

Integration of Absolute Orientation Measurements in the KinectFusion Reconstruction pipeline

Silvio Giancola, Jens Schneider, Peter Wonka and Bernard S. Ghanem

King Abdullah University of Science and Technology (KAUST), Saudi Arabia

{silvio.giancola, jens.schneider, peter.wonka, bernard.ghanem}@kaust.edu.sa

Abstract

In this paper, we show how absolute orientation measurements provided by low-cost but high-fidelity IMU sensors can be integrated into the KinectFusion pipeline. We show that integration improves both runtime, robustness and quality of the 3D reconstruction. In particular, we use this orientation data to seed and regularize the ICP registration technique. We also present a technique to filter the pairs of 3D matched points based on the distribution of their distances. This filter is implemented efficiently on the GPU. Estimating the distribution of the distances helps control the number of iterations necessary for the convergence of the ICP algorithm. Finally, we show experimental results that highlight improvements in robustness, a speed-up of almost 12%, and a gain in tracking quality of 53% for the ATE metric on the Freiburg benchmark.

1. Introduction

Automated 3D reconstruction of geometry from images is a highly versatile field of research with applications in archeology [11], topography [12], urban planning [27], robotics [19] and entertainment [31]. The ubiquity of this core vision task stems primarily from decades of advances in reconstruction algorithms [30, 39, 7, 16] and acquisition devices [1, 10, 8, 17, 37, 5]. With the latest advent of low-cost consumer-class RGB-D cameras, 3D scene understanding and reconstruction has become a pervasive mass market technology.

Simultaneous Localization and Mapping (SLAM) techniques are based on a registration step, i.e., they seek to properly align pairs of images or point clouds. This registration step consists of finding a 6-DoF rigid body transformation $\mathbf{M} \in \mathbb{SE}(3)$ composed of a rotation matrix $\mathbf{R} \in \mathbb{SO}(3)$ and a translation vector $\mathbf{t} \in \mathbb{R}^3$ such that overlapping parts of a scene can be superimposed, thereby minimizing the distances between pairs of matched points. It is clear that the process of finding the best-possible rigid transformation

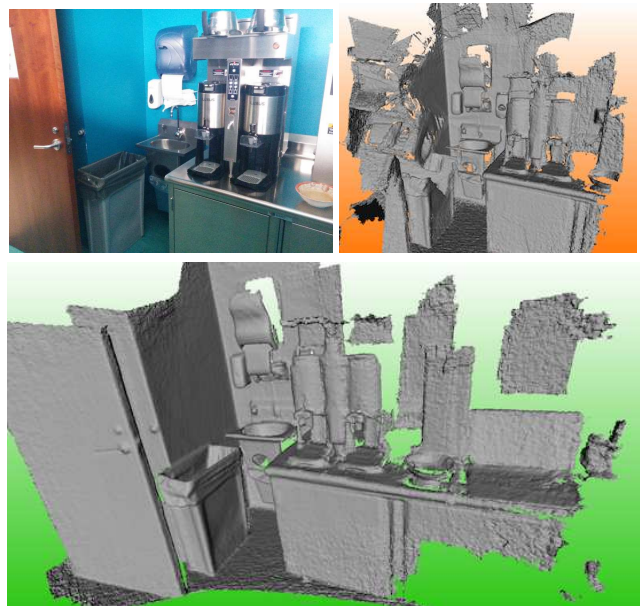


Figure 1. Reconstruction of a Kitchen environment using an absolute orientation prior to seed the ICP algorithm. **Top Left:** Photo of the Kitchen environment. **Top Right:** Original KinectFusion reconstruction. **Bottom:** IMU-seeded KinectFusion reconstruction.

$\mathbf{M} = (\mathbf{R}, \mathbf{t})$ is affected by noise, occlusion, matching quality, and outliers. Therefore, it remains intrinsically ambiguous and complex to solve (\mathbf{R}, \mathbf{t}) in a robust and computationally efficient fashion. The main reason is that solving for both \mathbf{R} and \mathbf{t} simultaneously is a non-convex problem. We believe that non-convexity is principally introduced by the orientation component. The problem can be linearized, however, if a rough estimate of the orientation is available.

Point cloud registration, be it frame-by-frame or frame-to-model, is the foundation for 3D reconstruction, since separate acquisitions have to be aligned. The original KinectFusion algorithm [28] makes use of a coarse-to-fine ICP method [2] with a fixed number of iterations, using a linear approximation [23] that minimizes for the point-to-plane metric [35].

In this work, we present a modified KinectFusion framework that integrates an absolute orientation measurement provided by an off-the-shelf IMU. The orientation prior is leveraged in the ICP method and the original problem is solved with a regularization term for the orientation component. Also, correspondences between closest points are filtered based on an efficient estimate of the median distance between pairs of points implemented on a GPU. As such, we observe improvements in quality, runtime, and robustness in the 3D reconstruction process. We characterize the uncertainty model for the orientation measurement according to the ISO JCGM 100:2008 [6], and we apply such model to the Freiburg dataset [38] in order to evaluate the quality of our reconstruction pipeline.

2. Related work

Visual-Inertial SLAM is a hot topic with direct applications to robotics, for which the setup usually includes cameras and IMUs. Here, we review works focusing on *Navigation* and *Mapping*, dual topics applied to robotics.

Navigation focuses on estimating the system’s ego-motion, i.e., its own trajectory in space and time. Many important navigation methods use *extended Kalman filter* (EKF) techniques [32, 42, 15, 21]. They define and update a system state based on visual and inertial measurements using advanced and complex Kalman filtering techniques [26, 25, 22]. Paul et al. [32] recently conducted a comparative study between monocular and stereo visual-inertial EKF-based techniques, but do not mention any RGB-D devices. Other methods such as [20, 9, 34] focus on solving a *global optimization problem* minimizing for a combined cost function. In those techniques, the translation is estimated by integrating the IMU acceleration measurements twice. However, since integration may amplify noise, such double integration generally produces significant drifts. Leutenegger et al. [20] modeled the noise in a probabilistic way, but residual drift inevitably affects the method’s stability negatively. In contrast, we show that orientation measurements can be estimated robustly and consistently using a 9-DoF IMU. In Visual-Inertial Navigation techniques, the IMU data can be *loosely-coupled*, i.e., an independent measurement in the problem [18, 41] or *tightly-coupled*, optimizing over all sensor measurements in the optimization [20] and EKF pipelines [26, 42]. Finally, it is worth noting that navigation techniques mostly make use of monocular or stereo cameras, but do not rely on direct depth measurements, mainly because the quality of the 3D reconstruction is usually out-of-scope for navigation tasks.

Mapping. While Navigation tries to solve for ego-motion, Mapping focuses attention on the 3D reconstruction of the surrounding environment. Zhang et al. [43] recently presented a pipeline that estimates ego-motion while building an accurate representation of the surrounding en-

vironment. Similar to Navigation, it is common practice to rely on EKF [3] or global optimization [24] techniques. Ma et al. [24] presented a solution for large-scale Visual-Inertial reconstruction using a volumetric representation similar to KinectFusion. Concha et al. [4] created a real-time, fully dense reconstruction based on an RGB camera and an IMU.

Visual Inertial KinectFusion. Nießner et al. [29] followed similar considerations by using an IMU from a smartphone to initialize both (\mathbf{R}, \mathbf{t}) in the ICP step of KinectFusion. While achieving some improvement in runtime, reliable position information cannot be extracted from an IMU alone, since the accelerometer’s noise leads to significant drift. Smartphones therefore use a fusion of multiple triangulation techniques based on measurements including GSM, WiFi, and GPS (outdoors). Nevertheless, the position information is of very low fidelity, especially indoors. We argue that since ICP is generally sensitive to its initialization, providing it with such noisy measurements can only lead to limited improvement (or even degradation) in many acquisition scenarios. In contrast to requiring an expensive, full sensor array as found in smartphones, our method only relies on a cheap and modular IMU that produces high-fidelity absolute orientation measurements in real-time.

Datasets. Pfrommer et al. [33] recently released a dataset for visual-inertial benchmarking. However, this dataset strictly focuses on odometry along long distances, but not on small range reconstructions as KinectFusion-based algorithms do. Sturm et al. [38] released multiple versions of their well-known Freiburg benchmark, albeit without orientation measurements. We will present a way to generate such measurements using the ground truth orientation and an empirical uncertainty model of our sensor.

Contributions. While the existing literature focuses on solving the problem for both position and orientation from the IMU raw measurements, we argue that orientation alone is sufficient to improve overall performance. Our contributions are: **(i)** We overcome the non-convexity of the joint problem induced by the unknown rotation \mathbf{R} by seeding the ICP algorithm with an orientation estimate from the IMU. **(ii)** We present a regularized point-to-plane metric in the coarse-to-fine ICP alignment. **(iii)** We use a novel filter that estimates the median distance between closest points efficiently on the GPU in order to reject outliers and control ICP convergence. **(iv)** We provide an uncertainty model of the IMU measurements and apply our model to the Freiburg dataset.

3. Proposed method

The original KinectFusion algorithm is based on a coarse-to-fine ICP between a newly acquired frame and the model reconstructed so far. The current frame is aligned by solving the non-linear least squares minimization problem shown in Eq. (1) using the point-to-plane metric to find the

optimal transformation matrix $\mathbf{M} \in \mathbb{SE}(3)$, composed of an orientation $\mathbf{R} \in \mathbb{SO}(3)$ and a translation $\mathbf{t} \in \mathbb{R}^3$.

$$\mathbf{M}_{opt} = \arg \min_{\mathbf{M}} \sum_i \left((\mathbf{M}\mathbf{p}_i - \mathbf{q}_i)^\top \mathbf{n}_i \right)^2, \quad (1)$$

where \mathbf{p}_i and \mathbf{q}_i are pairs of matched points belonging to the current frame and the model respectively. In our method, we use the orientation measurement from the IMU to pre-orient the point \mathbf{p}_i beforehand. Inspired by Low et al. [23], who proposed a linear solution to this problem by assuming small angles ($\alpha \simeq 0, \beta \simeq 0, \gamma \simeq 0$), transformation \mathbf{M} is approximated by $\widehat{\mathbf{M}}$ according to Eq. (2), intended to solve for small angle after applying the IMU seed.

$$\mathbf{M} \simeq \begin{pmatrix} 1 & -\gamma & \beta & t_x \\ \gamma & 1 & -\alpha & t_y \\ -\beta & \alpha & 1 & t_z \\ 0 & 0 & 0 & 1 \end{pmatrix} := \widehat{\mathbf{M}}. \quad (2)$$

In this small orientation setting, the optimization of Eq. (1) can be reformulated into the linear least square shown in Eq. (3), where $\mathbf{x} = (\alpha, \beta, \gamma, t_x, t_y, t_z)^\top$ and \mathbf{A} is a $n \times 6$ matrix. Please refer to the original implementation by Low et al. [23] for further details on \mathbf{A} and \mathbf{b} .

$$\mathbf{x}_{opt} = \min_{\mathbf{x}} \frac{1}{2} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 \quad (3)$$

To solve Eq. (3), $\mathbf{A}^\top \mathbf{A}$ and $\mathbf{A}^\top \mathbf{b}$ are efficiently computed on the GPU using a parallel reduction [13]. The resulting 6×6 linear system $\mathbf{A}^\top \mathbf{A}\mathbf{x} = \mathbf{A}^\top \mathbf{b}$ is efficiently solved using a Singular Value Decomposition (SVD) performed on the CPU.

3.1. Regularized formulation of the problem

After leveraging the orientation prior, we propose to solve the regularized formulation shown in Eq. (4).

$$\mathbf{x}_{opt} = \min_{\mathbf{x}} \frac{1}{2n} \|\mathbf{A}\mathbf{x} - \mathbf{b}\|_2^2 + \lambda \|\mathbf{P}\mathbf{x}\|_2^2. \quad (4)$$

We set $\mathbf{P} = [\mathbf{I}_3 | \mathbf{0}_3]$ to regularize only the 3 angles α, β and γ to be close to zero. Such a constraint enforces the small angle hypothesis presented in the original formulation (2), but also allows for small rotations to compensate for the noise in the IMU orientation measurement. The regularized linear system that we solve is presented in Eq. (5).

$$(\mathbf{A}^\top \mathbf{A} + 2\lambda n \mathbf{P}^\top \mathbf{P})\mathbf{x} = \mathbf{A}^\top \mathbf{b} \quad (5)$$

The regularized term $2\lambda n \mathbf{P}^\top \mathbf{P}$ is efficiently computed on the GPU. λ is a hyper-parameter that leverages the use of the IMU prior; it can be fixed (constant) or a function of n . The regularized formulation of the linearized problem is solved using an SVD on the CPU. The process is repeated in an ICP framework.

3.2. Distribution of closest-point distances

In the original problem, the number of *good* matches n in Eq. (4) is unknown. To both estimate n and identify inconsistent correspondences, we exploit the probability density function (PDF) of closest point pairs, which gathers the point-to-point distances into a histogram. The PDF and its cumulative (CDF) are efficiently computed on the GPU and are used in the next steps. Figure 2 shows a PDF extracted from a single iteration in the ICP framework.

3.3. Median-based filtering

We also propose a filtering method for the closest-point correspondences based on the median distance. This median is estimated from the CDF of the distances obtained up to the resolution of the used histogram. We argue that, by correcting the orientation, the distribution of the closest point distances should be clustered around the distance corresponding to the translation between the two point clouds. In other words, the distances should be gathered around $\bar{d} = \sqrt{t_x^2 + t_y^2 + t_z^2}$ after orientation correction.

The closest point estimation is performed using a normal shooting technique, with normals estimated using an integral image method [14]. Such a normal estimation is prone to error especially on the border of the depth frame and on the edges of objects in the scene. The algorithm may look for the closest point along a wrong direction. We thus filter out pairs of points that do not gather around the median distance, identifying them as erroneous. Figure 2 illustrates such median-based filtering. Here, the PDF of the distances is plotted with the good pair distances in green and the filtered out ones in red.

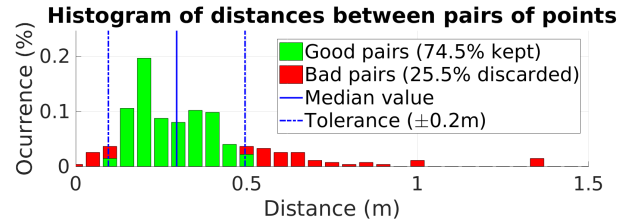


Figure 2. Median Filtering performed on a distance distribution. In this specific case, we are keeping around 75% of the pairs.

3.4. Convergence verification

The median distance used for the filtering is subsequently used to control the convergence of the ICP. Since we seed the ICP with the orientation, we observe that the coarse-to-fine ICP requires less iterations for convergence than before. To assess the proper alignment, we loop until the median converges instead of performing a fixed number of iterations. The convergence of a given level of the coarse-to-fine ICP is achieved when the median no longer changes within its resolution for 3 successive iterations.

4. Experiments

We use the state-of-the-art IMU BNO055 manufactured by Bosch, which is typically designed for embedded applications, such as flight control and motion capture. It has a small footprint of a few mm² and comes on a 1cm-sized board. We mounted the IMU on a Kinect V1 (Figure 3 right) and its pose has been calibrated with the reference system of the 3D camera using the hand-to-eye calibration method proposed by Tsai et al. [40]. The Kinect V1 grabber has been modified to include the IMU orientation in the point cloud. We used the KinectFusion implementation available in the Point Cloud Library [36].



Figure 3. BNO055 board on a Kinect V1 RGB-D Camera.

4.1. Evaluation of ICP Variants

We tested different versions of ICP using an illustrative *Desktop* example for the registration, shown in Figure 4. Here, we consider two consecutive frames in the *Desktop* sequence and evaluate how well each ICP variant registers the input point cloud to the target point cloud.

From Figure 4c, it is evident that ICP may converge to an undesirable local minimum if the IMU is not used, especially if the initial guess of the matched point correspondences is unreliable. Here, even though the IMU-Based metrics have a worse initial value, it lies in a more linear and more convex neighbourhood of the global optimal solution, hence it converges faster and more robustly to the desired solution.

Figure 5 compares the results of the alignment of two point clouds, as well as, the distribution of the closest point distance over iterations, using a plain ICP and our IMU-Based ICP both with and without regularization. In this example, the original ICP gets stuck in a local minimum and does not align the two point clouds properly (Figure 5a); the distribution of distances does not converge after 50 iterations. In Figures 5b and 5c we show how our system converges when initialized with the orientation measurement. Using a strong regularization (b), the orientation is not optimized, the distribution does converge, but, in this case, the translation \mathbf{t} is computed without taking into account the eventual noise in the orientation. Without regularization (c), the robustness is decreased, i.e. the ICP can converge to a local minimum, but an eventual error in the orientation measurement is corrected, which produces a narrower PDF.

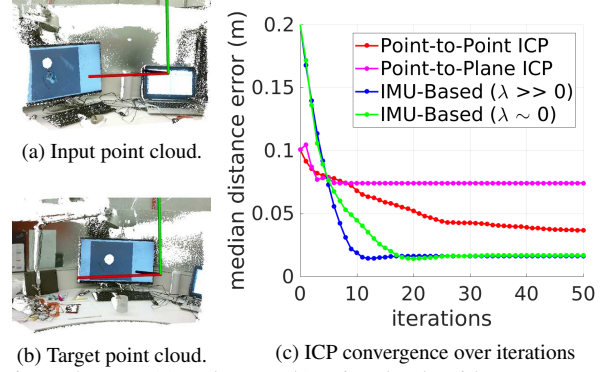


Figure 4. Input (a) and target (b) point clouds with IMU measurements, the convergence (c) of ICP metric (median) without IMU vs. IMU-Based methods.

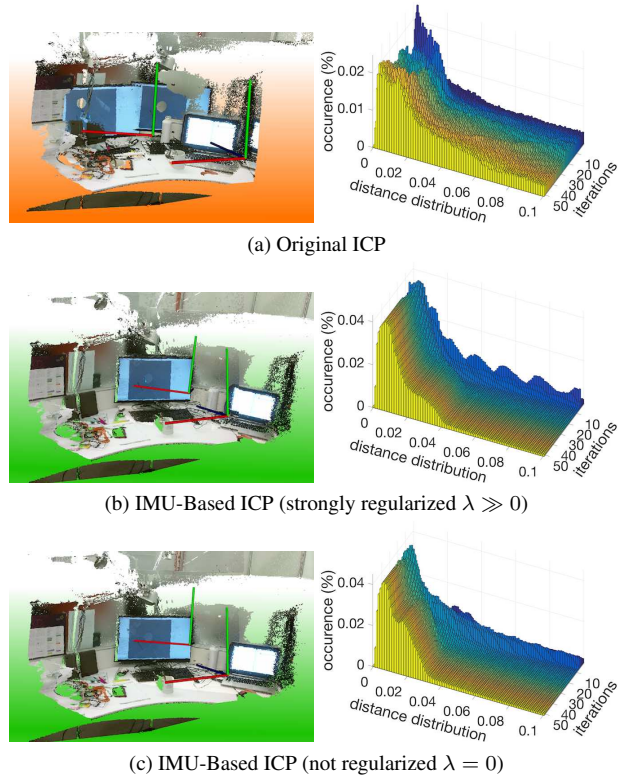


Figure 5. ICP alignment and distance distribution over iteration.

4.2. Angular convergence

In order to show the sensitivity of KinectFusion to angular motion, we considered 20 consecutive frames from the *Desktop* environment. Between pairs of frames, we rotate the 3D camera by a specific angle $\theta \in \{5^\circ, 10^\circ, \dots, 60^\circ\}$. The rationale is to simulate angular shifts in acquisition. We then feed the data to KinectFusion and inspect the resulting 3D reconstruction for obvious errors such as duplication of scene geometry. Figure 6 plots the median of the residual error over iterations, for a predefined set of intra-frame angles. Clearly, KinectFusion's original ICP cannot

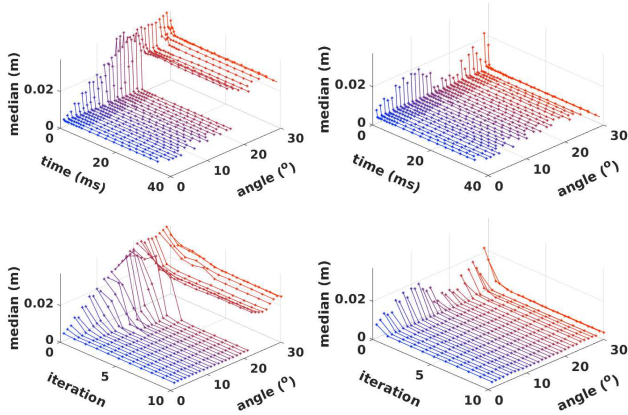


Figure 6. Original ICP (**left**) and our implementation (**right**) convergence over time (**top**) and iteration (**bottom**) for different intra-frame angles.

Table 1. Failure rate across intra-frame acquisition angles.

Desktop	5°	10°	15°	20°	25°	30°	40°	50°	60°
Original	0%	10%	25%	55%	75%	90%	100%	100%	100%
Ours	0%	0%	10%	0%	10%	0%	5%	10%	100%

cope with orientation changes of more than 20° between consecutive frames and shows an increasing degradation. In contrast, our IMU-based approach maintains the same convergence quality over angle.

Table 1 summarizes the failure rate (averaged across 20 pairs of frames) for each angular shift during the acquisition. Clearly, our IMU-based version of KinectFusion is more robust, since it has a much lower failure rate for most rotations. Our method degrades more gracefully than classical KinectFusion and total failure only occurs at $\theta = 60^\circ$, when the point clouds do not overlap anymore.

4.3. Improvements in reconstruction

For the same reason we have mentioned above, the reconstruction deteriorates with the scanning pace; a fast scanning at a fixed frame rate creates larger gaps between consecutive frames than a slow scanning. Figure 8 shows the ICP convergence for a slow and fast pace reconstruction. Clearly, for slow motion (small angles between consecutive frames), both versions of ICP tend to converge to similar results, with the IMU-based version being faster to converge (refer to Figure 8a). However, for fast motion (large angles between consecutive frames), ICP tends not to converge in the 19 iterations at its disposal (refer to Figure 8b). Because error rampantly propagates, a single instance where acquisition is too fast can lead to a completely erroneous model (refer to Fig. 1 (red background)). In comparison, KinectFusion with IMU-based ICP converges to a much better solution and in a shorter time for the same acquisition (refer to Fig. 1 (green background)).

In order to assess our hypothesis, we asked different

users to perform several reconstructions on multiple scenes with our setup. To prevent bias, the users were not given intricate instructions about the acquisition strategy they should use. Note that the users were visualizing the reconstruction as it was being incrementally created. As a result, they scanned the scene at a speed that they considered suitable for reconstruction.

Figure 1 shows some qualitative results for the *Kitchen* dataset, while Figure 7 shows results for the *People*, *Showcase* and *Desktop* environments. A red background indicates an erroneous reconstruction performed by the original KinectFusion algorithm and a green background shows the reconstruction performed by our implementation. In these particular examples, we only used the IMU as a seed to the ICP without any regularization, median filtering nor convergence check. We can already see an improvement in the reconstruction quality.

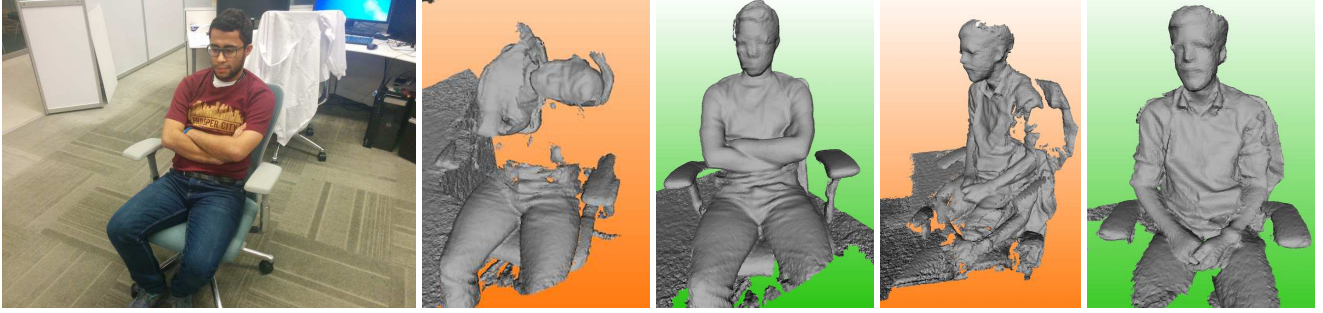
We notice for the *People* model in Figure 7a that a fast rotation along the camera principal axis (model on the middle) or an horizontal axis (model on the right) creates an inconsistent reconstruction. At a certain point, the coarse-to-fine ICP is not able to align the newly acquired point cloud with the model within the 19 iterations at its disposal. As a result, it duplicates the scene and overwrites the weight of the TSDF cube that stores the reconstruction. For the *Showcase* scene in Figure 7b, a fast motion aligned a frame on the wrong seat, creating a shift in the overall reconstruction, duplicating the objects. Using the orientation provided by the IMU helped in keeping track of the reconstruction. For the *Desktop* scene shown in Figure 7c, a bad alignment has created a duplication of the desktop in two different levels. The first part of the point cloud streaming reconstructed a model of the desk, but after a bad alignment occurred, a second model was built.

For all the previous examples, our implementation (green background) was able to cope with the difficulties of the different datasets, showing improved robustness in the presence of fast motion.

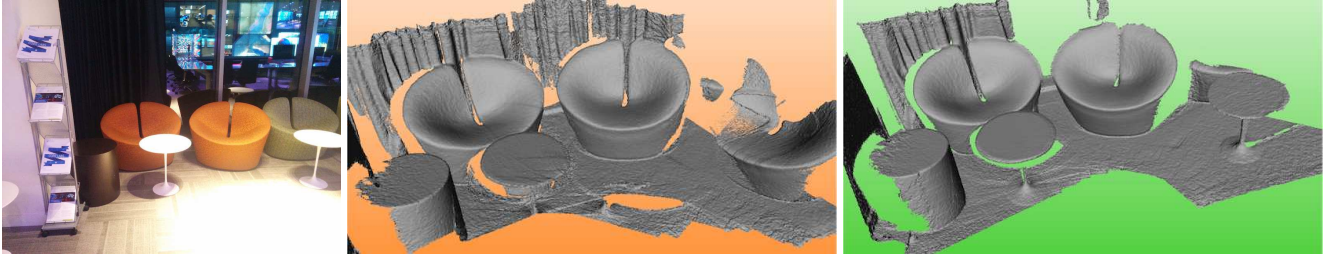
4.4. Considerations on convergence timings

The original KinectFusion method uses a fixed cadence (4, 5, 10) of coarse-to-fine ICP iterations, since computing the residual error on a GPU is an expensive *log-reduce* operation. By feeding the ICP with an initial orientation from the IMU, we can reduce the number of iteration to (2, 2, 3) without reducing the performance. Table 2 shows the improvement in time, which corresponds to 61%–70% of the time used for the ICP alone and 35%–43% of the overall KinectFusion pipeline.

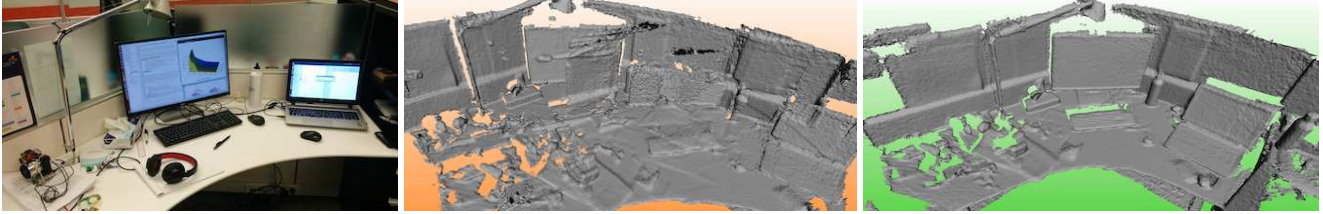
Since less ICP iterations result in improved speed, we can use this saved time to estimate the PDF of the closest-point distances. Such an operation can be time consuming if performed on the CPU, so we take advantage of the multiple



(a) People (~150 frames)

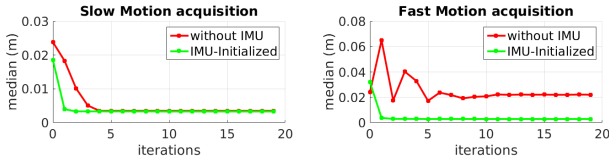


(b) Showcase (~300 frames)



(c) Desktop (~200 frames)

Figure 7. KinectFusion reconstruction of multiple scenes, each with (green background) and without IMU (red background).



(a) ICP convergence for slow motion (b) ICP convergence for fast motion
Figure 8. ICP convergence plot representative for more than 95% of our data, for slow (a) and fast (b) camera motion. Adding an IMU substantially improves both cases.

Table 2. KinectFusion runtime with custom number of iterations on NVIDIA K6000 and GTX850M GPUs.

	K6000		GTX850M	
	ICP	KinFu	ICP	KinFu
Original (4,5,10)	10.58 ms	15.74 ms	27.17 ms	43.41 ms
Ours (2,2,3)	4.07 ms	10.22 ms	8.22 ms	24.72 ms
Improvement	-61.53 %	-35.07 %	-69.75 %	-43.05 %

GPU cores available to build the distance PDF on the GPU as well. Figure 9 shows the timing for different versions of KinectFusion, while reconstructing the Freiburg scenes (a

total of 15,793 alignments). In blue, we show the distribution of the time required by the original KinectFusion, with an average of 15.55ms. Adding the estimation of the PDF at each ICP iteration (red distribution), the average time increases to 25.62ms. However, once we use the PDF to estimate its median value and stop the alignment after a convergence is detected, our method requires only 11.54 iterations on average to converge, whereas the original KinectFusion requires 19 iterations. Since our method converges in less iterations, it only requires an overall average of 13.71ms for alignment, which corresponds to an improvement of 11.8% on the original time. Furthermore, Table 3 shows the timing details for different datasets.

5. Performance on the Freiburg Benchmark

We evaluate our implementation on the Freiburg dataset [38], which is commonly used to benchmark large-scale SLAM techniques. Since this dataset does not provide any measured orientation, we first characterize the uncertainty model of our IMU in order to apply such a model to the ground truth orientation.

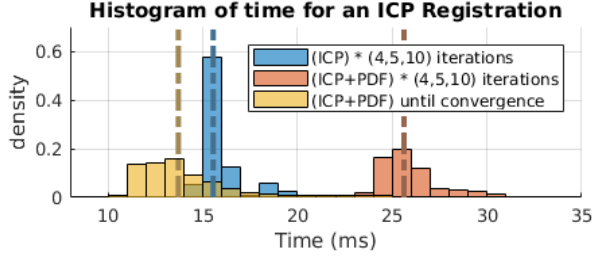


Figure 9. Computational time for the original KinectFusion (blue), KinectFusion with PDF estimation for distant point removal (red) and KinectFusion with PDF estimation and median value convergence check (yellow).

Table 3. KinectFusion runtime for the *original* KinectFusion, our version using a *fixed* number of iteration (4,5,10) and our version with a control of the *convergence*.

Dataset	nPC	Original KinFu	Fixed iteration	Until Converg.	Improvements
fr1/desk	571	16.0 ms	26.6 ms	12.8 ms	-19.6 %
fr1/desk2	609	16.1 ms	25.8 ms	14.1 ms	-12.6 %
fr1/plant	1117	15.4 ms	25.4 ms	14.9 ms	-3.1 %
fr1/room	1351	15.8 ms	25.7 ms	13.8 ms	-12.3 %
fr1/rpy	689	15.4 ms	25.8 ms	12.9 ms	-16.3 %
fr1/teddy	1398	15.4 ms	25.1 ms	13.1 ms	-15.1 %
fr1/xyz	789	15.8 ms	26.4 ms	13.1 ms	-17.0 %
fr2/desk	2173	15.7 ms	25.8 ms	14.4 ms	-7.9 %
fr2/dishes	2940	15.3 ms	25.1 ms	14.2 ms	-6.8 %
fr2/ms2	1834	15.6 ms	25.5 ms	13.1 ms	-15.8 %
fr3/teddy	2322	15.6 ms	25.3 ms	12.7 ms	-18.0 %
AVG	15793	15.6 ms	25.6 ms	13.7 ms	-11.8 %

5.1. Metrological Characterization of the IMU Orientation

To model the IMU noise and include it in the Freiburg dataset, we characterize the uncertainty of the IMU under static conditions according to the Guide to the Expression of Uncertainty in Measurement (GUM) [6]. We used a 6-DoF anthropomorphic robot ABB IRB 1200 to stress the IMU in rotation along the three main axes, within a $\pm 180^\circ$ range and a 5° angular step. The reference orientation has a 0.01° resolution on each axis, which is several orders of magnitude better in resolution than the expected IMU uncertainty.

For the static characterization, the IMU was positioned in space with a controlled orientation with respect to its reference system (Earth). The **Z**-axis vertical was aligned to *top*, and the **X**- and **Y**- axes orthogonal and horizontal. In particular, the **X**- axis was aligned with the *north* direction. After a stabilization time of 5 seconds, 1000 orientation measurements are performed at the IMU frequency of 100Hz. Uncertainty was measured as a combination of a *random error* represented by the standard deviation of the 1000 measurements and a *systematic error* measured by comparison to the ground truth orientation provided by the robotic arm.

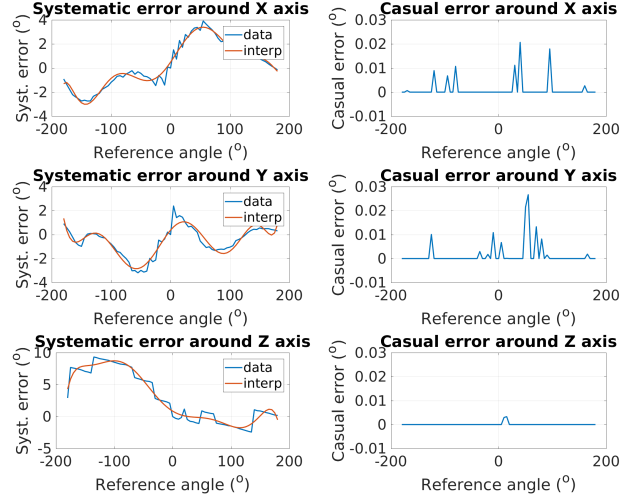


Figure 10. Uncertainty characterization of the BNO055 IMU along **X**-, **Y**- and **Z**- axes for the systematic (top) and random (bottom) errors. Note the different scales.

Results presented in Figure 10 show that the uncertainty has an important systematic error, but no random error. The absence of a random component in the uncertainty is due to a strong low pass filter occurring on the IMU for the orientation estimation. The systematic error is mainly due to the internal calibration between the sensors. For reference, the IMU orientation around the horizontal **X**- and **Y**- axes are extracted using the combination of a magnetometer, an accelerometer, and a gyroscope, which provide accuracy around $\pm 3^\circ$. Orientation around the vertical **Z**-axis is extracted using the magnetometer only, which provides lower quality measurement with uncertainty measured to be about $\pm 10^\circ$. The test was repeated several times with different initial orientations for the IMU leading to similar trends.

5.2. Methods and Results

We used multiple datasets from the Freiburg benchmark [38], focusing on 3D object reconstructions rather than large environments. The list of datasets we used is shown in Table 5. We adapted the parameters of the KinectFusion such that the reconstructions can fit inside a $512 \times 512 \times 512$ TSDF cube, with a maximum side length of 6m. We chose identical parameters to run instances of the original KinectFusion and our pipeline using the ground truth orientation, with and without the systematic error that we characterized in the uncertainty model.

We tested our model over different values of λ to figure out the best regularization. Table 4 shows the average improvement for Absolute Trajectory Error (ATE) metrics (in %) with respect to the original KinectFusion algorithm, over all the datasets used. Our results clearly indicate that larger regularization results in better accuracy.

Table 5 details the ATE and the Relative Pose Estimation

Table 4. Comparison of the average RMSE Absolute Trajectory Error (ATE) of our method respect to the original KinectFusion one using different regularizations parameter λ . Values are given in percentage of improvement over all the datasets.

$c =$	0.05	0.20	1.00	2.00	5.00	AVG
c	-30%	5%	-29%	-33%	-44%	-26%
c/\sqrt{n}	-26%	-31%	-33%	-46%	-46%	-36%
c/n	-6%	-11%	-51%	-41%	-29%	-28%
c/n^2	20%	-2%	-45%	-34%	-41%	-20%
$-c \ln(n)$	72%	-32%	-16%	-34%	-41%	-10%
AVG	6%	-14%	-35%	-38%	-40%	-24%

(RPE) metrics for the different datasets. We only show the improvement on the translational part, since the rotational part is directly measured. We can see that the average ATE error is reduced by 53% on the selected dataset and the RPE is reduced by 21%. Figure 11 plots the trajectories resulting from our method. They are consistent with the ground truth. When using the original KinectFusion method, part of the trajectory is not reconstructed due to tracking loss.

Table 5. ATE and RPE metrics for our method compared to the original KinectFusion. Here we used regularization ($\lambda = 5$), median filtering and convergence control

dataset	ATE (RMSE) (m)			RPE (RMSE) (m)		
	Orig.	IMU	+noise	Orig.	IMU	+noise
fr1/desk	0.073	0.030	0.044	0.738	0.745	0.739
fr1/desk2	1.162	0.293	0.444	1.288	0.838	0.870
fr1/plant	0.569	0.134	0.176	0.707	0.655	0.632
fr1/room	0.533	0.292	0.370	0.573	0.568	0.576
fr1/rpy	0.144	0.026	0.039	0.190	0.147	0.155
fr1/teddy	0.772	0.033	0.136	1.107	0.562	0.552
fr1/xyz	0.022	0.018	0.032	0.471	0.455	0.443
fr2/desk	1.399	1.087	1.107	0.513	0.325	0.411
fr2/dishes	1.118	0.268	0.619	0.394	0.239	0.323
fr2/ms2	1.126	0.141	0.142	0.426	0.321	0.329
fr3/teddy	0.396	0.097	0.292	0.430	0.392	0.397
AVG	0.665	0.220	0.309	0.622	0.477	0.494
Improv.	-	-67%	-53%	-	-23%	-21%

6. Conclusion

In this work, we have demonstrated the benefits of integrating a cheap and modular IMU into the popular KinectFusion reconstruction pipeline. We used the IMU orientation to seed the ICP algorithm, which allows us to linearize/convexify the original problem more faithfully. We used a regularized point-to-plane metric that constrains the orientation within boundaries. We made sure the chosen correspondences are consistent and free of outliers by exploiting their median distance as a basis for outlier removal. We showed qualitative and quantitative improvements in the robustness of our modified KinectFusion pipeline over the original KinectFusion. In addition to improved reconstruction quality, speed is also improved by almost 12%. This

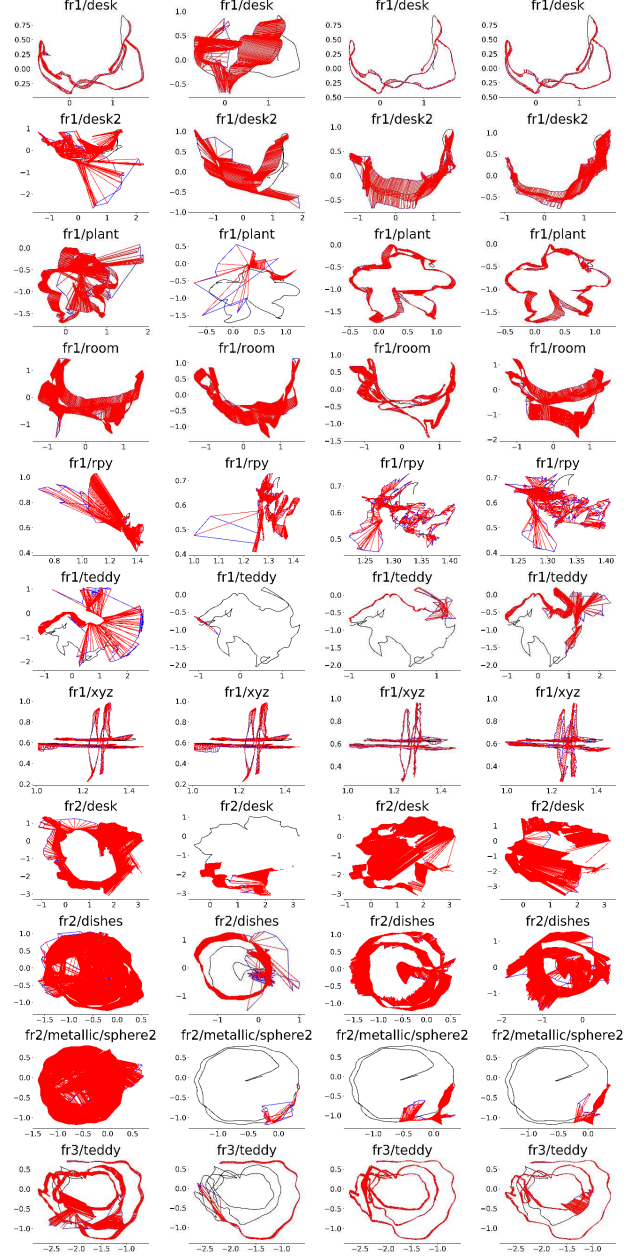


Figure 11. **From top to bottom:** Trajectories for some Freiburg datasets. (**black:** ground truth, **blue:** estimated, **red:** difference) **From left to right:** Original KinectFusion / our method without regularization ($\lambda = 0$) / adding regularization ($\lambda = 5$) and convergence check / using noisy orientation.

is due to a significant reduction in the number of ICP iterations needed for convergence. Finally, results on the Freiburg benchmark show that the overall quality of the tracking/reconstruction is improved by a factor of 2.

Acknowledgments. This work was supported by the King Abdullah University of Science and Technology (KAUST) Office of Sponsored Research and the Visual Computing Center (VCC).

References

- [1] P. Axelsson. Processing of laser scanner data-algorithms and applications. *ISPRS Journal of Photogrammetry and Remote Sensing*, 54(2):138–147, 1999. 1
- [2] P. J. Besl and N. D. McKay. Method for registration of 3-d shapes. In *Robotics-DL tentative*, pages 586–606. International Society for Optics and Photonics, 1992. 1
- [3] N. Brunetto, S. Salti, N. Fioraio, T. Cavallari, and L. Stefano. Fusion of inertial and visual measurements for rgb-d slam on mobile devices. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 1–9, 2015. 2
- [4] A. Concha, G. Loianno, V. Kumar, and J. Civera. Visual-inertial direct slam. In *Robotics and Automation (ICRA), 2016 IEEE International Conference on*, pages 1331–1338. IEEE, 2016. 2
- [5] A. Corti, S. Giancola, G. Mainetti, and R. Sala. A metrological characterization of the kinect v2 time-of-flight camera. *Robotics and Autonomous Systems*, 75:584–594, 2016. 1
- [6] B. I. des Ponds et Mesures. Evaluation of measurement data guide for the expression of uncertainty in measurement. JCGM 100: 2008. 2, 7
- [7] H. Durrant-Whyte and T. Bailey. Simultaneous localization and mapping: part i. *Robotics & Automation Magazine, IEEE*, 13(2):99–110, 2006. 1
- [8] S. Foix, G. Alenya, and C. Torras. Lock-in time-of-flight (tof) cameras: a survey. *Sensors Journal, IEEE*, 11(9):1917–1926, 2011. 1
- [9] C. Forster, L. Carlone, F. Dellaert, and D. Scaramuzza. Imu preintegration on manifold for efficient visual-inertial maximum-a-posteriori estimation. Georgia Institute of Technology, 2015. 2
- [10] J. Geng. Structured-light 3d surface imaging: a tutorial. *Advances in Optics and Photonics*, 3(2):128–160, 2011. 1
- [11] S. Giancola, D. Piron, P. Poppa, and R. Sala. A Solution for Crime Scene Reconstruction using Time-of-Flight Cameras. *ArXiv e-prints*, Aug. 2017. 1
- [12] H. Giberti, M. Tarabini, F. Cheli, and M. Garozzo. Accuracy enhancement of a device for automated underbridge inspections. In *Structural Health Monitoring, Damage Detection & Mechatronics, Volume 7*, pages 59–66. Springer, Cham, 2016. 1
- [13] M. Harris, S. Sengupta, and J. D. Owens. Parallel prefix sum (scan) with cuda. *GPU gems*, 3(39):851–876, 2007. 3
- [14] S. Holzer, R. B. Rusu, M. Dixon, S. Gedikli, and N. Navab. Adaptive neighborhood selection for real-time surface normal estimation from organized point cloud data using integral images. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 2684–2689. IEEE, 2012. 3
- [15] G. Huang, M. Kaess, and J. J. Leonard. Towards consistent visual-inertial navigation. In *Robotics and Automation (ICRA), 2014 IEEE International Conference on*, pages 4926–4933. IEEE, 2014. 2
- [16] S. Izadi, D. Kim, O. Hilliges, D. Molyneaux, R. Newcombe, P. Kohli, J. Shotton, S. Hodges, D. Freeman, A. Davison, et al. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th annual ACM symposium on User interface software and technology*, pages 559–568. ACM, 2011. 1
- [17] K. Khoshelham and S. O. Elberink. Accuracy and resolution of kinect depth data for indoor mapping applications. *Sensors*, 12(2):1437–1454, 2012. 1
- [18] K. Konolige, M. Agrawal, and J. Sola. Large-scale visual odometry for rough terrain. In *Robotics research*, pages 201–212. Springer, 2010. 2
- [19] K. Konolige, E. Rublee, S. Hinterstoisser, T. Straszheim, G. Bradski, and H. M. Strasdat. Detection and reconstruction of an environment to facilitate robotic interaction with the environment, Apr. 25 2017. US Patent 9,630,320. 1
- [20] S. Leutenegger, S. Lynen, M. Bosse, R. Siegwart, and P. Furgale. Keyframe-based visual-inertial odometry using nonlinear optimization. *The International Journal of Robotics Research*, 34(3):314–334, 2015. 2
- [21] M. Li and A. I. Mourikis. High-precision, consistent ekf-based visual-inertial odometry. *The International Journal of Robotics Research*, 32(6):690–711, 2013. 2
- [22] G. Ligorio and A. M. Sabatini. Extended kalman filter-based methods for pose estimation using visual, inertial and magnetic sensors: Comparative analysis and performance evaluation. *Sensors*, 13(2):1919–1941, 2013. 2
- [23] K.-L. Low. Linear least-squares optimization for point-to-plane icp surface registration. *Chapel Hill, University of North Carolina*, 4, 2004. 1, 3
- [24] L. Ma, J. M. Falquez, S. McGuire, and G. Sibley. Large scale dense visual inertial slam. In *Field and Service Robotics*, pages 141–155. Springer, 2016. 2
- [25] A. Martinelli. Vision and imu data fusion: Closed-form solutions for attitude, speed, absolute scale, and bias determination. *IEEE Transactions on Robotics*, 28(1):44–60, 2012. 2
- [26] A. I. Mourikis and S. I. Roumeliotis. A multi-state constraint kalman filter for vision-aided inertial navigation. In *Robotics and automation, 2007 IEEE international conference on*, pages 3565–3572. IEEE, 2007. 2
- [27] P. Musialski, P. Wonka, D. G. Aliaga, M. Wimmer, L. v. Gool, and W. Purgathofer. A survey of urban reconstruction. In *Computer graphics forum*, volume 32, pages 146–177. Wiley Online Library, 2013. 1
- [28] R. A. Newcombe, S. Izadi, O. Hilliges, D. Molyneaux, D. Kim, A. J. Davison, P. Kohli, J. Shotton, S. Hodges, and A. Fitzgibbon. Kinect-fusion: Real-time dense surface mapping and tracking. In *Mixed and augmented reality (ISMAR), 2011 10th IEEE international symposium on*, pages 127–136. IEEE, 2011. 1
- [29] M. Nießner, A. Dai, and M. Fisher. Combining inertial navigation and icp for real-time 3d surface reconstruction. In *Eurographics (Short Papers)*, pages 13–16, 2014. 2
- [30] D. Nistér. Preemptive ransac for live structure and motion estimation. *Machine Vision and Applications*, 16(5):321–329, 2005. 1
- [31] D. Pagliari, F. Menna, R. Roncella, F. Remondino, and L. Pinto. Kinect fusion improvement using depth camera calibration. *The International Archives of Photogrammetry, Remote Sensing and Spatial Information Sciences*, 40(5):479, 2014. 1
- [32] M. K. Paul, K. Wu, J. A. Hesch, E. D. Nerurkar, and S. I. Roumeliotis. A comparative analysis of tightly-coupled monocular, binocular, and stereo vins. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 165–172. IEEE, 2017. 2
- [33] B. Pfrommer, N. Sanket, K. Daniilidis, and J. Cleveland. Penco-syvio: A challenging visual inertial odometry benchmark. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 3847–3854. IEEE, 2017. 2
- [34] U. Qayyum, J. Kim, et al. Inertial-kinect fusion for outdoor 3d navigation. In *Australasian Conference on Robotics and Automation (ACRA)*, 2013. 2
- [35] S. Rusinkiewicz and M. Levoy. Efficient variants of the icp algorithm. In *3-D Digital Imaging and Modeling, 2001. Proceedings. Third International Conference on*, pages 145–152. IEEE, 2001. 1
- [36] R. B. Rusu and S. Cousins. 3D is here: Point Cloud Library (PCL). In *IEEE International Conference on Robotics and Automation (ICRA)*, Shanghai, China, May 9-13 2011. 4
- [37] J. Smisek, M. Jancosek, and T. Pajdla. 3d with kinect. In *Consumer Depth Cameras for Computer Vision*, pages 3–25. Springer, 2013. 1
- [38] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers. A benchmark for the evaluation of rgb-d slam systems. In *Intelligent Robots and Systems (IROS), 2012 IEEE/RSJ International Conference on*, pages 573–580. IEEE, 2012. 2, 6, 7
- [39] S. Thrun et al. Robotic mapping: A survey. *Exploring artificial intelligence in the new millennium*, 1:1–35, 2002. 1
- [40] R. Y. Tsai and R. K. Lenz. A new technique for fully autonomous and efficient 3d robotics hand/eye calibration. *IEEE Transactions on robotics and automation*, 5(3):345–358, 1989. 4

- [41] S. Weiss, M. W. Achtelik, S. Lynen, M. Chli, and R. Siegwart. Real-time onboard visual-inertial state estimation and self-calibration of mavs in unknown environments. In *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, pages 957–964. IEEE, 2012. 2
- [42] K. Wu, A. Ahmed, G. A. Georgiou, and S. I. Roumeliotis. A square root inverse filter for efficient vision-aided inertial navigation on mobile devices. In *Robotics: Science and Systems*, 2015. 2
- [43] J. Zhang and S. Singh. Enabling aggressive motion estimation at low-drift and accurate mapping in real-time. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pages 5051–5058. IEEE, 2017. 2