

## Abstract

This work proposes a robust visual odometry method for structured environments that combines point features with line and plane segments, extracted through an RGB-D camera. Noisy depth maps are processed by a probabilistic depth fusion framework based on Mixtures of Gaussians to denoise and derive the depth uncertainty, which is then propagated throughout the visual odometry pipeline. Probabilistic 3D plane and line fitting solutions are used to model the uncertainties of the feature parameters and pose is estimated by combining the three types of primitives based on their uncertainties.

Performance evaluation on RGB-D sequences collected in this work and two public RGB-D datasets: TUM and ICL-NUIM show the benefit of using the proposed depth fusion framework and combining the three feature-types, particularly in scenes with low-textured surfaces, dynamic objects and missing depth measurements.

*Keywords:* Feature-based Visual Odometry, Probabilistic Plane and Line Extraction, Depth Fusion, Depth Uncertainty, Structured Environments

## 1. Introduction

Point, line and plane primitives allow a minimalistic, yet comprehensive representation of structured environments, which is more appealing than dense representations [1], in terms of efficiency. While feature points can be insufficient for visual odometry in low textured environments, combining them with planes

---

*Corresponding author

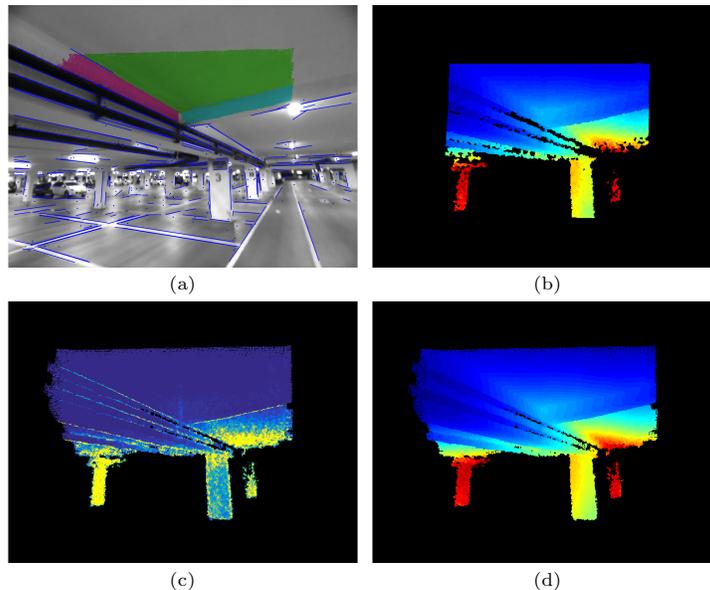*Email address:* p.proenca@surrey.ac.uk ( Pedro F. Proença )

Figure 1: Raw and processed RGB-D frame by our system, in a challenging environment with low textured areas and missing depth measurements from a structured-light sensor. (a) Detected features overlaid on the intensity image. (b) Raw depth map. (c) Depth uncertainty estimated by the proposed depth filter and (d) the respective fused depth map. To expand the camera FOV, a wide-angle lens was mounted on the RGB camera while the depth fusion allows propagation of depth measurements beyond the narrow FOV of the depth camera.

and line segment features may lead to more robustness to plain planar surfaces [2, 3, 4], blur caused by sudden motion [2] and light variations [5]. Therefore, this work proposes a feature-based odometry method that combines points, lines and planes for visual odometry by relying on an RGB-D camera to capture densely the scene geometry and texture.

However, the depth measurements captured by active depth sensors is affected by significant error [6, 7], which in turn affects the estimation of 3D feature parameters. In particular, lines tend to be detected on depth discontinuities (see Fig. 1), where noise is more severe. Moreover, features may have missing depth measurements due to the range and field-of-view (FOV) limitations of these cameras. Thus, we propose a depth fusion framework to: (i) denoise the raw depth map, (ii) model the depth uncertainty and (iii) recover temporally missing depth measurements. Since the effective depth error depends on the scene properties, besides considering the systematic depth sensor error, the proposed framework captures the observed uncertainty by assessing the spatial and tem-

poral distribution of depth measurements. This uncertainty is then propagated throughout the visual odometry pipeline. Specifically, probabilistic 3D line and plane fitting solutions, based on weighted linear least squares, are used to model the uncertainty of these primitives and then pose is estimated by taking into account these uncertainties. The motivation for representing the uncertainties of these primitives, is that their impact on the pose estimation should depend on the precision of their estimated parameters, which depend on the number, uncertainties and distribution of their samples. The key contributions of this paper are the following:

- A probabilistic depth fusion framework based on Mixture of Gaussians that models depth uncertainty. The code is available as open-source[1].

- Extend our recently developed visual odometry method [2] based on points and planes to line segments.

- A probabilistic analytical solution to 3D line fitting.

- Evaluate the system on public and author-collected[2] RGB-D datasets, and demonstrate the benefit of modelling temporally depth uncertainty and combining points, planes and lines in low textured and dynamic environments.

The remainder of this paper is organized as follows. Section 2 reviews related work focusing on methods that use depth cameras and features (e.g. points, planes and lines). Section 3 and Section 4 describe respectively our depth filter framework and visual odometry method. Our results are reported and discussed in Section 5. Finally, Section 6 concludes and discusses the limitations of the approach.

## 2. Related work

Active depth sensors have been widely adopted in many computer vision tasks, e.g., reconstruction, segmentation, egomotion estimation, object recognition, human pose estimation and scene understanding [8]. However, extensive

---

[1]https://github.com/pedropro/OMG_Depth_Fusion

[2]A video of the experiments is available at: https://youtu.be/YOT2_ghlng0

analysis [6, 9, 7, 10, 11] of these consumer-grade sensors have revealed their limitations. It is now well known that structured-light sensors, used by the first version of Kinect cameras, suffer from severe quantization and consequently the depth error grows quadratically with the distance to the sensor. A theoretical model for this source of error has been derived in [6]. Later, time-of-flight (ToF) sensors, used by the second version of Kinect cameras, have been empirically analyzed and compared to the first version in [7, 10]. Although, the ToF depth error proved to be significantly less affected by the distance to the sensor, ToF sensors suffer from other sources of error: flying pixels arising from depth descontinuities, non-Lambertian surfaces (e.g. black surfaces), multipath interference and the depth error grows with the image distance to the principal point. Besides these sensor-specific issues, depth sensors, in general, have limited range and FOV compared to LIDAR sensors, as shown in Fig. 1.

KinectFusion [1] was the first work to reconstruct dense models from the noisy and incomplete depth maps captured by these sensors. This system uses raw depth maps to update a global volumetric model based on cumulative moving average updates of voxel states, which are represented as Truncated Signed Distance Functions (TSDF), while pose is estimated by using Iterative Closest Point (ICP) algorithm. To perform ICP, the model is raycasted and the depth maps are first denoised by a bilateral filter [12]. Since then, several works [11, 13, 14] have extended KinectFusion to achieve better quality reconstructions. In [11], a depth noise model that takes into account both the sensor lateral and axial noise, was empirically derived and incorporated into the KinectFusion pipeline. Specifically, the depth uncertainty was used to weight the ICP and the voxel TSDF updates. Due to the GPU memory requirements and voxel discretization of these volumetric methods, a selective point-based fusion method was instead proposed in [13] to reconstruct denoised 3D models in dynamic environments. More recently, temporal depth map fusion has been used to denoise depth maps either by using the median [14] or the moving average [15]. However, in these works, depth uncertainty is neither modelled nor explicitly used for fusion.

The dense RGB-D odometry method, termed DVO [16], which is based on the minimization of the photometric and geometric error, in [17], has also been improved, in [18, 19], by considering the depth error. [18] proposed using the inverse depth to parameterize the geometric error, whereas [19] proposed using

the image derivatives to weight the residuals of the minimization problem.

In feature-based SLAM methods, to address the low depth resolution of structured light cameras, specific sensor depth uncertainty models [6, 9] were adopted for point and line odometry [20] and for our recently proposed point and plane odometry [2]. Moreover, [21] proposed using a Mixture of Gaussians convolution to assess the uncertainty of a single depth map. Such framework represents the uncertainty around the object edges better than the sensor error models, proposed in [6, 11], thus our proposed depth filter builds on this framework. The resulting uncertainty was further used, in [21], to update a sparse model of feature points, through Kalman filter correction equations. An experimental study comparing dense vs. feature point based VO and ICP variations was performed in [22] and showed no clear winner since their relative performance depends on the particular environment characteristics.

Beside features points, line primitives are becoming increasingly popular in monocular [3, 4, 23] and RGB-D Odometry [5, 20, 24]. A non-linear 3D line fitting was proposed in [5] to fit depth measurement samples and their uncertainties, however it is not efficient to cast each line fitting as an iterative problem, considering the typical high number of detected 2D lines (e.g. 100). This problem was more recently simplified in [3], where an analytical method was devised for monocular odometry by exploiting the fact that a 3D line is projected as a plane. While, this solution is more appealing than the previous one, pixel samples from a 2D line may not lie all on the same plane, e.g., the pixels crossed by an oblique line. Moreover, one may wish to sample depth measurements from the 2D line neighbourhood, due to a lack of measurements, or fit a 3D line to measurements obtained from multiple line observations. Therefore, the analytical 3D line fitting solution, proposed here, is more general, as it supports all these cases. Basic 3D line fitting based on PCA was employed in [24] but this neglects the depth uncertainties.

Plane primitives have also been widely exploited by SLAM methods: A planar method was proposed in [25] to use data from both a 2D LIDAR and a depth camera as these complement each other in terms of FOV and operating range. Points and planes were initially combined by a SLAM system, in [26], to avoid the geometric degeneracy of planes. The system used a RANSAC framework for mixed 3D registration of both point-to-point and plane-to-plane matches by sampling any triplets formed by these matches. To cope with miss-

ing depth measurements, this framework was later extended in [27] to include also 2D-to-3D point matches as a triplet hypothesis and 2D-to-2D matches for checking the hypothesis consensus. To reduce the computational cost of plane extraction, plane tracking was proposed in [28], however faster plane extraction algorithms have been recently developed [29, 30]. Alternatively, in [31], feature points have been enhanced with planar patches, for a small overhead, to improve feature matching and increase the constraints imposed by feature points on pose estimation, but nevertheless this approach does not exploit featureless planar patches. On the contrary, in [32], the Direct SLAM method [17] was combined with global plane model tracking through an EM framework to reduce drift. Later, a more efficient alternative was developed in [15], without using GPU. The uncertainty in plane extraction was analyzed thoroughly in [33] by comparing direct and iterative plane fitting methods, in terms of accuracy and speed. In [2], we proposed an RGB-D Odometry method for points and planes that modelled and propagated the depth uncertainty throughout the system pipeline. Here, we extend this approach to lines and propose a better depth model.

## 3. Probabilistic Depth Filter

The proposed depth filter, outlined in Fig. 2, can be split into three stages: (i) Given the raw depth map of the current frame, depth uncertainty is assessed according to a specific sensor model, (ii) Based on this depth uncertainty, the raw depth map is convolved with the Gaussian Mixture (GM) kernel proposed in [21], to capture the uncertainty within the pixels neighbourhood, (iii) The depth estimates and uncertainties resulting from this GM convolution are then combined with estimates from past frames by using our proposed Optimal-GM fusion method. In order to do so, depth estimates from a sliding window of frames are maintained and updated as 3D measurements by a point cloud registration module. A detailed explanation of these modules is given in the following subsections.

### 3.1. Depth Sensor Error Model

In our experimental work, we have used structured-light depth sensors based on active stereo, which suffer inherently from disparity quantization, therefore, we adopted the theoretical error model of [6], which accounts for the propagation
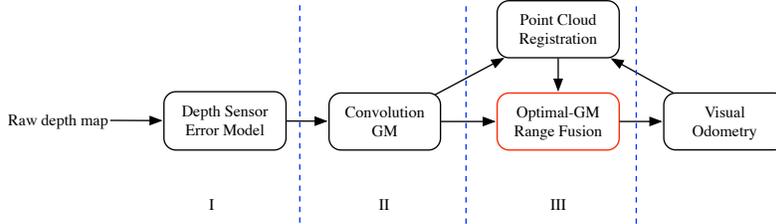
6

Figure 2: Overview of the proposed probabilistic depth filter

of disparity random error $\sigma_d$ in Kinect structured-light sensors, so that the depth uncertainty is given by: $\sigma_z = \sigma_d(\frac{m}{fb})Z^2$ where $f$ is the focal length, $b$ is the baseline between the IR projector and camera and $m$ is a normalizing parameter. Setting $\sigma_d = 0.5$, as in [6], yields the following expression:

$$\sigma_z = 1.425 \times 10^{-6} z^2 \quad [mm] \tag{1}$$

which fits well the planar residuals in [6] and is consistent with the axial noise in [11]. However, this simple expression does not comprehend many other sources of depth error, e.g., lateral noise [11], ambient background light and temperature drift [7]. Although, a more comprehensive model could be developed, it is extremely difficult to model the actual depth error, since this depends also on the properties of the observed object surfaces. Thus, we refrain from doing so and instead look at the spatial and temporal distribution of depth samples, through the next consecutive modules.

### 3.2. Convolution of Gaussian Mixtures

To address the lateral error, Dryanovski et al. [21] proposed to quantify the uncertainty of depth pixels based on the depth values of their image neighbourhoods through a GM formulation. Let the probability density function of depth in a $3 \times 3$ local window, centered at pixel $p$, be given by the following $N$ mixture of Gaussians:

$$f(z) = \frac{1}{S} \sum_{i=1}^{N} w_i \mathcal{N}(z_i, \sigma_{z_i}^2) \tag{2}$$

where each Gaussian corresponds to a pixel of the local window with a variance given by the depth sensor error model, $S$ is a normalizing constant and the weights $w_i$ are assigned to the local window according to the following kernel:

$$W = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix} \tag{3}$$
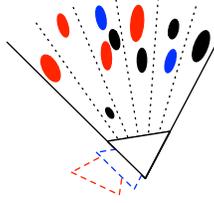
7

Figure 3: Illustration of the problem of fusing point measurements with their uncertainties from three frames

Then, the new estimated depth for $p$ takes the value of the mean of (2):

$$\bar{z} = \frac{1}{S} \sum_{i=1}^{N} w_i z_i \tag{4}$$

and the GM uncertainty can be found by expressing the variance in terms of moments:

$$\mathrm{Var}(f(z)) = \frac{1}{S} \sum_{i=1}^{N} w_i(z_i^2 + \sigma_{z_i}^2) - \bar{z}^2 \tag{5}$$

To account for pixels with missing depth values, we simply represent their depth and uncertainty as 0 and flag them with an indicator function $y$, such that the normalizing constant is given by:

$$S = \sum_{i=1}^{N} w_i y_i \tag{6}$$

Effectively, the resulting variance allows assigning high uncertainty to outliers (e.g. flying pixels) and depth discontinuity locations. One motivation for the latter, is that the 2D coordinates of features detected on the RGB images are also subject to error and moreover RGB images may not be perfectly aligned with the depth map due to errors in the extrinsic calibration and temporal synchronization, consequently image features corresponding to foreground may be associated to background. This information should be taken into account during both the pose estimation and the temporal fusion to reduce the impact of wrong depth associations.

### 3.3. Optimal Gaussian Mixture for Temporal Fusion

Given the 3D measurements of past frames, which are transformed to the current frame by the point cloud registration, these are projected to the image

grid, as illustrated in Fig. 3, such that each pixel will have a set of points with their respective uncertainties, estimated by the GM convolutions. Due to the transformation between frames, ideally the $3 \times 3$ covariances of the points should be rotated as well, however this is computationally expensive, thus we work on the range space which is invariant to camera rotation unlike the depth. Specifically, the depth uncertainties given by the GM convolutions are propagated to range uncertainties and then stored, as follows:

$$\sigma_r^2 = \frac{\sigma_z^2}{\cos^2 \alpha} \tag{7}$$

where $\alpha$ is the angle of incidence of the projection ray on the image plane and a matrix of cosines, for all pixels, can be pre-computed according to the camera intrinsic parameters. Once the points are converted to range as well, the previous GM framework can be applied to obtain a new range estimate for each pixel, although in this case we use the range uncertainties to weight the GM, such that the resulting range is:

$$\bar{r} = \frac{1}{\sum_{i=1}^{M} \sigma_{r_i}^{-2}} \sum_{i=1}^{M} \frac{r_i}{\sigma_{r_i}^2} \tag{8}$$

for $M$ projected points. This expression corresponds in fact to the Maximum Likelihood Estimator (MLE) [34], but more notably because this framework does not make independence assumptions, the uncertainty given by the expression (5) takes into account both the intra-group and inter-group variances, unlike the least squares formulation. Thus, pixels with inconsistent range measurements will have high uncertainty and the gross errors, e.g., flying pixels, that were *a-priori* modelled by the GM convolutions will be penalized during the range fusion by the weighting function in (8).

The fused range image and uncertainty can then be converted back to depth using: $z = r \cos \alpha$ and (7). As can be seen in Fig. 4 and 5, this framework removes significant noise from the raw depth maps. Additionally, one can see in Fig. 5 that the depth uncertainty represents well the respective depth error in a synthetic dataset, though our sensor model differs from the one used to generate the depth noise in [35].

### 3.4. Depth Fusion Constraints

Temporal fusion assumes that the transformations (i.e. stereo poses) between the frames of the sliding window are sufficiently good to fuse measure-
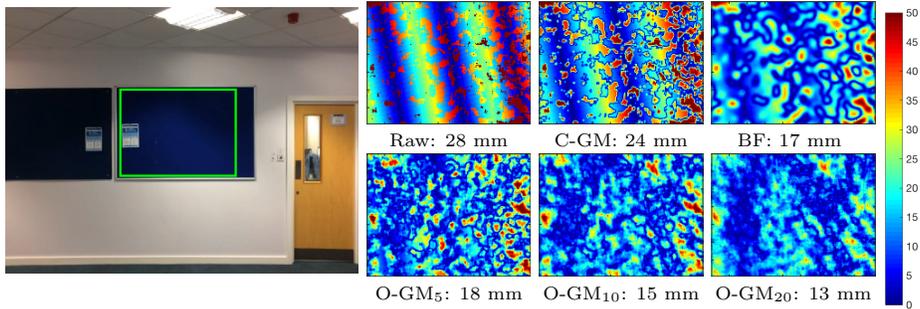
Figure 4: Planar fitting residuals for three depth filtering methods: C-GM is the method described in Section 3.2; BF is the widely used bilateral filter [12] (e.g. used in KinectFusion [1]) and O-GM$_n$ is our temporal fusion method with measurements from $n$ frames. Maps of point-to-plane distances and the respective RMSE are depicted for the region highlighted on the left image. O-GM$_n$ and BF are able to remove two types of error, which are revealed in the raw depth map: random errors and quantization errors (appearing as stripes).

ments from the same point in space. Therefore, besides using the pose estimated by the visual odometry to bring the registered point cloud to the current frame, the uncertainty of the transformations is monitored using the method described in Section 4.6.1. If the uncertainty of a transformation exceeds a given threshold, measurements from the respective frame are removed from the registered point cloud. As a result, the length of the sliding window of frames is dynamic.

Furthermore, the temporal fusion is not intended for long durations and wide baselines due to: memory and computation time requirements, dynamic objects and occlusions. Parallax, due to camera translation, causes incorrect fusion of measurements from occluded background with more recent foreground measurements. To reject measurements from occlusions, we enforce a consistency constraint to old depth measurements, i.e., during the range fusion, point are projected from newest to oldest, if a pixel receives at least $k = 5$ points, more points are only accepted if their ranges are within the margin: $\bar{r} \pm 3\sigma_r$, where $\bar{r}$ and $\sigma_r$ correspond to the current pixel state of range and uncertainty.

Fig. 6 shows the behaviour of the filter in a dynamic environment. Although dynamic objects could be addressed by segmentation, as in [13, 37], our temporal fusion framework already assigns high uncertainty to the depth values that are affected by the motion of moving objects, which implicitly will downweight the features arising from the moving objects. Our results, in Section 5, support this idea.
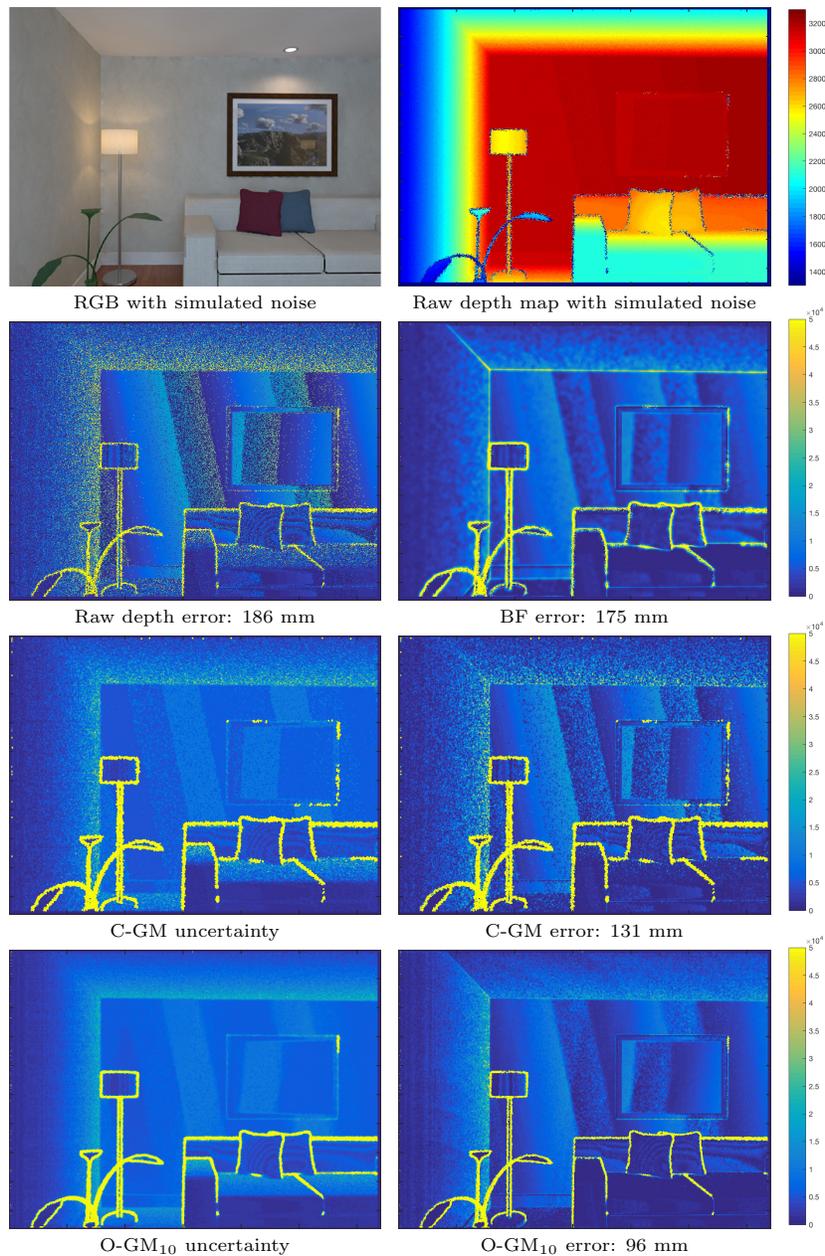
10

Figure 5: Depth error and uncertainty on one frame from the synthetic ICL-NUIM dataset [35], with simulated noise, along with the respective RMSE for each depth filter. Depth error is measured by checking the available noiseless version of the depth maps. Despite the ability of the bilateral filter (BF) to preserve edges, errors are still introduced on the room edges. Notice how some of the errors on the picture frame and on the sofa are removed after applying the O-GM fusion.

Raw depth map       Fused depth map

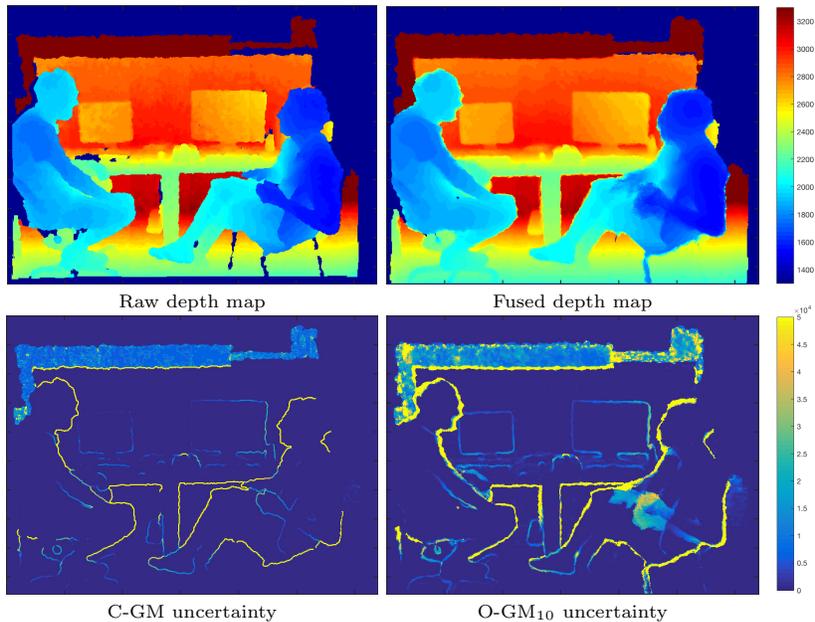C-GM uncertainty       O-GM$_{10}$ uncertainty

Figure 6: Depth fusion and uncertainty in a dynamic scene, captured in [36], where two people talk and gesticulate. Although depth fusion degrades the depth map around the human contours, these errors are captured by the uncertainty model (notice the hands), which has the benefit, for pose estimation, of assigning lower weight to any features detected in these regions.

## 4. RGB-D Odometry based on Points, Planes and Lines

The proposed visual odometry method, outlined in Fig. 8, starts by detecting points, lines and planes from the current RGB-D frame. While 2D points and lines, along with their feature descriptors, are extracted from the intensity of the RGB channel, planes are extracted from an organized point cloud back-projected from the depth map, after applying the first and second stages of the depth filter (i.e. the depth sensor error model and the GM convolution) in order to obtain the 3D point uncertainties, which are then used by a weighted least squares plane fitting. The extracted primitives are then matched against the ones extracted from the previous frame. Resulting 3D-to-2D point and line matches and 3D-to-3D plane matches are subsequently used jointly to estimate the frame-to-frame pose, according to their uncertainties. Once the pose is estimated, a depth fused map is obtained, using the third stage of the depth filter, described in Section 3. Given this new depth map, the 3D coordinates of the current point and

line features are finally obtained for the next frame-to-frame pose estimation through backprojection and a weighted 3D line fitting method, which also takes into account the depth uncertainty. Furthermore, plane detection and fitting is repeated to obtain presumably better plane estimates. These modules are described in further detail below.
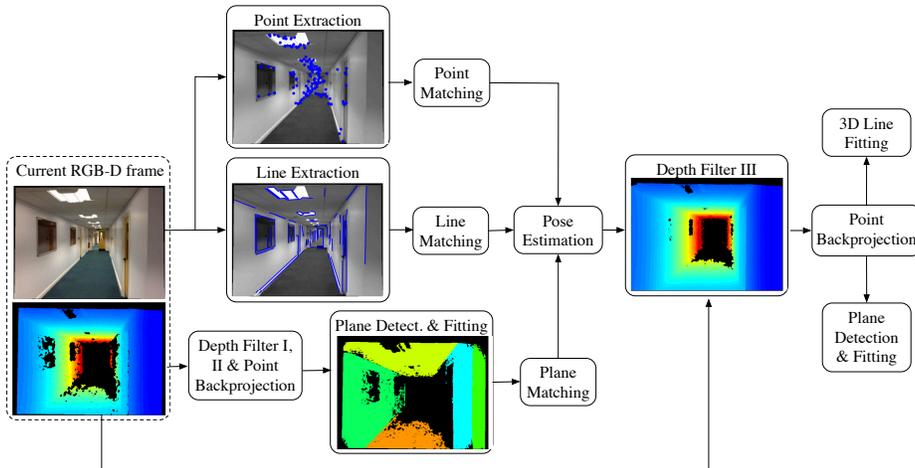


Figure 7: Visual odometry system overview

### 4.1. Extraction of Points, Lines and Planes

Image points are detected by relying on SURF features, whereas for lines, the LSD [38] method is used to detected line segment endpoints and then binarized LBD [39] descriptors (implemented in OpenCV) are extracted from the respective lines. For plane extraction, once the depth map is backprojected, we make use of the method proposed in [29], which processes efficiently organized point clouds in real-time, the result is a segmented point cloud (as depicted in Fig. 8). For each point cloud segment, a plane model is fit and its uncertainty is derived using the method described in Section 4.3.

### 4.2. Point Backprojection

Assuming that the depth image is mapped to the RGB reference frame, given the extrinsic calibration, the 3D coordinates $P = \begin{bmatrix} X, Y, Z \end{bmatrix}^{\top}$ corresponding to a pixel $p = \{u, v\}$ on either depth or RGB image can be obtained through

backprojection:

$$P = Z \begin{bmatrix} (u - c_x)/f_x \\ (v - c_y)/f_y \\ 1 \end{bmatrix} \tag{9}$$

where $Z$ is the value of the depth pixel, and $\{f_x, f_y\}$ and $\{c_x, c_y\}$ are respectively the focal length and principal point of the RGB camera. The uncertainty of $P$ can be obtained by the first order error propagation:

$$\Sigma_P = J_P \begin{bmatrix} \Sigma_p & 0 \\ 0 & \sigma_Z^2 \end{bmatrix} J_P^\top \tag{10}$$

where $J_P$ is the Jacobian of (9) with respect to $p$ and $Z$, $\sigma_Z^2$ is the uncertainty of the depth value given by the depth filter and $\Sigma_p$ is a $2 \times 2$ identity matrix times the pixel coordinate uncertainty $\sigma_p^2$, which accounts for the pixel quantization error. Let this error be modelled by a uniform PDF of length equal to 1 pixel, then its variance is $\sigma_p^2 = 1/12$.

### 4.3. WLS Plane Fitting

For plane fitting, we employ our recently proposed weighted least squares method [2]. For the sake of completeness, we describe here the method and then propose a modification to derive a more accurate plane uncertainty.

It is efficient to express planes as infinite planes in the Hessian normal form: $\theta = \{N_x, N_y, N_z, d\}$. However, such representation is overparameterized, thus the estimation of these parameters by unconstrained linear least squares is degenerate. This issue has been solved in [33] by using constrained optimization and in [40] by using a minimal plane parameterization. Similarly to [40], we use a minimal plane representation: $\theta_m = \left[ N_x, N_y, N_z \right]/d$, as an intermediate parameterization. Since, a plane with $d = 0$ implies detecting a plane that passes through the camera center (i.e. projected as a line), it is safe to use this parameterization. The new parameters are then estimated by minimizing the point-to-plane distances through the following weighted least-squares problem:

$$E = \sum_{i=1}^{n} \frac{w_i (\theta_m P_i + 1)^2}{2} \tag{11}$$

where the scaling weights were chosen to be the inverse of the point depth uncertainties: $w_i = \sigma_{Z_i}^{-2}$, which represent well the point-to-plane distance uncertainties when the detected plane is approximately parallel to the image plane.

By setting the derivative of (11), with respect to $\theta_m$, to zero, we arrive at the solution of the form: $\theta_m^\top = A^{-1}b$, where $A = \sum_{i=1}^n w_i P_i P_i^\top$ and $b = -\sum_{i=1}^n w_i P_i$.

Following the Fisher observed information [34], the covariance of $\theta_m$ is given by the inverse Hessian matrix of $E$, i.e., $\Sigma_{\theta_m} = H^{-1}$ where $H$ is simply $A$. However, the residuals are scaled by an heuristic choice of weights, and as a result $E$ is a just a scaled approximation of the negative log-likelihood function. This fact was neglected in [2], and as a result the uncertainty was overestimated. Therefore the weights need to be first updated with the actual uncertainty of the plane residuals $\Sigma_{r_i}$. These can be found at the solution $\theta_m$ by propagating the point uncertainties $\Sigma_{P_i}$ as follows: $\Sigma_{r_i} = \theta_m \Sigma_{P_i} \theta_m^\top$. The uncertainty of $\theta_m$ is then derived from the updated matrix $A$. Finally, the Hessian normal form can be recovered by:

$$\theta = \frac{\begin{bmatrix} \theta_m & 1 \end{bmatrix}}{\|\theta_m\|} \tag{12}$$

and the respective uncertainty is obtained via first order error propagation: $\Sigma_\theta = J_\theta \Sigma_{\theta_m} J_\theta^\top$, where $J_\theta$ is the Jacobian of (12).

### 4.4. WLS Line Fitting

The proposed solution to 3D line fitting is illustrated in Fig. 8. First, as in [20], depth pixels are sampled uniformly across the 2D line segments (maximum 100 pixels per line). The pixels with available depth are backprojected to 3D points and then these are processed by a RANSAC loop based on 3D point-to-line Euclidean distances to remove outliers. We believe that, in this work, the Euclidean distance is more adequate to remove outliers than the Mahalanobis metric proposed in [20], due to the high uncertainties given by the GM convolution at depth discontinuities (see Fig. 6). The final consensus set of 3D points is denoted as $P = \{P_1, ..., P_n\}$.

The problem of fitting a 3D line to 3D points can be solved non-iteratively by casting it as 2D vector estimation problem. The key idea is to exploit the fact that the optimal line passes in the center of mass, denoted as $O$, by estimating a line pinpointed at $O$, as depicted in Fig. 8. The centroid $O$ corresponds to the MLE: $(\sum_{i=1}^n W_i)^{-1} \sum_{i=1}^n W_i P_i$ where $W_i = \Sigma_{P_i}^{-1}$ and the MLE variance is $(\sum_{i=1}^n W_i)^{-1}$. But for efficiency, we instead approximate $O$ as the mean of the $n$ points weighted by the inverse of their depth variances, in order to avoid inverting the $n$ covariance matrices, required by the MLE. This approximation
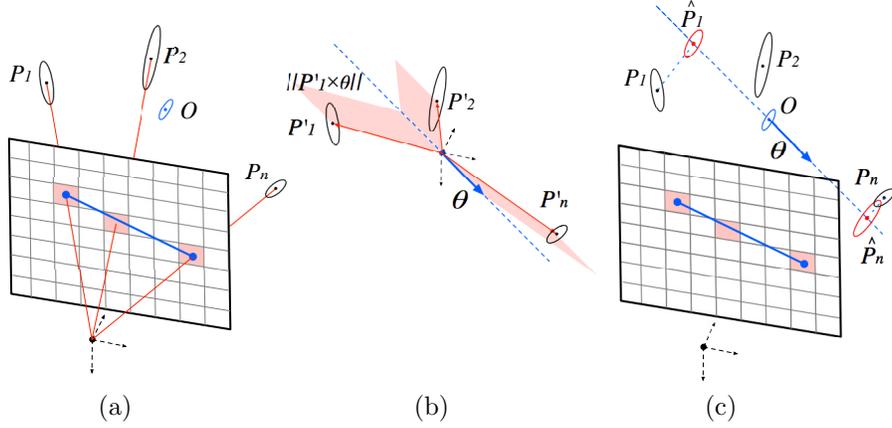
Figure 8: Illustration of the solution proposed to 3D line fitting. (a) After sampling and backprojecting depth pixels across the 2D line, the center of mass $O$ is determined. (b) Points are translated so $O$ is in the origin, and the vector $\theta$ is estimated by minimizing the sum of the shaded areas, which correspond to the cross product norms. (c) Line endpoints are selected by projecting the extreme points $P_1$ and $P_n$ on $\theta$.

implies a covariance isotropic assumption. Therefore, $O = \sum_{i=1}^{n} w_i P_i$ where $w_i$ is a normalized weight: $\sigma_{Z_i}^{-2} / \sum_{j=1}^{n} \sigma_{Z_j}^{-2}$ and the respective covariance is $\Sigma_O = \sum_{i=1}^{n} w_i^2 \Sigma_{P_i}$.

The point samples are then translated by subtracting $O$: $P_i' = P_i - O$ in order to estimate a vector $\theta$ by minimizing the magnitudes of the cross products between $\theta$ and the vectors $\overrightarrow{OP_i}$ through the following weighted least squares cost function:

$$E = \sum_{i=1}^{n} \frac{w_i \| P_i' \times \theta \|^2}{2} \tag{13}$$

Once again, the chosen weights are $w_i = \sigma_{Z_i}^{-2}$. Since the 3D vector $\theta$ is over-parameterized, we reduce it to 2D by fixing one of its dimensions $\theta^{(k)}$ at 1. This dimension cannot be chosen arbitrary, as the optimal vector may have zero entries. Thus, we select the dimension where the range of samples is the highest. Given the resulting 2D parameterization $\theta_m$ and by setting the partial derivatives of (13) equal to zero, we arrive at a solution of the form $\theta_m = A^{-1}b$ with three possible results for $A$ and b depending on the fixed dimension:

$(\theta^{(1)} = 1)$

$$A = \begin{bmatrix} \sum_{i=1}^{n} w_i(X_i^2 + Z_i^2) & -\sum_{i=1}^{n} w_i Y_i Z_i \\ -\sum_{i=1}^{n} w_i Y_i Z_i & \sum_{i=1}^{n} w_i(X_i^2 + Y_i^2) \end{bmatrix}, \ b = \begin{bmatrix} \sum_{i=1}^{n} w_i X_i Y_i \\ \sum_{i=1}^{n} w_i X_i Z_i \end{bmatrix} \tag{14}$$

16

$(\theta^{(2)} = 1)$

$$A = \begin{bmatrix} \sum_{i=1}^{n} w_i(Y_i^2 + Z_i^2) & -\sum_{i=1}^{n} w_i X_i Z_i \\ -\sum_{i=1}^{n} w_i X_i Z_i & \sum_{i=1}^{n} w_i(X_i^2 + Y_i^2) \end{bmatrix}, \ b = \begin{bmatrix} \sum_{i=1}^{n} w_i X_i Y_i \\ \sum_{i=1}^{n} w_i Y_i Z_i \end{bmatrix} \quad (15)$$

$(\theta^{(3)} = 1)$

$$A = \begin{bmatrix} \sum_{i=1}^{n} w_i(Y_i^2 + Z_i^2) & -\sum_{i=1}^{n} w_i X_i Y_i \\ -\sum_{i=1}^{n} w_i X_i Y_i & \sum_{i=1}^{n} w_i(X_i^2 + Z_i^2) \end{bmatrix}, \ b = \begin{bmatrix} \sum_{i=1}^{n} w_i X_i Z_i \\ \sum_{i=1}^{n} w_i Y_i Z_i \end{bmatrix} \quad (16)$$

where, here, $P_i' = \{X_i, Y_i, Z_i\}$ for readability. To obtain $\Sigma_{\theta_m}$, as explained in the last section, the weights need to be rectified with the inverse of the uncertainties of the residuals $r_i = \|P_i' \times \theta\|$ through first order error propagation: $\Sigma_{r_i} = J_{r_i} \Sigma_{P_i} J_{r_i}^\top$. When $r_i$ is exactly zero, $J_{r_i}$ is indeterminate, thus we add a small perturbation to $P_i'$ to avoid such case. Once A is rectified, $\Sigma_\theta$ is found by restructuring $\Sigma_{\theta_m} = A^{-1}$ as a 3×3 matrix, where the entries corresponding to the fixed dimension are 0.

Finally, estimated line endpoints $\{\hat{P}_1, \hat{P}_n\}$ can be sampled through interpolation as follows:

$$\hat{P}_i = O + \lambda_i \theta \quad (17)$$

where $\lambda_i$, denoting the interpolation factor for each endpoint, is obtained by projecting the measured line endpoint onto the estimated line: $\lambda_i = \theta^\top P_i' / \|\theta\|^2$. The endpoint uncertainties are then given by propagating $\Sigma_\theta$ and $\Sigma_O$ through (17).

### 4.5. Matching Points, Lines and Planes

For 2D points and lines, feature correspondences are established between successive frames by matching their descriptors using a $k$-NN search and then select the strongest match per query that satisfies an image geometric distance constraint: the image coordinates of point matches must be within a certain Euclidean distance and 2D line matches must have a similar slope angle and distance to origin (i.e. image top-left corner), according to their line Hessian normal parameterization.

Planes are matched between successive frames using the approach proposed in [2], as follows: First, 1-to-N candidate matches are obtained by enforcing the following constraints:
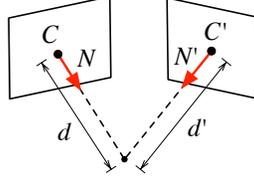
Figure 9: Geometry of two planes and their representation as points: $\{C, C'\}$

- Projection overlap: The projections of two planes, defined as the image segments covered by the inliers of the planes, must have an overlap of at least 50% the number of plane inliers of the smallest plane. This can be checked efficiently by using bitmask operations after checking the geometric constraint.

- Geometric constraint: Given the Hessian plane equations of two planes: $\{N, d\}$ and $\{N', d'\}$, the angle between the plane normals: $\arccos(N \cdot N')$ must be less than $10°$ and the distance: $|d - d'|$ must be less than 10 cm.

To select the best plane match between the plane candidates, we select the plane candidate that yields the minimum plane-to-plane distance, a concept introduced in [2], described as follows: Let $\{N', d'\}$ and $\{N, d\}$ be the equations of two planes then the distance between the two planes is expressed by:

$$\|C - C'\| = \|d'N' - dN\| \tag{18}$$

where $C$ and $C'$ represent points on the planes, as shown in Fig. 9.

### 4.6. Pose Estimation

Pose, defined as the 3D rigid body transformation: $\{R, t \mid R \in SO(3), t \in \mathbb{R}^3\}$, is estimated by jointly minimizing the point and line 3D-to-2D reprojection errors and the 3D plane-to-plane distances. While, pose could be alternatively estimated, as shown in [2] and [20], by minimizing the 3D Mahalanobis distance between point matches, the unidirectional reprojection error tolerates missing depth values and it could be less affected by a depth error that is incorrectly modelled by the uncertainty.

Given a 3D-to-2D point match $\{P, p'\}$, the reprojection error is expressed in the vector form as follows:

$$\widetilde{p} = (p' - \pi(RP + t))^\top \tag{19}$$

whereas the residual of a 3D-to-2D correspondence of line segments is expressed as the point-to-line distance between a 2D line: $l'$ and the projection of the corresponding 3D line endpoints $\{P_1, P_2\}$:

$$\widetilde{l} = l' \begin{bmatrix} \pi(RP_1 + t) & \pi(RP_2 + t) \\ 1 & 1 \end{bmatrix} \tag{20}$$

For two plane matches: $\{N, d\}$ and $\{N', d'\}$, we make use of the plane-to-plane distance, defined in (18), such that, the residual can be derived, in the vector form, as:

$$\widetilde{C} = N'R(N't + d') - dN \tag{21}$$

Given a set of point matches: $S_1$, a set of line matches: $S_2$ and a set of plane matches: $S_3$, we minimize the following joint cost function, by using a Levenberg-Marquart algorithm:

$$E = \sum_{i=1}^{S_1} \widetilde{p_i}^2 w(\widetilde{p_i}) + \sum_{i=1}^{S_2} \widetilde{l_i}^2 w(\widetilde{l_i}) + \alpha \sum_{i=1}^{S_3} \widetilde{C_i}^2 w(\widetilde{C_i}) \tag{22}$$

where $w(r_i)$ is a function that computes a vector of weights for a given residual $r_i$ based on its uncertainty $\Sigma_{r_i}$. Since $\Sigma_{r_i}$ depends on the pose parameters, the weights are recomputed in an iteratively re-weighted least-squares fashion. The residual uncertainties are derived through first order error propagation of (19, 20 and 21) given the uncertainties of the respective extracted primitives (i.e. 3D points, line endpoints and plane equations). Then, $w(r_i)$ returns simply the inverse of the diagonal entries of $\Sigma_{r_i}$. Although, this means that the covariances between dimensions are neglected for points and planes, this allows maintaining the residuals as vectors in the least squares problem, which we have found to improve the convergence, and it does not require inverting the covariance matrices.

Despite the residual weighting, we found necessary to use a fixed scaling factor $\alpha$ to tune the impact of the plane residuals on the pose optimization. In this work, we have used the trade-off $\alpha = 0.02$, based on coarse tuning. Furthermore, an M-estimator with Tukey weights is used to further reweight the point and line residuals in order to down-weight the impact of outliers, whereas plane matching outliers are already addressed by the plane matching method and plane matches are typically too few to rely on statistics.

### 4.6.1. Pose Uncertainty

Assessing the uncertainty of the estimated pose allows: (i) to detect degenerate feature configurations and reject pose estimates under such configurations and (ii) to avoid inaccurate depth fusion due to pose errors. The derivation of the pose uncertainty is described below.

During the pose optimization, the rotation is parameterized locally as a three-dimensional representation of an unit quaternion: $\{q_1, q_2, q_3\}$ such that $q_4 = \sqrt{1 - q_1^2 - q_2^2 - q_3^2}$. Let the pose parameters be: $\xi = \{t_x, t_y, t_z, q_1, q_2, q_3\}$, then its uncertainty can be approximated by back-propagation [41]:

$$\Sigma_\xi = (J_r W J_r^\top)^{-1} \tag{23}$$

where $J_r$ is the stacked Jacobian matrix of the residuals with respect to the pose parameters and $W$ is a diagonal matrix that contains the weights assigned to the residuals. To validate the pose estimate, we simply check if the largest eigenvalue of the matrix block corresponding to the translation vector is larger than a given threshold, if so, the optimized pose is ignored and a decaying velocity model is used instead. To further assess the pose drift between frame 1 and $k + 1$: $\Sigma_\xi^{(k+1|1)}$, the transformation uncertainty can be propagated using the EKF state covariance propagation:

$$\Sigma_\xi^{(k+1|1)} = F \Sigma_\xi^{(k|1)} F^\top + G Q G^\top \tag{24}$$

where $Q$, known as the process noise covariance, is given by (23), and $F$ and $G$ are the Jacobian matrices of the first 6 columns of the following state transition equation, with respect to $\xi^{(k|1)}$ and $\xi^{(k+1|k)}$, respectively:

$$f = \begin{bmatrix} t^{(k+1|1)} \\ q^{(k+1|1)} \end{bmatrix} = \begin{bmatrix} R^{(k+1|k)} t^{(k|1)} + t^{(k+1|k)} \\ q^{(k+1|k)} \otimes q^{(k|1)} \end{bmatrix} \tag{25}$$

where $\otimes$ denotes the quaternion product. As mentioned in Section 3.4, this framework is used to validate the propagation of depth measurements from the sliding window of frames. Specifically, the uncertainties of the transformations between the current frame and the frames where the depth measurements were taken are continuously updated using (24) and validated based on the largest eigenvalue criterion (described above).

## 5. Experiments and Results

In order to evaluate the proposed dataset, we tested our method on various sequences from two public RGB-D datasets: the TUM benchmark [36] and the synthetic ICL-NUIM [35] benchmark. The results of this evaluation are reported and discussed in the next section. Additionally, we have captured four RGB-D sequences in structured environments using the setup shown in Fig. 11 and evaluated the trajectory estimated by our method, in Section 5.2. Throughout the experiments, we fixed the maximum length of the depth fusion sliding window at 10 frames. Timing results are reported and discussed in Section 5.3.

### 5.1. Public datasets

For the sake of diversity, the following sequences were selected from the TUM dataset for evaluation: *fr1/desk* and *fr1/360* captured from a textured office; *fr3/struct_no_text_far* and *fr3/cabinet* collected from low textured and structured scenes; *fr3/walking_static* captured from a dynamic environment with people walking; and *fr2/360_hemisphere* captured in a warehouse. The ICL-NUIM dataset contains two versions of sequences rendered from realistic models of an office and a living room, one without any noise and another with simulated RGB and depth noise, as show in Fig. 5. Two sequences were selected respectively: *kt0 (lr)* from the living room and *kt0 (or)* from the office. While the RGB noise does not seem significant, the depth noise introduced around the object boundaries, seen in Fig. 5, is significantly worse than the observed depth noise of real structured-light cameras (see Fig. 14). Furthermore, the depth map boundaries are corrupted with dense noise, thus we removed all depth values within a margin of 5 pixels. Both the relative pose error (RPE) per second and the absolute trajectory error (ATE) are reported, as RMSEs, for these sequences and the estimated trajectories are compared against the ground-truth in Fig. 10.

Table 1 and 2 shows how the performance is improved by introducing new feature-types. The performance gain of using the three geometric primitives is consistent and significant, especially in low textured environments and in the *fr1/360*, where several RGB images are blurred due to sudden rotations, causing few detected feature points.

21

| Features | Points | Points & Lines | Points & Planes | All |
|---|---|---|---|---|
| fr1/desk | 34 mm 2.4 deg | 30 mm 2.2 deg | 28 mm 1.9 deg | 23 mm 1.7 deg |
| fr1/360 | 88 mm 4.4 deg | 69 mm 3.1 deg | 76 mm 3.5 deg | 64 mm 2.7 deg |
| fr3/struct_no_text_far | Fail | 32 mm 0.9 deg | 28 mm 0.9 deg | 19 mm 0.7 deg |
| fr3/cabinet | 112 mm 4.5 deg | 70 mm 2.8 deg | 40 mm 1.8 deg | 39 mm 1.8 deg |
| fr3/walking_static | 87 mm 1.1 deg | 69 mm 1.0 deg | 86 mm 1.1 deg | 68 mm 0.7 deg |
| fr2/360_hemisphere | 78 mm 1.4 deg | 72 mm 1.1 deg | 79 mm 1.7 deg | 69 mm 1.1 deg |
| kt0 (lr) | 8 mm 0.6 deg | 9 mm 0.6 deg | 8 mm 0.6 deg | 7 mm 0.5 deg |
| kt0 (lr) w/ noise | 8 mm 0.6 deg | 8 mm 0.7 deg | 7 mm 0.5 deg | 6 mm 0.5 deg |
| kt0 (or) | 9 mm 0.5 deg | 7 mm 0.5 deg | 7 mm 0.5 deg | 7 mm 0.5 deg |
| kt0 (or) w/ noise | 6 mm 0.5 deg | 5 mm 0.5 deg | 7 mm 0.5 deg | 6 mm 0.5 deg |

Table 1: RPE on TUM and ICL_NUIM datasets for different combinations of geometric primitives.

| Features | Points | Points & Lines | Points & Planes | All |
|---|---|---|---|---|
| fr1/desk | 64 mm | 53 mm | 50 mm | 40 mm |
| fr1/360 | 116 mm | 109 mm | 91 mm | 91 mm |
| fr3/struct_no_text_far | Fail | 80 mm | 63 mm | 54 mm |
| fr3/cabinet | 437 mm | 241 mm | 195 mm | 200 mm |
| fr3/walking_static | 200 mm | 181 mm | 199 mm | 179 mm |
| fr2/360_hemisphere | 237 mm | 238 mm | 350 mm | 203 mm |
| kt0 (lr) | 496 mm | 446 mm | 76 mm | 99 mm |
| kt0 (lr) w/ noise | 428 mm | 281 mm | 83 mm | 59 mm |
| kt0 (or) | 197 mm | 31 mm | 99 mm | 27 mm |
| kt0 (or) w/ noise | 237 mm | 199 mm | 134 mm | 167 mm |

Table 2: ATE on TUM and ICL_NUIM datasets for different combinations of geometric primitives.

Table 3 compares the performance between using the different depth models (i.e. stages) described in Section 3. Overall, fusing the depth maps using the Optimal-GM framework decreases significantly the odometry error. Interestingly enough, the performance in the sequence captured in the dynamic environment is significantly improved by using the full depth filter framework, which indicates that modelling temporally the depth uncertainty helps reducing the impact of moving objects. In the ICL-NUIM captures, the introduction of simulated noise increases overall the ATE, however, the virtual camera in *kt0*

| Error | RPE | | | ATE | | |
|---|---|---|---|---|---|---|
| Depth model | Sensor model | C-GM | O-GM | Sensor model | C-GM | O-GM |
| fr1/desk | 32 mm 2.3 deg | 33 mm 2.4deg | 23 mm 1.7 deg | 60 mm | 65 mm | 40 mm |
| fr1/360 | 67 mm 2.8 deg | 66 mm 2.9 deg | 64 mm 2.7 deg | 127 mm | 123 mm | 91 mm |
| fr3/struct_no_text_far | 27 mm 0.9 deg | 29 mm 0.9 deg | 19 mm 0.7 deg | 82 mm | 95 mm | 54 mm |
| fr3/cabinet | 56 mm 2.3 deg | 62 mm 2.5 deg | 39 mm 1.8 deg | 239 mm | 275 mm | 200 mm |
| fr3/walking_static | 149 mm 2.0 deg | 144 mm 1.9 deg | 68 mm 0.7 deg | 417 mm | 407 mm | 179 mm |
| fr2/360_hemisphere | 77 mm 1.1 deg | 73 mm 1.1 deg | 69 mm 1.1 deg | 198 mm | 193 mm | 203 mm |
| kt0 (lr) | 7 mm 0.5 deg | 7 mm 0.5 deg | 7 mm 0.5 deg | 198 mm | 97 mm | 99 mm |
| kt0 (lr) w/ noise | 6 mm 0.6 deg | 6 mm 0.5 deg | 6 mm 0.5 deg | 303 mm | 113 mm | 59 mm |
| kt0 (or) | 7 mm 0.5 deg | 7 mm 0.5 deg | 7 mm 0.5 deg | 25 mm | 36 mm | 27 mm |
| kt0 (or) w/ noise | 5 mm 0.5 deg | 6 mm 0.5 deg | 6 mm 0.5 deg | 359 mm | 220 mm | 167 mm |

Table 3: RMSE on TUM and ICL_NUIM datasets for different depth uncertainty models. As described in Section 3, the sensor model corresponds to the first stage of the proposed filter method, the C-GM uses the two first stages and the O-GM uses the full depth fusion method.

(lr) faces a texture-less wall in the middle of the sequence, which causes the pose estimation to fail for a few frames, thus the ATE is affected by the employed velocity model. It is worth noting that although C-GM by itself does not seem advantageous in most sequences, when large amount of noise is present in the ICL-NUIM, we observe the contrary in terms of ATE.

Our method is compared to state-of-the-art visual odometry methods in Tables 4 and 5 for each respective dataset. For the sake of fairness, our comparison does not include full SLAM systems that perform map optimization or loop closure detection. In terms of RPE, our method achieves state-of-the-art results in the TUM dataset.

| | Ours | | State-of-the-art (VO) | |
|---|---|---|---|---|
| Error | RPE | ATE | RPE | ATE |
| fr1/desk | 23 mm | 40 mm | 25 mm [18] | 32 mm [18] |
| fr1/360 | 64 mm | 91 mm | 73 mm [2] | - |
| fr3/struct_no_text | 19 mm | 54 mm | 43 mm [3] | 19 mm [42] |
| fr3/cabinet | 39 mm | 200 mm | 80 mm [2] | 268 mm [42] |
| fr3/walking_static | 68 mm | 179 mm | 111 mm [37] | - |
| fr2/360_hemisphere | 69 mm | 203 mm | 66 mm [24] | - |

Table 4: Comparison of visual odometry methods on TUM dataset.

|                    | Ours    | DVO    | FOVIS    | [18]    |
|--------------------|---------|--------|----------|---------|
| kt0 (lr)           | 99 mm   | 114 mm | 1931 mm  | 10 mm   |
| kt0 (lr) w/ noise  | 59 mm   | 291 mm | 2051 mm  | 6 mm    |
| kt0 (or)           | 27 mm   | 398 mm | 3396 mm  | 4 mm    |
| kt0 (or) w/ noise  | 167 mm  | 335 mm | 3296 mm  | 15 mm   |

Table 5: Comparison of absolute trajectory errors obtained by several visual odometry methods on ICL-NUIM dataset, according to the results published in [35] and [18]. For [18], we report the results for the best overall visibility ratio threshold used in the keyframe selection.



Figure 10: Trajectories estimated by the proposed method versus the groundtruth

Notably, we outperform the point and line odometry methods [3, 20], and the recently proposed [37], which addresses explicitly dynamic environments. In the ICL-NUIM dataset, our method is compared against: the DVO [16], which minimizes densely the photometric error; the feature point-based FOVIS [43]; and the method proposed in [18], which is essentially an improved version of the DVO that minimizes also the geometric error by adopting an inverse depth parameterization and performs keyframe-to-frame alignment. Although, our method outperforms DVO and FOVIS, it performs significantly worse than the remarkable performance published in [18]. Nevertheless, the frame-to-frame version of this method, shown in [18], compares well with our results, suggesting that this observed discrepancy is partially due the use of a keyframe-to-frame strategy, which is known to be a good way to avoid the accumulation of pose errors.

## 5.2. Author-collected dataset

Here, we evaluate the visual odometry method on four closed-loop trajectories, recorded, whilst walking, by the hand-held RGB-D setup shown in Fig. 11. The employed depth sensor is also based on structured-light and follows, along with Kinect 1, the same design as Primesense cameras [44] with an equal projector-camera baseline of 75 mm, thus we used the same depth sensor noise model. As shown in Fig. 14, all sequences were collected in structured environments, where low textured surfaces are predominant. To avoid degenerate scene configurations due to the lack of textures and depth information, a wide-angle lens was mounted on the color camera to expand the FOV and the depth fusion was used to recover depth values outside the FOV of the depth camera.

The drawback of expanding the FOV is that the pixel resolution is reduced and consequently mapping the depth measurements to the color image downsamples significantly the depth map, since the FOV of the RGB camera is more than the double of the depth camera FOV. As this results in the projection multiple depth measurements to a pixel, analogous to temporal depth fusion, we used, here, the O-GM method in place of the C-GM, at the second stage, to obtain simultaneously the aligned depth image and its uncertainty.

Table 6 reports the final trajectory error for different feature-type combinations, while Fig. 12 shows the respective estimated trajectories. The translational and angular error were measured in these closed-loop trajectories by

Figure 11: Mobile RGB-D capture setup used in this work. An Occipital Structure sensor is used to collect depth while RGB is collected by a wide-angle lens (Moment wide lens) mounted on an iPhone camera. The effective angle-of-views are respectively $85° \times 70°$ and $58° \times 45°$ for the color and the infrared camera and the baseline between them is around 38 mm. Both cameras operate at 30 fps with VGA resolution.

| Sequence (distance) | Points | Points & Lines | Points & Planes | All |
|---|---|---|---|---|
| lab (51 m) | 2.3 m | 1.1 m | 2.4 m | 1.2 m |
| | 20 deg | 12 deg | 19 deg | 12 deg |
| parking garage 1 (56 m) | 1.6 m | 1.9 m | 4.2 m | 1.3 m |
| | 17 deg | 23 deg | 14 deg | 18 deg |
| parking garage 2 (53 m) | 7.9 m | 1.7 m | 3.0 m | 0.8 m |
| | 51 deg | 10 deg | 18 deg | 7 deg |
| corridor (110 m) | 23 m | 5.7 m | 13.8 m | 3.6 m |
| | 47 deg | 26 deg | 37 deg | 25 deg |

Table 6: Final trajectory errors for the RGB-D sequences collected in this work.

using a marker [45] and expressing the angular error through the angle-axis representation. The combination of points, lines and planes shows consistently better results than the other versions, whereas the point odometry yields poor results and loses tracking for several frames. Although the final error in the *parking garage 1* indicates that using just points is better than combining points and lines, the estimated trajectory, shown in Fig. 12, is coarser when using just points. Additionally, we observed that both versions fail tracking in this sequence, thus switch temporally the pose estimation to the velocity model, contrary to the full combination.

### 5.3. Processing Time

All sequences were processed offline with a single thread on an Intel Core i5-6500 CPU 3.20 GHz. The method was implemented on MATLAB with C++ mex functions for certain modules as indicated in Table 7. There are many opportunities for optimization: As shown in Fig. 8, point, line and plane processing (i.e. detection, extraction and matching) can be parallelized in three
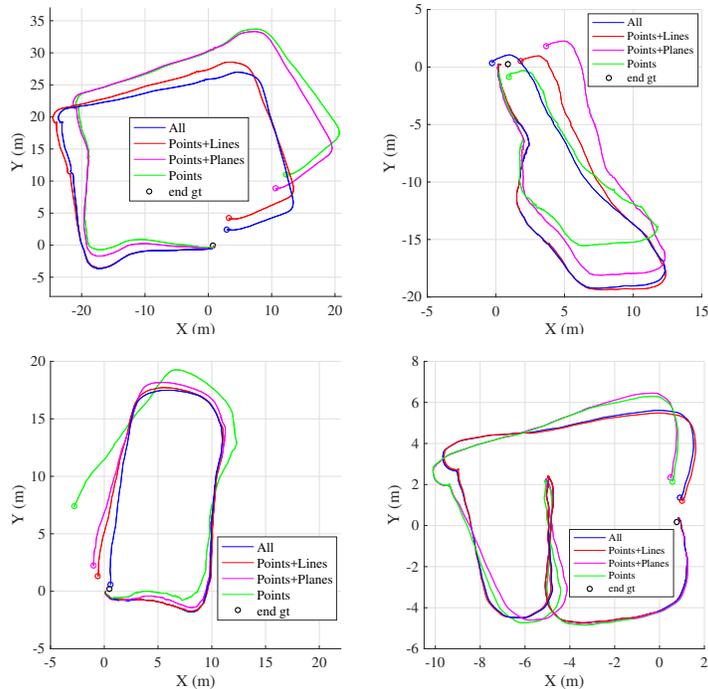
Figure 12: Top view of the trajectories estimated by the proposed method for different feature combinations versus the final position groundtruth (end gt). Clockwise from top left: corridor, parking garage 1, lab and parking garage 2.

threads. However the plane extraction must be called once more after the depth fusion. As shown in [29], the plane extraction can be speeded-up significantly by avoiding the per-pixel refinement and using a coarser graph. In terms of feature points, faster alternatives to SURF are well known [46], whereas for detection of line segments, unfortunately to our knowledge, there is a lack of good alternatives, thus one can either do line detection at half resolution (QVGA) or adopt a line tracking approach as in [4].

As shown in Fig. 13, increasing the size of the sliding window for the temporal fusion beyond 10 frames can improve even further the RPE performance, however the the cost of depth fusion grows linearly with the number of frames.

## 6. Conclusion and Future Work

Combining points with lines and planes proves to improve the robustness of visual odometry. Our results show no redundancy between the different
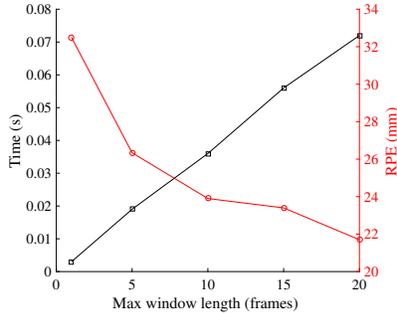
27

Figure 13: Maximum number of frames used for temporal depth fusion vs. O-GM runtime (black) and RPE (red) on the *fr1/desk* sequence.

|  | Time |
|---|---|
| O-GM$_{10}$ | 35 ms |
| 2D point processing | 31 ms |
| 2D line processing | 40 ms |
| 3D line RANSAC sampling | 2.3 ms |
| 3D line Fitting (per line) † | 0.1 ms |
| Plane extraction | 29 ms |
| Plane fitting (per plane) † | 1.5 ms |
| Plane matching † | 3 ms |
| Pose estimation † | 33 ms |

Table 7: Timing average results on *fr1/desk*. The three stages of the depth filter were timed both for a sliding window of 5 frames, as O-GM$_5$, and 10 frames, as O-GM$_{10}$, used in our experiments. Both the 2D point and line processing include feature detection, extraction and matching. The processes marked with a † are implemented on MATLAB, while the rest is implemented on C++ and integrated through mex functions.

primitives. The proposed method achieves state-of-the-art results, in terms of RPE, between frame-to-frame odometry methods. However, in the ICL-NUIM dataset, the ATE yield by our visual odometry is still inferior to some model-to-frame [1] and keyframe-to-frame [18] based approaches, thus, extending our method to a SLAM version, such as in [32], is promising research direction in order to reduce the pose drift.

Furthermore, the visual odometry performance is improved significantly by using the proposed depth fusion framework. While this framework shows its capability to denoise the raw depth maps, errors may be still introduced and propagated due to flying pixels, occlusions and pose errors, therefore modelling the depth uncertainty is necessary to capture these errors. The full system was

additionally evaluated on RGB-D video sequences, captured with a wide-angle RGB camera, where the depth fusion framework plays also the role of recovering old depth measurements that are no longer inside the current FOV of the depth camera. These past measurements can only be maintained for a small number of frames, specified by the depth fusion framework. A more flexible strategy could extend this duration by switching old pixels from depth fusion to a simple hole-filling mode, where points would only be maintained and used if they had an unique pixel projection.
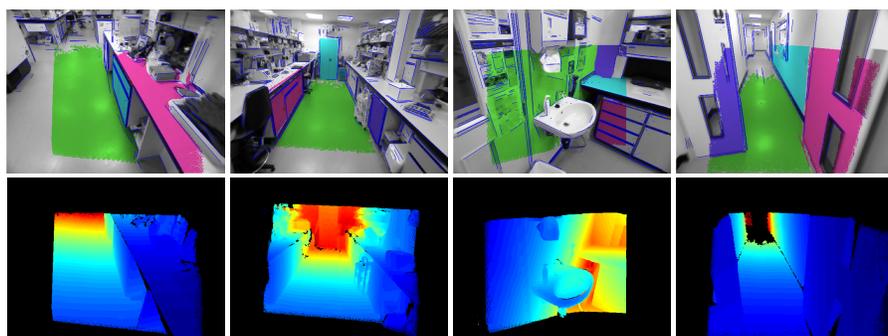
## Acknowledgement

## References

[1] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohi, J. Shotton, S. Hodges, A. Fitzgibbon, Kinectfusion: Real-time dense surface mapping and tracking, in: International symposium on Mixed and augmented reality (ISMAR), IEEE, 2011, pp. 127–136.

[2] P. Proença, Y. Gao, Probabilistic combination of noisy points and planes for rgb-d odometry, in: Towards Autonomous Robotic Systems, Springer, 2017, pp. 340–350.

[3] S. Yang, S. Scherer, Direct monocular odometry using points and lines, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2017.

[4] R. Gomez-Ojeda, J. Briales, J. González-Jiménez, Pl-svo: Semi-direct monocular visual odometry by combining points and line segments, in: Proc. International Conference on Intelligent Robots and Systems (IROS), IEEE, 2016.

[5] Y. Lu, D. Song, Robustness to lighting variations: An rgb-d indoor visual odometry using line segments, in: Proc. International Conference on Intelligent Robots and Systems (IROS), IEEE, 2015, pp. 688–694.

[6] K. Khoshelham, S. O. Elberink, Accuracy and resolution of kinect depth data for indoor mapping applications, Sensors 12 (2) (2012) 1437–1454.

[7] H. Sarbolandi, D. Lefloch, A. Kolb, Kinect range sensing: Structured-light versus time-of-flight kinect, Computer Vision and Image Understanding 139 (2015) 1 – 20.

[8] J. Han, L. Shao, D. Xu, J. Shotton, Enhanced computer vision with microsoft kinect sensor: A review, IEEE transactions on cybernetics 43 (5) (2013) 1318–1334.

[9] J. Smisek, M. Jancosek, T. Pajdla, 3d with kinect, in: Consumer depth cameras for computer vision, Springer, 2013, pp. 3–25.

[10] O. Wasenmüller, D. Stricker, Comparison of kinect v1 and v2 depth images in terms of accuracy and precision, in: Proc. Asian Conference on Computer Vision Workshop (ACCV workshop), Springer, 2016.

[11] C. V. Nguyen, S. Izadi, D. Lovell, Modeling kinect sensor noise for improved 3d reconstruction and tracking, in: Proc. International Conference on 3D Imaging, Modeling, Processing, Visualization Transmission, 2012, pp. 524–530.

[12] C. Tomasi, R. Manduchi, Bilateral filtering for gray and color images, in: International Conference on Computer Vision, IEEE, 1998, pp. 839–846.

[13] M. Keller, D. Lefloch, M. Lambers, S. Izadi, T. Weyrich, A. Kolb, Real-time 3d reconstruction in dynamic scenes using point-based fusion, in: International Conference on 3D Vision – 3DV, 2013.

[14] O. Wasenmller, M. Meyer, D. Stricker, Augmented reality 3d discrepancy check in industrial applications, in: International Symposium on Mixed and Augmented Reality (ISMAR), IEEE, 2016, pp. 125–134.

[15] M. Hsiao, E. Westman, G. Zhang, M. Kaess, Keyframe-based dense planar slam, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2017.

[16] C. Kerl, J. Sturm, D. Cremers, Robust odometry estimation for rgb-d cameras, in: Proc. International Conference on Robotics and Automation (ICRA), 2013, pp. 3748–3754.
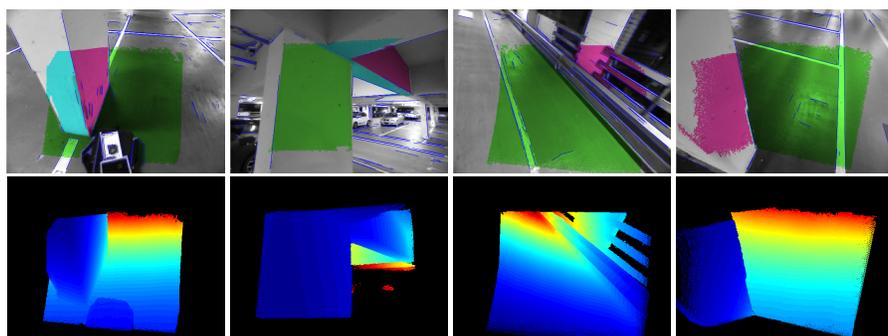
[17] C. Kerl, J. Sturm, D. Cremers, Dense visual slam for rgb-d cameras, in: Proc. International Conference on Intelligent Robots and Systems (IROS), IEEE, 2013, pp. 2100–2106.

[18] D. Gutierrez-Gomez, W. Mayol-Cuevas, J. J. Guerrero, Dense rgb-d visual odometry using inverse depth, Robotics and Autonomous Systems 75 (2016) 571–583.

[19] O. Wasenmüller, M. D. Ansari, D. Stricker, Dna-slam: Dense noise aware slam for tof rgb-d cameras, in: Asian Conference on Computer Vision Workshop (ACCV workshop), Springer, 2016.

[20] Y. Lu, D. Song, Robust rgb-d odometry using point and line features, in: International Conference on Computer Vision (ICCV), IEEE, 2015.

[21] I. Dryanovski, R. G. Valenti, J. Xiao, Fast visual odometry and mapping from rgb-d data, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2013, pp. 2305–2310.

[22] Z. Fang, S. Scherer, Experimental study of odometry estimation methods using rgb-d cameras, in: Proc. International Conference on Intelligent Robots and Systems (IROS), IEEE, 2014, pp. 680–687.

[23] A. Pumarola, A. Vakhitov, A. Agudo, A. Sanfeliu, F. Moreno-Noguer, PL-SLAM: Real-Time Monocular Visual SLAM with Points and Lines, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2017.

[24] P. Proença, Y. Gao, Splode: Semi-probabilistic point and line odometry with depth estimation from rgb-d camera motion, in: Proc. International Conference on Intelligent Robots and Systems (IROS), IEEE, 2017, pp. 1594–1601.

[25] A. J. Trevor, J. G. Rogers, H. I. Christensen, Planar surface slam with 3d and 2d sensors, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2012, pp. 3041–3048.

[26] Y. Taguchi, Y.-D. Jian, S. Ramalingam, C. Feng, Point-plane slam for hand-held 3d sensors, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2013, pp. 5182–5189.

[27] E. Ataer-Cansizoglu, Y. Taguchi, S. Ramalingam, Pinpoint slam: A hybrid of 2d and 3d simultaneous localization and mapping for rgb-d sensors, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2016, pp. 1300–1307.

[28] E. Ataer-Cansizoglu, Y. Taguchi, S. Ramalingam, T. Garaas, Tracking an rgb-d camera using points and planes, in: Proceedings of the IEEE International Conference on Computer Vision Workshops, 2013, pp. 51–58.

[29] C. Feng, Y. Taguchi, V. R. Kamat, Fast plane extraction in organized point clouds using agglomerative hierarchical clustering, in: IEEE International Conference on Robotics and Automation (ICRA), 2014, pp. 6218–6225.

[30] D. Holz, S. Holzer, R. B. Rusu, S. Behnke, Real-time plane segmentation using rgb-d cameras, in: Robot Soccer World Cup, Springer, 2011, pp. 306–317.

[31] R. Li, Q. Liu, J. Gui, D. Gu, H. Hu, A novel rgb-d slam algorithm based on points and plane-patches, in: International Conference on Automation Science, IEEE, 2016, pp. 1348–1353.

[32] L. Ma, C. Kerl, J. Stückler, D. Cremers, Cpa-slam: Consistent plane-model alignment for direct rgb-d slam, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2016, pp. 1285–1291.

[33] K. Pathak, N. Vaskevicius, A. Birk, Uncertainty analysis for optimum plane extraction from noisy 3d range-sensor point-clouds, Intelligent Service Robotics 3 (1) (2010) 37–48.

[34] B. Efron, D. V. Hinkley, Assessing the accuracy of the maximum likelihood estimator: Observed versus expected fisher information, Biometrika (1978) 457–482.

[35] A. Handa, T. Whelan, J. McDonald, A. Davison, A benchmark for RGB-D visual odometry, 3D reconstruction and SLAM, in: Proc. International Conference on Robotics and Automation (ICRA), IEEE, 2014.

[36] J. Sturm, N. Engelhard, F. Endres, W. Burgard, D. Cremers, A benchmark for the evaluation of rgb-d slam systems, in: Proc. International Conference on Intelligent Robots and Systems (IROS), IEEE, 2012.
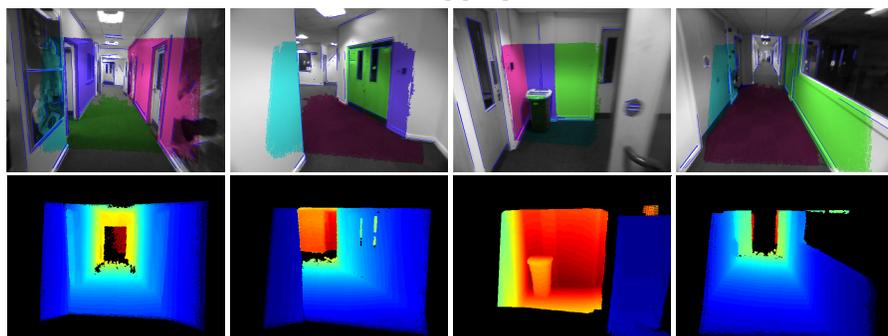
[37] M. Jaimez, C. Kerl, J. Gonzlez-Jimnez, D. Cremers, Fast odometry and scene flow from rgb-d cameras based on geometric clustering, in: Proc. International Conference on Robotics and Automation (ICRA), 2017.

[38] v. G. R. Grompone, J. Jakubowicz, J.-M. Morel, G. Randall, Lsd: a fast line segment detector with a false detection control., IEEE transactions on pattern analysis and machine intelligence 32 (4) (2010) 722–732.

[39] L. Zhang, R. Koch, An efficient and robust line segment matching approach based on lbd descriptor and pairwise geometric consistency, Journal of Visual Communication and Image Representation 24 (7) (2013) 794–805.

[40] J. W. Weingarten, G. Gruener, R. Siegwart, Probabilistic plane fitting in 3d and an application to robotic mapping, in: Proc. International Conference on Robotics and Automation (ICRA), Vol. 1, IEEE, 2004, pp. 927–932.

[41] R. Hartley, A. Zisserman, Multiple view geometry in computer vision, Cambridge university press, 2003.

[42] X. Wang, D. Wei, M. Zhou, R. Li, H. Zha, Edge enhanced direct visual odometry, in: Proc. British Machine Vision Conference (BMVC), 2016.

[43] A. S. Huang, A. Bachrach, P. Henry, M. Krainin, D. Maturana, D. Fox, N. Roy, Visual odometry and mapping for autonomous flight using an rgb-d camera, in: International Symposium on Robotics Research (ISRR), 2011, pp. 1–16.

[44] P. Zanuttigh, G. Marin, C. Dal Mutto, F. Dominio, L. Minto, G. M. Cortelazzo, Time-of-Flight and Structured Light Depth Cameras, Springer, 2016.

[45] W. Daniel, S. Dieter, Artoolkitplus for pose tracking on mobile devices, in: Proc. of Computer Vision Winter Workshop, 2007.

[46] M. Calonder, V. Lepetit, C. Strecha, P. Fua, Brief: Binary robust independent elementary features, Computer Vision–ECCV 2010 (2010) 778–792.

Lab



Parking garage



Corridor

Figure 14: Frames from the dataset sequences collected and evaluated in this work. Intensity images with overlaid detected features are shown on the top, while the respective fused and aligned depth maps are shown on the bottom.