



HAL
open science

Robust Data Fusion of Multi-modal Sensory Information for Mobile Robots

Vladimír Kubelka, Lorenz Oswald, François Pomerleau, Francis Colas, Tomas Svoboda, Michal Reinstein

► **To cite this version:**

Vladimír Kubelka, Lorenz Oswald, François Pomerleau, Francis Colas, Tomas Svoboda, et al.. Robust Data Fusion of Multi-modal Sensory Information for Mobile Robots. *Journal of Field Robotics*, 2015, 32 (4), pp.447–473. hal-01143471

HAL Id: hal-01143471

<https://hal.science/hal-01143471v1>

Submitted on 17 Apr 2015

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Robust Data Fusion of Multi-modal Sensory Information for Mobile Robots

Vladimír Kubelka

Center for Machine Perception, Dept. of Cybernetics
Faculty of Electrical Engineering, Czech Technical University in Prague
Technicka 2, 166 27, Prague 6, Czech Republic
kubelka.vladimir@fel.cvut.cz

Lorenz Oswald

ETH Zurich
Tannenstrasse 3, 8092 Zurich, Switzerland
loswald@student.ethz.ch

François Pomerleau

ETH Zurich
Tannenstrasse 3, 8092 Zurich, Switzerland
francois.pomerleau@mavt.ethz.ch

Francis Colas

ETH Zurich
Tannenstrasse 3, 8092 Zurich, Switzerland
francis.colas@mavt.ethz.ch

Tomáš Svoboda

Center for Machine Perception, Dept. of Cybernetics
Faculty of Elect. Eng., CTU in Prague
Technicka 2, 166 27, Prague 6, Czech Republic
svobodat@fel.cvut.cz

Michal Reinstein

Center for Machine Perception, Dept. of Cybernetics
Faculty of Elect. Eng., CTU in Prague
Technicka 2, 166 27, Prague 6, Czech Republic
reinstein.michal@fel.cvut.cz

Abstract

Urban Search and Rescue missions for mobile robots require reliable state estimation systems resilient to conditions given by the dynamically changing environment. We design and evaluate a data fusion system for localization of a mobile skid-steer robot intended for USAR missions. We exploit a rich sensor suite including both proprioceptive (inertial measurement unit and tracks odometry) and exteroceptive sensors (omnidirectional camera and rotating laser rangefinder). To cope with the specificities of each sensing modality (such as significantly differing sampling frequencies), we introduce a novel fusion scheme based on Extended Kalman filter for 6DOF orientation and position estimation. We demonstrate the performance on field tests of more than 4.4 km driven under standard USAR conditions. Part of our datasets include ground truth positioning; indoor with a Vicon motion capture system and outdoor with a Leica theodolite tracker. The overall median accuracy of localization—achieved by combining all the four modalities—was 1.2 % and 1.4 % of the total distance traveled, for indoor and outdoor environments respectively. To identify the true limits of the proposed data fusion we propose and employ a novel experimental evaluation procedure based on failure case scenarios. This way we address the common issues like: slippage, reduced camera field of view, limited laser rangefinder range, together with moving obstacles spoiling the metric map. We believe such characterization of the failure cases is a first step towards identifying the behavior of state estimation under such conditions. We

release all our datasets to the robotics community for possible benchmarking.

1 Introduction

Mobile robots are sought to be deployed for many tasks, from tour-guide robots to autonomous cars. With the rapid advance in sensor technology, it has been possible to embed richer sensor suites and extend the perception capabilities. Such sensor suites provide multi-modal information that naturally ensure perception robustness, allowing also better means of self-calibration, fault detection and recovery—given that appropriate data fusion methods are exploited. Independently from the application, a key issue of mobile robotics is state estimation. It is crucial for both perception, like mapping, and action, like avoiding obstacles or terrain adaptation.

In this paper, we address the problem of data fusion for localization of an Unmanned Ground Vehicle (UGV) intended for Urban Search and Rescue (USAR) missions. There has been a significant effort presented in the field of USAR for robot localization that mostly aims for a minimal suitable sensing setup; exploiting usually the inertial measurements aided with either vision or laser data. Having a sufficient on-board computational power, we therefore aim for a richer sensors suite and hence for better robustness and reliability. Therefore, our UGV used in this work (see Figure 1) embeds track encoders, an Inertial Measurement Unit (IMU), an omnidirectional camera, and a rotating laser range-finder.

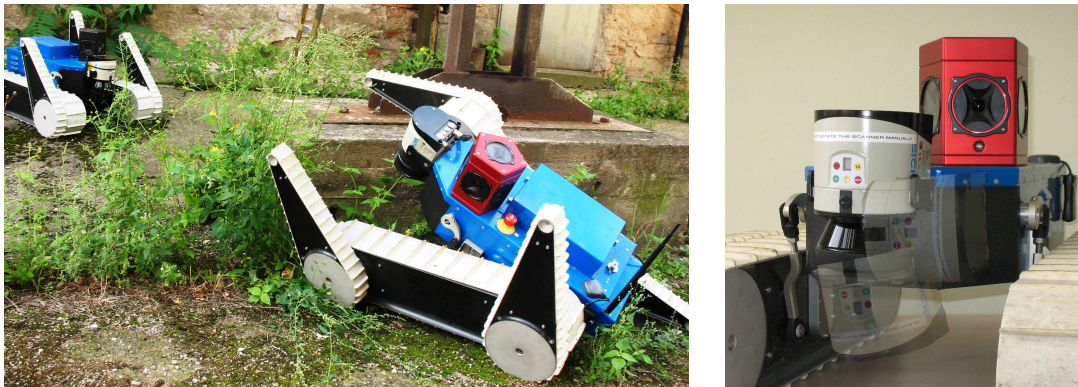


Figure 1: Picture of two USAR UGVs used for experimental evaluation (FP7-ICT-247870 NIFTi project) and a detail of the sensor setup (a PointGrey Ladybug 3 omnicaamera and a rotating SICK LMS-151 laser range finder). See Section 3.1 for more details.

Our first contribution lies in the development of a model for such multi-modal data fusion using the Extended Kalman Filter (EKF), especially in the way we incorporate sensors with slow and fast measurement update rates. In order to cope with such significant difference in the update rates of various sensor modalities, we concentrated the model design on integrating the slow laser and visual odometry with the faster IMU and track odometry measurements. For this purpose, we propose and investigate three different possible methods—one of them, the *trajectory approach* (see Section 4.3.3 for further details), is our contribution that we compare it to the *velocity approach*, which is a common state-of-the-art practice. We show that a standard EKF designed with the *velocity approach* does not cope well with such significant differences in the frequency, whether our proposed *trajectory approach* does.

The context of USAR missions implicitly defines challenges and limitations of our application. The environment is often unstructured (collapsed buildings) and unstable (moving objects or other ongoing changes, deformable terrain causing high slippage). Robots need to cope with indoor-outdoor transitions (change from confined to open spaces), bad lighting conditions with rapid changes and sometimes decreased visibility (smoke and fire). These are essentially the main challenges that come with the sensor data we process.

Therefore, our main contribution lies in the actual experimental evaluation and analysis of limits of the proposed filter. We review the different sensing modalities and their expected failure cases to assess the impact of possible data degradation (or outage) on the overall precision of localization. We believe that the field deployment of state estimation for multi-modal data fusion needs to be characterized both under standard expected conditions as well as for case of partial or full failures of sensing modalities. Indeed, robustness to sensor data outage or degradation is a key element to the scaling up of a field robotics system. Therefore, we evaluate our filter using several hours and kilometers of experimental data validated by indoor or outdoor ground truth measurements. In order to share this contribution to the robotics community, we release all the captured datasets (including the ground truth measurements) to be used as benchmarks.¹

The state of the art of sensor fusion for state estimation is elaborated in Section 2. In Section 3, we present the hardware and software used in this work before describing in details the design of our data fusion algorithm (Section 4). In Section 5, we explain our experimental evaluation including our fail-case methodology before a discussion and conclusion (Section 6).

2 Related work

In general, the information obtained from various sensors can be classified as either proprioceptive (inertial measurements, joint sensors, motor or wheel encoders, etc.) or exteroceptive (Global Positioning System (GPS), cameras, laser range finder, ultrasonic sensors, magnetic compass etc.). Exteroceptive sensors that acquire information from the environment can be also used to perceive external landmarks that are necessary for long-term precision in navigation tasks. In modern mobile robots, a popular solution lies usually in the combination of a proprioceptive component in the form of Inertial Navigation System (INS) (Titterton and Weston, 1997), that captures the body dynamics at high frequency, and an external source of aiding, using vision (Chowdhary et al., 2013) or range measurements (Bachrach et al., 2011). The key issue lies in the appropriate integration of the different characteristics of the different sensor modalities.

As it was repeatedly shown, the combination of an IMU with wheel odometry is a popular technique to localize a mobile robot in a dead reckoning manner. It generally allows very high sampling frequency as well as processing rate, usually without excessive computational load. Dead reckoning can be used for short term navigation without any necessity of perceiving surrounding environment via exteroceptive sensors. In real outdoor conditions, dynamically changing environment often causes signal degradation or even outage of exteroceptive sensors. However, proprioceptive sensing, in principle, is too prone to accumulating errors to be used as a standalone solution. Computational and environmental errors as well as errors caused by misalignment and instrumentation cause the dead reckoning system to drift quickly with time. Moreover, motor encoders do not reflect the true path, especially heading of the vehicle, in case of frequent wheel slip. In (Yi et al., 2007) and (Anousaki and Kyriakopoulos, 2004), an improvement through skid-steer model of a 4-wheel robot is presented, based on a Kalman filter estimating trajectory using velocity constraints and slip estimate. An alternative method appears in (Endo et al., 2007) where the IMU and odometry are used to improve tracked vehicle navigation via slippage estimates. We addressed this problem in (Reinstein et al., 2013). Substantial effort has also been made to investigate the odometry derived constraints (Dissanayake et al., 2001), or innovation of the motion models (Galben, 2011). Concerning all the references so far, localization of the navigated object via dead reckoning was performed only in 2D. There exist solutions providing real 3D odometry derived from the rover-type multi-wheel vehicle design (Lamon and Siegwart, 2004). Nevertheless, the error is still about one order of magnitude higher than what we aim to achieve (below 2% of the total distance traveled).

However, if long-term precision and reliability is to be guaranteed, dead-reckoning solutions require other exteroceptive aiding sensor systems. In the work of (Shen et al., 2011), it is shown that a very low-cost IMU and odometry dead-reckoning system can be realized and successfully combined with visual odometry (VO) (Sakai et al., 2009; Scaramuzza and Fraundorfer, 2011) to produce a reliable navigation system. With

¹The datasets are available as *bagfiles* for ROS at <https://sites.google.com/site/kubelvla/public-datasets>

the increasing on-board computational power, visual odometry is becoming very popular even for large-scale outdoor environments. Most solutions are based on the Extended Kalman filter (EKF) (Oskiper et al., 2010; Civera et al., 2010; Konolige et al., 2011; Chowdhary et al., 2013) or a dimensional-bounded EKF with landmark classifier introduced in (Jesus and Ventura, 2012). However, in (Rodriguez F et al., 2009) it is pointed out that a trade-off between precision and execution time has to be examined. Moreover, VO degrades due to high rotational speed movements and it is susceptible to illumination changes and lack of sufficient scene texture (Scaramuzza and Fraundorfer, 2011).

Another typically used 6 DOF aiding source is a laser rangefinder, which is used for estimating vehicle motion by matching consecutive laser scans and creating a 3D metric map of the environment (Suzuki et al., 2010; Yoshida et al., 2010). Examples of successful application can be found for both indoor—without IMU but combined with vision (Ellekilde et al., 2007)—as well as outdoor—relying on the IMU (Bachrach et al., 2011). As in case of the visual odometry, solutions using EKF are often proposed (Morales et al., 2009; Bachrach et al., 2011). The most popular approach of scan matching is based on the Iterative Closest Point (ICP) algorithm first proposed by (Besl and McKay, 1992) and in parallel by (Chen and Medioni, 1991). More recently, (Nuchter et al., 2007) proposed a 6D Simultaneous Localization and Mapping (SLAM) system relying mainly on ICP. Closer to USAR applications, (Nagatani et al., 2011) demonstrated the use of ICP in exploration missions and used a pose graph minimization scheme to handle multi-robot mapping. (Kohlbrecher et al., 2011) proposed a localization system combining a 2D laser SLAM with a 3D IMU/odometry-based navigation subsystem. Combination of 3D-landmark-based SLAM and multiple proprioceptive sensors is also presented in (Chiu et al., 2013), their work aims mainly on low latency solution while estimating the navigation state by means of Sliding-Window Factor Graph. The problem of utilizing several sensors for localization that may provide contradictory measurements is discussed in (Sukumar et al., 2007). The authors use Bayes filters to estimate sensor measurement uncertainty and sensor validity to intelligently choose a subset of sensors that contribute to localization accuracy. As opposed to the later publications realized in the context of SLAM, we only consider the results of the ICP algorithm as a local pose measurement similarly to (Almeida and Santos, 2013) who use the ICP algorithm to extract steering angle and linear velocity of a car-like vehicle to update its non-holonomic model of motion. In our approach, the 3D reconstruction of the environment is considered locally coherent and neither loop detection nor error propagation is used.

As stated in (Kelly et al., 2012), it is the right time to address concerning issues of the state-of-the-art in long-term navigation and autonomy. In this respect, benefits and challenges of repeatable long-range driving were addressed in (Barfoot et al., 2012). In this context, we believe that bringing more insight into multi-modality state estimation algorithms is an important step for long-term stability of an USAR system evolving in a complex range of environments.

Regarding multi-modal data fusion, we built on our previous work concerning complementary filtering (Kubelka and Reinstein, 2012), odometry modeling (Reinstein et al., 2013), and design of EKF error models (Reinstein and Hoffmann, 2013), even though the later work applied to a legged robot.

3 System description

Our system aims at high state estimation accuracy while ensuring robust performance against rough terrain navigation and obstacle traversals. We selected four modalities to achieve this goal: the inertial measurements (*IMU*), odometry data (*OD*), visual odometry (*VO*) and laser rangefinder data (*ICP*) processed by the ICP algorithm. This section explains the motion capabilities of the Search & Rescue platform and the preprocessing computation applied to its sensors in order to extract meaningful inputs for the state estimation. These explanations provide a motivation for a list of states to be estimated by the EKF described in Section 4.

3.1 Mobile Robotic Platform

Figure 1 presents the UGV designed for USAR mission that we use in this paper. As described in (Kruijff et al., 2012), this platform was deployed multiple times in collaboration with various rescue services (Fire Department of Dortmund/Germany, Vigili del Fuoco/Italy). It has two bogies linked by a differential that allows a passive adaptation to the terrain. On each of the tracks, there are two independent flippers that can be position-controlled in order to increase the mobility in difficult terrain. For example, they can be unfolded to increase the support polygon which helps in overcoming gaps and being more stable on slopes. They can also be raised to help with climbing over higher obstacles. Given that the robot was designed to operate in 3D unstructured environments, the state estimation system needs to provide a 6 DOF localization.

Encoders are placed on the differential, giving the angle between the two bogies and the body; on the tracks to give their current velocity; and on each flipper to give its position with respect to its bogies. Inside the body, vertical to the center of the robot, lies the Xsens MTi-G IMU providing angular velocities and linear acceleration along each of the three axes. The IMU data capture the body dynamics at high rate of 90 Hz. GPS is not taken into account due to the low availability of the signal indoors or in close proximity with building. Magnetic compass is also easily disturbed by metallic masses, pipes, and wires, which make it highly unreliable and hence we do not use it.

The exteroceptive sensors of the robot consist of an omnidirectional camera and a laser rangefinder. The omnidirectional camera is the PointGrey Ladybug 3 and produces a 12 megapixels stitched omni-directional images at 5-6 Hz. The omni-directionality of the sensor provides a stronger stability of rotation estimation at the expense of scale estimation, which would be better handled by a stereocamera. The laser rangefinder used is the Sick LMS-151 mounted on a rolling axis in front of the robot. The laser spins left and right alternately, taking a full 360° scan at approximately 0.3 Hz to create a point cloud of around 55,000 points.

3.2 Inertial data processing

Though the precision and reliability of the IMU measurements is sufficient in short term, in long term the information provided suffers from random drift that, together with integrated noise, cause unbounded error growth. To cope with these errors all the 6 sensor biases have to be estimated (see Section 4.1 for more details). Therefore, we have included sensor biases in the state space of the proposed EKF estimator. Furthermore, correct calibration of the IMU output and its alignment with respect to the robot's body frame has to be assured.

3.3 Odometry for skid-steer robots

Our platform is equipped with caterpillar tracks and therefore steering is realized by setting different velocities for each of the tracks (*skid-steering*). The encoders embedded in the tracks of the platform measure the left and right track velocities at approximately 15 Hz. However, in contradistinction to differential robots, the odometry for skid-steering vehicles has significant uncertainties. Indeed, as soon as there is a rotation, the tracks must either deform or slip significantly. The slippage is affected by many parameters including the type and local properties of the terrain. To keep the computation complexity low, we assume only a simple odometry model and we do not model the slippage. Instead, we take advantage of the exteroceptive modalities in our data fusion to observe the true motion dynamics using different sources of information. Hence, the fusion compensates for cases when the tracks are slipping because the surface is slippery or because of an obstacle blocking the robot. Another advantage of using caterpillar tracks odometry lies in the opportunity to exploit nonholonomic constraints. Further explanations on those constraints are given in Section 4.3.

3.4 ICP-based localization

Using as *Input* the current 3D point cloud, a registration process is used to estimate the pose of the robot with respect to a global representation called *Map*. We used a derivation of the point-to-point ICP algorithm introduced by (Chen and Medioni, 1991) combined with the trimmed outlier rejection presented by (Chetverikov et al., 2002).

The implementation uses `libpointmatcher`², an open-source library fast enough to handle real-time processing while offering modularity to cover multiple scenarios as demonstrated in (Pomerleau et al., 2013). The complete list of modules used with their main parameters can be found in Table 1. In more details, the configuration of the rotating laser produced a high density of points in front of the robot, which was desirable to predict collision but not beneficial to the registration minimization. Thus, we forced the maximal density to 100 points per m³ after having randomly subsampled the point cloud in order to finish the registration and the map maintenance within 2s. We expected the error on pre-alignment of the 3D scans to be less than 0.5m based on the velocity of the platform and the number of ICP per second that was to be executed. So we used this value to limit the matching distance. We also removed paired points with an angle difference larger than 50° to avoid the reconstruction of both sides of walls to collapse when the robot was exploring different rooms. The surface normal vector used for the *outlier filtering* and for the *error minimization* are computed using 20 Nearest Neighbors (NN) of every point within a single point cloud. As for the global map, we maintained a density of 100 points per m³ every time a new input scan was merged in it. A maximum of 600,000 points were kept in memory to avoid degradation of the computation time when exploring a larger environment than expected. However, the only output of the ICP algorithm we consider is the robot’s localization, i.e. position and orientation relative to its inner 3D point-cloud map. We do not aim at creating a globally consistent map and we do not exploit the map in any other way than for analysis of the ICP performance (no map corrections or loop closures are performed).

Table 1: Configurations of ICP chains for the NIFTi mapping applications.

	<i>Step</i>	<i>Module</i>	<i>Description</i>
Input	Read. filtering	SimpleSensorNoise SamplingSurfaceNormal ObservationDirection OrientNormals MaxDensity	SickLMS keep 80 %, surface normals based on 20 NN add vector pointing toward the laser orient surface normals toward the obs. direction subsample to keep point with density of 100 pts/m ³
	Ref. filtering	-	processing from the rows Map
	Read. filtering	-	processing from the rows Input
	Data association	KDTree	kd-tree matching with 0.5 m max. distance, $\epsilon = 3.16$
Registration	Outlier filtering	TrimmedDist SurfaceNormal	keep 80 % closest points remove paired normals angle > 50°
	Error min.	PointToPlane	point-to-plane
	Trans. checking	Differential Counter Bound	min. error below 0.01 m and 0.001 rad iteration count reached 40 transformation fails beyond 5.0 m and 0.8 rad
	Ref. filtering	SurfaceNormal MaxDensity MaxPointCount	Update normal and density, 20 NN, $\epsilon = 3.16$ subsample to keep point with density of 100 pts/m ³ subsample 70 % if more than 600,000 points
	Map		

There is one ICP-related issue observed with our platform. Although the ICP creates a locally precise metric map, the map as whole tends to slightly twist or bend (we do not perform any loop-closure). This is the reason why the position and the attitude estimated by the ICP odometry collide with other position information sources. Another limitation is the refresh rate of the pose measurements limited to 0.3 Hz. This rate is far from our fastest measurement (i.e., the IMU at 90 Hz), which poses a linearization problem. For these reasons, we investigated three different types of measurement models; see Section 4.3.3 for details.

Furthermore, the true bottleneck of the ICP-based localization lies in the way it is realized on our platform

²<https://github.com/ethz-asl/libpointmatcher>

and hence prone to mechanical issues. As the laser rangefinder has to be turning to provide full 3D point cloud, in environment with high vegetation such mechanism is easily struck, causing this modality to fail. Large open spaces, indoor / outdoor transitions, or significantly large moving obstacles can also cause the ICP to fail updating the metric map. Since this modality is very important, we analyzed these failure cases in Section 5.4.

3.5 Visual odometry

Our implementation of visual odometry generally follows the usual scheme (Tardif et al., 2008; Scaramuzza and Fraundorfer, 2011). The VO computation runs solely on the robot on-board computer and estimates the pose at the frame rate 2-3Hz which, compared to the robot speed, is sufficient. It does search for correspondences (i.e., image matching) (Rublee et al., 2011), landmark reconstruction and sliding bundle adjustment (Kummerle et al., 2011; Fraundorfer and Scaramuzza, 2012), which refines the landmark 3D positions and the robot poses. The performance essentially depends on the visibility and variety of landmarks. The more variant landmarks are visible at more positions, the more stable and precise is the pose estimation. The process uses panoramic images constructed from spherical approximation of the Ladybug camera model. The Ladybug camera is approximated as one central camera. The error of the approximation is acceptable for landmarks which are few meters from the robot.

The visual odometry starts with detecting and matching features in two consecutive images. We use OpenCV implementation of the Orb keypoint detector and descriptor (Rublee et al., 2011). Only the matches, which are distinctive above certain threshold, survive. The initial matching is supported by a guided matching which uses an initial estimate of the robot movement. The robot movement is estimated by the 5-point solver (Li and Hartley, 2006) encapsulated in RANSAC iterations. As the error measure we use the angular deviation of points from epipolar planes. This is less precise than the usual distance from epipolar lines. However, as we work with spherical projection we have epipolar curves. Computing angular deviations is faster than computing distance to the epipolar curve. The movement estimate projects already known landmarks and we can actively search around the projection. The feature tracks are updated and associated with landmarks if they pass an observation consistency test. The landmark 3D position is triangulated from all possible observations and the complete estimate of landmark and robot positions are refined by a bundle adjustment (Kummerle et al., 2011).

Using an almost omnidirectional camera for the robot motion estimation is geometrically advantageous (Brodsky et al., 1998; Svoboda et al., 1998). The scale estimation however, depends on the precision of 3D reconstruction where the omnidirectionality does not really help. It is also important to note the omnidirectional camera we use sits very low above the terrain (below 0.5 m) and directly on the robot body. This makes a huge difference compared to, e.g. (Tardif et al., 2008), where the camera is more than 2 m above the terrain and sees the ground plane much better than our camera. Estimation of the yaw angle is still well conditioned since it relies mostly on the side correspondences. The pitch estimation however, would sometimes need more landmarks on the ground plane. The pitch part of the motion induces largest disparity of the correspondences in the front and back cameras. Unfortunately the back view is significantly occluded by the battery cover. This is especially problem in the street scenes where the robot moves along the street, see e.g., Figure 11. The front cameras see the street level better however, the uniform texture of the tar surface often generates only a few reliable correspondences. The search for correspondences is further complicated by the tilting flippers which occlude the field of view and induce outliers. Second problem is the agility of the robot combined with relatively low frequency of the visual odometry. The robot can turn on spot very quickly, much quicker than an ordinary wheeled car. Even worse, the quick turn is the usual way how the movement direction is changed. This all makes correspondence search difficult. In the future versions of the visual odometry we want to improve the landmark management in order to resolve the problem of too few landmarks surviving the sudden turn. We also think about replacing the approximate spherical model by reformulating it in a multiview model.

4 Multi-modal data fusion

The core of the data fusion system is realized by an error state EKF inspired by the work of (Weiss, 2012). The description of the multi-modal data fusion solution we propose can be divided into two parts. First is the process error model for the EKF, that shows how we model the errors, which we aim to estimate and use for corrections. Second part is the measurement model, that couples the sensory data coming at different rates.

The overall scheme of our proposed approach is shown in Figure 2. Raw sensor data are preprocessed and used as measurements in the error state EKF (the *FUSION* block). There is no measurement rejection implemented; based on the assumption that fusion of several sensor modalities should deal with anomalous data inherently—for details see Section 5 and Section 6—this however will be subjected to a future work. As apparent from the Figure 2, measurement rates significantly differ among the sensor modalities—main difference is especially between the IMU at 90Hz and the ICP output at 0.3Hz. Having the update rate of the EKF at 90Hz the experiments have proven that this issue is crucial and has to be resolved as part of the filter design to ensure reliable output from the fusion process (see Section 5.3.3). In our case, this problem concerns mainly the ICP-based localization that provides measurements at very low rate of 0.3Hz—too low to capture the motion dynamics as the IMU does (i.e. the motion dynamics spectrum gets sub-sampled). During these 3 seconds, real-world disturbances (which are often non-Gaussian and difficult to model and predict, e.g. tracks slippage) accumulate. This was the motivation to investigate various ways of fusing measurements at significantly different rates. Three proposed approaches that incorporate the ICP measurements are described in the Section 4.3.3.

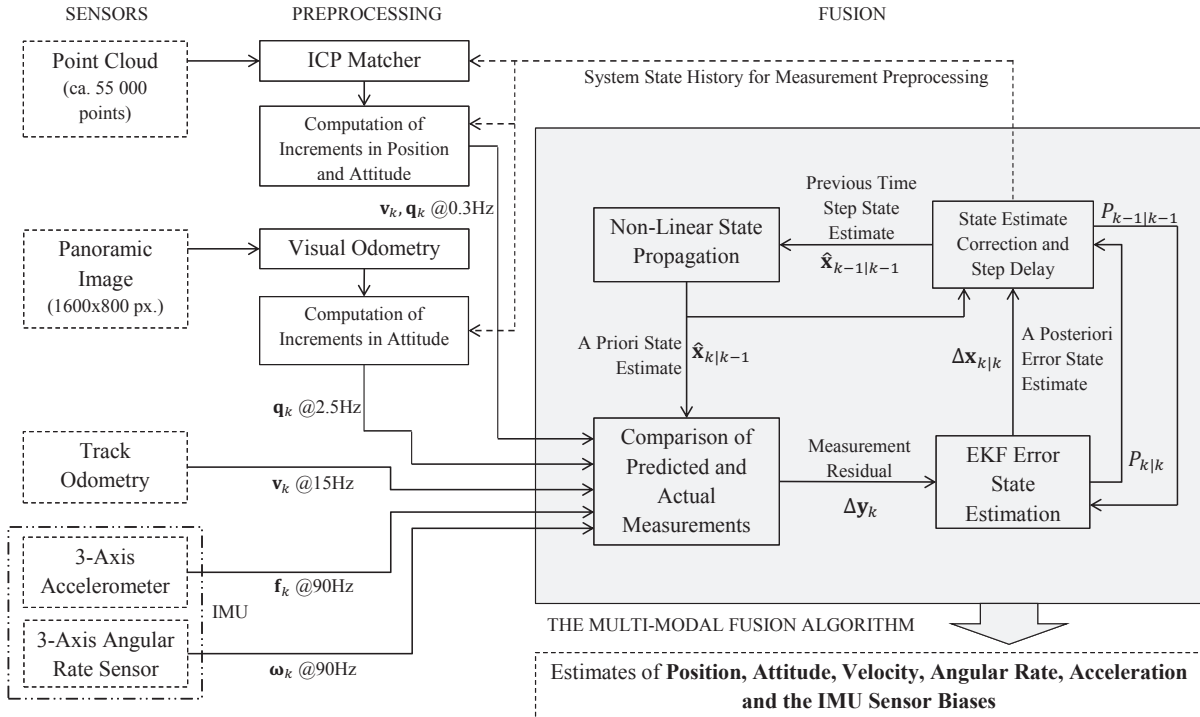


Figure 2: The scheme of the proposed multi-modal data fusion system ($\boldsymbol{\omega}$ is angular velocity, \mathbf{f} is specific force (Savage, 1998), \mathbf{v} is velocity and \mathbf{q} is quaternion representing attitude).

4.1 Process error model

For the purpose of localization, we model our robot as a rigid body with constant angular rate and constant rate of change of velocity ($\dot{\omega} = 0, \dot{v} = \text{const.}$). Presence of constant gravitational acceleration is expected and incorporated into system model; no dissipative forces are considered.

We define four coordinate frames: $R(\text{obot})$ frame coincides with center of the robot, $I(MU)$ frame represents the Inertial Measurement Unit coordinate frame as defined by the manufacturer, $O(\text{ometry})$ frame represents the tracked gear-frame, and $N(\text{avigation})$ frame represents the world frame. In all these frames, the North-West-Up axes convention is followed, with the x -axis pointing forwards (or to the North in the N-frame), the y -axis pointing to the left (or to the West), and the z -axis pointing upwards. Rotations about each axis follow the *right-hand rule*. The fundamental part of the system design are the differential equations describing development of the states in time. The state space with the corresponding errors is defined as:

$$\mathbf{x} = \begin{bmatrix} \mathbf{p}_N \\ \mathbf{q}_N^R \\ \mathbf{v}_R \\ \omega_R \\ \mathbf{f}_R \\ \mathbf{b}_{\omega,I} \\ \mathbf{b}_{f,I} \end{bmatrix} \quad \Delta \mathbf{x} = \begin{bmatrix} \Delta \mathbf{p}_N \\ \delta \theta \\ \Delta \mathbf{v}_R \\ \Delta \omega_R \\ \Delta \mathbf{f}_R \\ \Delta \mathbf{b}_{\omega,I} \\ \Delta \mathbf{b}_{f,I} \end{bmatrix} \quad (1)$$

where \mathbf{p}_N is position of the robot in the N-frame, \mathbf{q}_N^R is unit quaternion representing its attitude, \mathbf{v}_R is velocity expressed in the R-frame, ω_R is angular rate, \mathbf{f}_R is specific force (Savage, 1998), $\mathbf{b}_{\omega,I}$ and $\mathbf{b}_{f,I}$ are accelerometer and angular rate sensor IMU-specific biases expressed in the I-frame.

The error state $\Delta \mathbf{x}$ is defined—following the idea of (Weiss, 2012, eq. 3.25)—as difference between the system state and its estimate $\Delta \mathbf{x} = \mathbf{x} - \hat{\mathbf{x}}$ except for attitude, where rotation error vector $\delta \theta$ is the vector part of the error quaternion $\delta \mathbf{q} = \mathbf{q} \otimes \hat{\mathbf{q}}^{-1}$ multiplied by 2; \otimes represents quaternion multiplication as defined in (Breckenridge, 1999).

The states and the error states of the robot, modeled as a rigid body movement, propagate in time according to the following equations:

$$\dot{\mathbf{p}}_N = C_{(\mathbf{q}_N^R)}^T \mathbf{v}_R \quad \Delta \dot{\mathbf{p}}_N \approx C_{(\mathbf{q}_N^R)}^T \Delta \mathbf{v}_R - C_{(\mathbf{q}_N^R)}^T \delta \theta \quad (2)$$

$$\dot{\mathbf{q}}_N^R = \frac{1}{2} \Omega(\omega_R) \mathbf{q}_N^R \quad \delta \dot{\theta} \approx -[\hat{\omega}_R] \delta \theta + \Delta \omega_R + \mathbf{n}_\theta \quad (3)$$

$$\dot{\mathbf{v}}_R = \mathbf{f}_R - C_{(\mathbf{q}_N^R)} \mathbf{g}_N + [\mathbf{v}_R] \omega_R \quad \Delta \dot{\mathbf{v}}_R \approx \Delta \mathbf{f}_R - [C_{(\mathbf{q}_N^R)} \mathbf{g}_N] \delta \theta + [\hat{\mathbf{v}}_R] \Delta \omega_R - [\hat{\omega}_R] \Delta \mathbf{v}_R + \mathbf{n}_v \quad (4)$$

$$\begin{aligned} \dot{\omega}_R &= 0 & \dot{\mathbf{f}}_R &= 0 & \dot{\mathbf{b}}_{\omega,I} &= 0 & \dot{\mathbf{b}}_{f,I} &= 0 \\ \Delta \dot{\omega}_R &= \mathbf{n}_\omega & \Delta \dot{\mathbf{f}}_R &= \mathbf{n}_f & \Delta \dot{\mathbf{b}}_{\omega,I} &= \mathbf{n}_{b,\omega} & \Delta \dot{\mathbf{b}}_{f,I} &= \mathbf{n}_{b,f} \end{aligned} \quad (5)$$

where derivation of the left part of (3) can be found in (Trawny and Roumeliotis, 2005, eq. 110) and the left part of (4) is based on (Nemra and Aouf, 2010, eq. 5); the difference from the original is caused by different ways of expressing attitude. The right parts of (2-4) can be derived by neglecting higher-order error terms and by approximation of the error in attitude by the rotation error vector $\delta \theta$ following (Weiss, 2012, eq. 3.44). We define $\mathbf{g}_N = [0, 0, g]^T$, $\mathbf{n}_{(\cdot)}$ are the system noise terms and $\Omega(\omega_R)$ in (3) is a matrix representing quaternion and vector product operation (Trawny and Roumeliotis, 2005, eq. 108). It is constructed as

$$\Omega(\omega) = \begin{bmatrix} 0 & \omega_3 & -\omega_2 & \omega_1 \\ -\omega_3 & 0 & \omega_1 & \omega_2 \\ \omega_2 & -\omega_1 & 0 & \omega_3 \\ -\omega_1 & -\omega_2 & -\omega_3 & 0 \end{bmatrix} \quad (6)$$

In (5), time derivations of angular rates and specific forces are equal to zero—usually, they are considered rather as input than state. However, we included them into the state vector to be updated by the EKF. The error model equations can be expressed in compact matrix form:

$$\Delta \dot{\mathbf{x}} = F_c \Delta \mathbf{x} + G_c \mathbf{n} \quad (7)$$

where F_c is continuous time state transition matrix, G_c is noise coupling matrix and \mathbf{n} is noise vector composed of all the $\mathbf{n}_{(\cdot)}$ terms; the F_c matrix is as follows:

$$F_c = \begin{bmatrix} \emptyset_3 & -C_{(\hat{q}_N^R)}^T & C_{(\hat{q}_N^R)}^T & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & -[\hat{\omega}_R] & \emptyset_3 & I_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & -[C_{(\hat{q}_N^R)} g_N] & -[\hat{\omega}_R] & [\hat{v}_R] & I_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \end{bmatrix} \quad (8)$$

and the $G_c \mathbf{n}$ term is

$$G_c \mathbf{n} = \begin{bmatrix} \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ I_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & I_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & I_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & I_3 & \emptyset_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & I_3 & \emptyset_3 \\ \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & \emptyset_3 & I_3 \end{bmatrix} \begin{bmatrix} \mathbf{n}_\theta \\ \mathbf{n}_v \\ \mathbf{n}_\omega \\ \mathbf{n}_f \\ \mathbf{n}_{b,\omega} \\ \mathbf{n}_{b,f} \end{bmatrix} \quad (9)$$

The noise coupling matrix describes, how particular noise terms affect the system state. Each $\mathbf{n}_{(\cdot)}$ term is a random variable with Normal probability distribution. Properties of these random variables are described by their covariances in the system noise matrix Q_c . Since they are assumed independent, the matrix Q_c is diagonal $Q_c = \text{diag}(\sigma_{\theta_x}^2, \sigma_{\theta_y}^2, \sigma_{\theta_z}^2, \sigma_{v_x}^2, \sigma_{v_y}^2, \dots)$, where σ is standard deviation.

In order to implement the proposed model, we have to transform the continuous time equations to discrete time domain. We use the Van Loan discretization method (Loan, 1978) instead of explicitly expressing values of the discretized matrices. We substitute into matrix M defined by Van Loan

$$M = \begin{bmatrix} -F_c & G Q_c G^T \\ \emptyset & F_c^T \end{bmatrix} \Delta t \quad (10)$$

and evaluate the matrix exponential

$$e^M = \begin{bmatrix} \cdot & F_d^{-1} Q_d \\ \emptyset & F_d^T \end{bmatrix} \quad (11)$$

The result of the matrix exponential contains the discretized system matrix F_d in the bottom-right part and the discretized system noise matrix Q_d left multiplied by the inversion of F_d in the top-right part. The discretized system matrix F_d can be easily extracted; Q_d can be obtained by left multiplying the upper right part of e^M by F_d .

4.2 State prediction and update using EKF

The Extended Kalman filter (Smith et al., 1962; McElhoe, 1966), is a modification of the Kalman filter (Kalman, 1960), i.e. optimal observer minimizing variances of the observed states. Since the error state EKF is used in our approach, the state of the system is expressed as sum of current best estimate ($\hat{\mathbf{x}}$) and some small error ($\Delta \mathbf{x}$). The only difference compared to a standard EKF is that the linearised system matrices F and Q describe only the error state and the error state covariance propagation in time, rather than the whole state and state covariance propagation in time. This is mainly beneficial from the computational point of view since

it simplifies linearisation of the system equations. Flowchart describing the error state EKF computation is shown in Figure 3 and can be decomposed into a series of steps that describe the actual implementation. As new measurements arrive, state estimate ($\hat{\mathbf{x}}$) and its error covariance matrix (P) are available from the previous time-step (or as initialized during first iteration). This state estimate $\hat{\mathbf{x}}$ is propagated in time using the nonlinear system equations. The continuous time F_c and G_c matrices are evaluated based on the current value of $\hat{\mathbf{x}}$. Van Loan discretization method is used to obtain discrete forms of F_d and Q_d . Then the error state covariance matrix P is propagated in time. Expected measurements are compared to the incoming ones and their difference is expressed in form of measurement residual $\Delta\mathbf{y}$. Innovation matrix H , expressing the measurement residual as a linear combination of the error state components, is evaluated. Using the a priori estimate of P , H and the variance of the sensors signals expressed as R , the Kalman gain matrix K is computed. The error state $\Delta\mathbf{x}$ is updated using the Kalman gain and the measurement residual; the a posteriori estimate of the error state covariance matrix P , is evaluated as well. Finally, the a priori state estimate $\hat{\mathbf{x}}$ is corrected using the estimated error $\Delta\mathbf{x}$.

Although this EKF cycle can be repeated each time measurements arrive, yet, for performance reasons, we have chosen to group the incoming measurements to the highest frequency measurement, i.e. the IMU data. Hence, each time any non-IMU measurement arrives, it is slightly delayed until the next IMU measurement is available. The maximum possible sampling error caused by this grouping approach is $1/(2 \cdot 90)$ s and thus it can be neglected compared to the significantly longer sampling periods of the non-IMU data sources. The update rate of the EKF is then equal to the IMU sampling rate, i.e. 90 Hz.

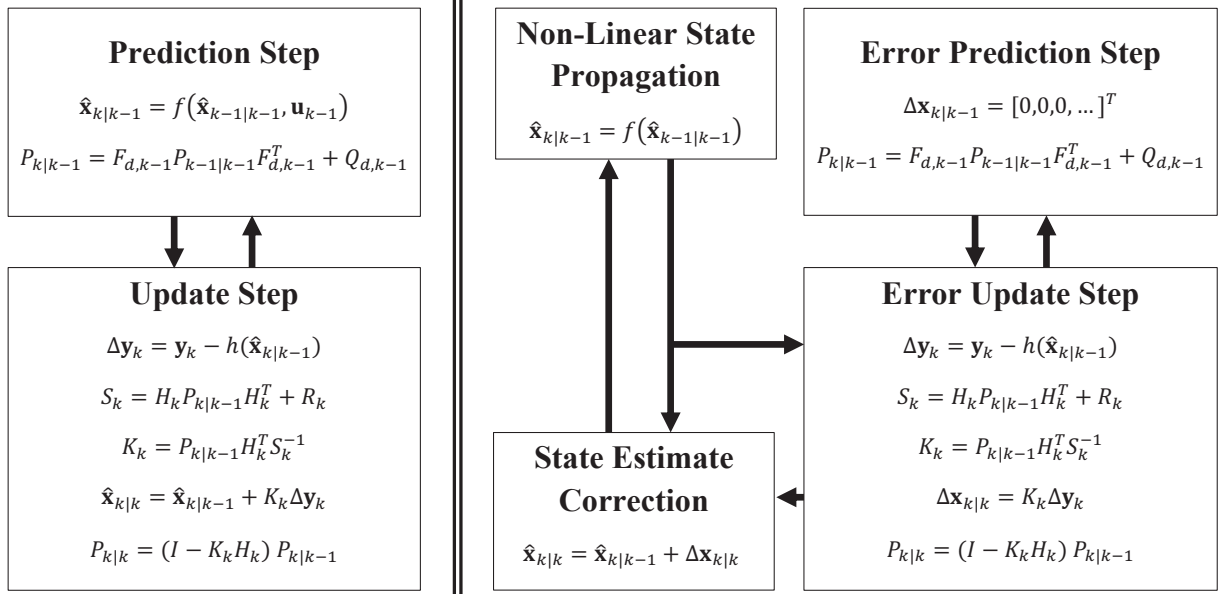


Figure 3: Standard EKF (left) computation flowchart compared to the error state EKF computation flowchart (right): in the error state EKF prediction step, the a priori state is estimated using the non-linear system equation $f()$ and the covariances are estimated using F_d (linearized matrix form of the error state propagation equations). In the update step, measurement residual $\Delta\mathbf{y}$ is obtained by comparing the incoming measurement \mathbf{y} with its predicted counterpart. The residual covariance S and the Kalman gain K are evaluated and used to update the state and covariance matrix to obtain the a posteriori estimates. Note that in the case of the error state EKF, Q_d and H_k couple system noise and measurements with the error state $\Delta\mathbf{x}$ rather than $\hat{\mathbf{x}}$.

4.3 Measurement error model

In general, the measurement vector \mathbf{y} can be described as sum of measurement function $h(\mathbf{x})$ of the state \mathbf{x} and of some random noise \mathbf{m} due to properties of the individual sensors:

$$\mathbf{y} = h(\mathbf{x}) + \mathbf{m} \quad (12)$$

Using the function h , we can predict the measured value based on current knowledge about the system state:

$$\hat{\mathbf{y}} = h(\hat{\mathbf{x}}) \quad (13)$$

There is a difference $\Delta\mathbf{y} = \hat{\mathbf{y}} - \mathbf{y}$ caused by the modeling imperfections in the state estimate as well as by the sensor errors. This difference can be expressed in terms of the error state $\Delta\mathbf{x}$:

$$\begin{aligned} \Delta\mathbf{y} &= \mathbf{y} - \hat{\mathbf{y}} = h(\mathbf{x}) - h(\hat{\mathbf{x}}) + \mathbf{m} \\ &= h(\hat{\mathbf{x}} + \Delta\mathbf{x}) - h(\hat{\mathbf{x}}) + \mathbf{m} \end{aligned} \quad (14)$$

If function h is linear, (14) becomes

$$\Delta\mathbf{y} = h(\Delta\mathbf{x}) + \mathbf{m} \quad (15)$$

Although the condition of linearity is not always met we still can approximate the behavior of h in some close proximity to the current state $\hat{\mathbf{x}}$ by a similar function h' , which is linear in elements of $\hat{\mathbf{x}}$ such that

$$h(\hat{\mathbf{x}} + \Delta\mathbf{x}) - h(\hat{\mathbf{x}}) \approx h'(\Delta\mathbf{x})|_{\hat{\mathbf{x}}} = H_{\hat{\mathbf{x}}}\Delta\mathbf{x} \quad (16)$$

where $H_{\hat{\mathbf{x}}}$ is the innovation matrix projecting observed differences in measurements onto the error states.

4.3.1 IMU measurement model

The inertial measurement unit (IMU) is capable of measuring specific force (Savage, 1998) in all three dimensions as well as angular rates. The specific force measurement is a sum of acceleration and gravitational force, but it also contains biases—constant or slowly changing value independent of the actual acting forces—and sensor noise, which is expected to have zero mean normal probability. All the values are measured in the I-frame.

$$\mathbf{y}_{f,I} = \mathbf{f}_I + \mathbf{b}_{f,I} + \mathbf{m}_{f,I} \quad (17)$$

where $\mathbf{y}_{f,I}$ is the measurement, \mathbf{f}_I is the true specific force, $\mathbf{b}_{f,I}$ is sensor bias and $\mathbf{m}_{f,I}$ is sensor noise.

Since the interesting value $\mathbf{y}_{f,I}$ is expressed in the I-frame, we define a constant rotation matrix C_R^I of R-frame to I-frame. Translation between the I- and R-frames does not affect the measured values directly; thus, it is not considered. Since the IMU is placed close to the R-frame origin, we neglect centrifugal force induced by rotation of R-frame and conditioned by non-zero translation between R- and I-frames. Using this rotation matrix, we express the measurement as:

$$\mathbf{y}_{f,I} = C_R^I \mathbf{f}_R + \mathbf{b}_{f,I} + \mathbf{m}_{f,I} \quad (18)$$

where both \mathbf{f}_R and $\mathbf{b}_{f,I}$ are elements of the system state. If we compare the measured value and the expected measurement, we can express the h function, which is—in this case—equal to the h' :

$$\begin{aligned} \mathbf{y}_{f,I} - \hat{\mathbf{y}}_{f,I} &= \Delta\mathbf{y}_{f,I} = C_R^I \mathbf{f}_R + \mathbf{b}_{f,I} - C_R^I \hat{\mathbf{f}}_R - \hat{\mathbf{b}}_{f,I} + \mathbf{m}_{f,I} \\ &= C_R^I \Delta\mathbf{f}_R + \Delta\mathbf{b}_{f,I} + \mathbf{m}_{f,I} \end{aligned} \quad (19)$$

and hence can be expressed in $H_{\hat{\mathbf{x}}}\Delta\mathbf{x}$ form as

$$\Delta\mathbf{y}_{f,I} = [\emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad C_R^I \quad \emptyset_3 \quad I] \Delta\mathbf{x} + \mathbf{m}_{f,I} \quad (20)$$

where the error state $\Delta\mathbf{x}$ was defined in (1).

The angular rate measurement is treated identically; the output of the sensor is

$$\mathbf{y}_{\omega,I} = \omega_I + \mathbf{b}_{\omega,I} + \mathbf{m}_{\omega,I} \quad (21)$$

where ω_I is angular rate, $\mathbf{b}_{\omega,I}$ is sensor bias and $\mathbf{m}_{\omega,I}$ is sensor noise.

Similarly, the measurement residual is obtained:

$$\mathbf{y}_{\omega,I} - \hat{\mathbf{y}}_{\omega,I} = \Delta\mathbf{y}_{\omega,I} = C_R^I \Delta\omega_R + \Delta\mathbf{b}_{\omega,I} + \mathbf{m}_{\omega,I} \quad (22)$$

which can be expressed in the matrix form

$$\Delta\mathbf{y}_{\omega,I} = [\emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad C_R^I \quad \emptyset_3 \quad \emptyset_3 \quad I] \Delta\mathbf{x} + \mathbf{m}_{\omega,I} \quad (23)$$

4.3.2 Odometry measurement model

Our platform is equipped with caterpillar tracks and therefore, steering is realized by setting different velocities to each of the tracks (*skid-steering*). The velocities are measured by incremental optical angle sensors at 15 Hz. Originally, we implemented a complex model introduced in (Endo et al., 2007), which exploits angular rate measurements to model the slippage to further improve the odometry precision. However, with respect to our sensors, no improvement was observed. Moreover, since the slippage is inherently corrected via the proposed data fusion, we can neglect it in the odometry model, assuming only a very simple but sufficient model:

$$v_{O,x} = \frac{v_r + v_l}{2} \quad (24)$$

where $v_{O,x}$ is the forward velocity, v_l and v_r are track velocities measured by incremental optical sensors—the velocities in the lateral and vertical axes are set to zero. Since the robot position is obtained by integrating velocity expressed in R-frame, we define a rotation matrix C_R^O :

$$\mathbf{v}_O = C_R^O \mathbf{v}_R \quad (25)$$

which expresses the \mathbf{v}_R in the O-frame.

During experimental evaluation, we observed a minor misalignment between these two frames, which can be described as rotation about the lateral axis by approximately 1 degree. Although relatively small, this rotation caused the position estimate in the vertical axis to grow at constant rate while the robot was moving forward. To compensate for this effect, we handle the C_R^O as constant—its value was obtained by means of calibration. The measurement equation is then as follows:

$$\mathbf{y}_{v,O} = C_R^O \mathbf{v}_R + \mathbf{m}_{v,O} \quad (26)$$

where $\mathbf{y}_{v,O}$ is linear velocity measured by the track odometry, expressed in O-frame. Since this relation is linear, the measurement innovation is

$$\begin{aligned} \mathbf{y}_{v,O} - \hat{\mathbf{y}}_{v,O} &= \Delta\mathbf{y}_{v,O} = \\ &= C_R^O \mathbf{v}_R - C_R^O \hat{\mathbf{v}}_R + \mathbf{m}_{v,O} \\ &= C_R^O \Delta\mathbf{v}_R + \mathbf{m}_{v,O} \end{aligned} \quad (27)$$

and expressed in the matrix form

$$\Delta\mathbf{y}_{v,O} = [\emptyset_3 \quad \emptyset_3 \quad C_R^O \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3] \Delta\mathbf{x} + \mathbf{m}_{v,O} \quad (28)$$

4.3.3 ICP-based localization measurement model

The ICP algorithm is used to estimate translation and rotation between each new incoming laser scan of the robot surroundings and a metric map created from the previously registered laser scans. In course of our work, three approaches processing the output of the ICP were proposed and tested. The first approach treats the ICP-based localization as movement in the R-frame in between two consecutive laser scans in form of a position increment (the *incremental position approach*). The idea of measurements expressed in a form of some $\Delta \mathbf{p}$ can be, for example, found in (Ma et al., 2012). In our case, the increment is obtained as:

$$\Delta \mathbf{p}_{R,ICP,i} = C_{(\mathbf{q}_{N,ICP,i-1}^R)}(\mathbf{p}_{N,ICP,i} - \mathbf{p}_{N,ICP,i-1}) \quad (29)$$

where both the position $\mathbf{p}_{N,ICP}$ and attitude $\mathbf{q}_{N,ICP}^R$ are outputs of the ICP algorithm. The increment $\Delta \mathbf{p}_{R,ICP,i}$ is added to the position estimated by the whole fusion algorithm at time-step $i - 1$ to be used as a direct measurement of position. The same idea is applied in the case of attitude (an increment in attitude is extracted by means of quaternion algebra). The purpose is to overcome the ICP world frame drift. However, it is impossible to correctly discretize the system equations respecting the laser scan sampling frequency ($\frac{1}{3}$ Hz). Also, the assumption of measurements being independent is violated by utilizing a previously estimated state to create a new measurement. Thus, corrections that propagate to the system state from this measurement tend to be inaccurate.

The second approach treats the ICP output as velocity in the R-frame (the *velocity approach*). We consider it a state-of-the-art practice utilized, for example, by (Almeida and Santos, 2013). The velocity is expressed in the N-frame first:

$$\mathbf{v}_{N,ICP} = \frac{\mathbf{p}_{N,ICP,i} - \mathbf{p}_{N,ICP,i-1}}{t(i) - t(i-1)} \quad (30)$$

where $t(i)$ is time corresponding to a time-step i . To express the velocity in the R-frame:

$$\mathbf{v}_{R,ICP}(t) = C_{(\mathbf{q}_{R',ICP}^R(t) \otimes \mathbf{q}_{N,ICP,i-1}^R)} \mathbf{v}_{N,ICP} \quad (31)$$

it is necessary to interpolate the attitude between $\mathbf{q}_{N,ICP,i-1}^R$ and $\mathbf{q}_{N,ICP,i}^R$ in order to obtain the increment $\mathbf{q}_{R',ICP}^R(t)$. Angular velocity is assumed to be constant between the two laser scans. The velocity $\mathbf{v}_{R,ICP}$ and the constant angular velocity obtained from the interpolation can be directly used as measurements which are independent of the estimated state and because of the interpolation, they can be generated with arbitrary frequency and thus, there is no problem with discretization (compared to the previous approach). However, this approach expects the robot to move in a line between the two ICP scans. This is a too strong assumption and also a major drawback of this approach that results in incorrect trajectory estimates.

Therefore, we propose the third approach, the *trajectory approach*, which overcomes the assumption of the *velocity approach* by (sub-optimal) use of the estimated states in order to approximate possible behavior of the system between each two consecutive ICP scans. This *trajectory approach* proved to be the best for pre-processing the output of the ICP algorithm; for details see Section 5.4.5.

The *trajectory approach* assumes that the first estimate of the trajectory (without the ICP measurement) is locally very similar to the true trajectory (up to the effects of drift). Thus, when a new ICP measurement arrives the trajectory estimated since the previous ICP measurement is stored to be used as the best guess around the previous ICP pose. The ICP poses at time-steps i and $i - 1$ are aligned with the N-frame so the ICP pose at time-step $i - 1$ coincides with the first pose of the stored trajectory. This way the ICP world frame drift is suppressed. Then, the stored trajectory is duplicated and aligned with the new ICP pose to serve as the best guess around the new ICP pose; see Figure 4. The resulting trajectory is obtained as weighted average of the original and the duplicated trajectories:

$$\hat{\mathbf{p}}_{N,weighted,k} = \hat{\mathbf{p}}_{N,k} w_k + \hat{\mathbf{p}}'_{N,k} w'_k \quad (32)$$

where $\hat{\mathbf{p}}_{N,k}$ are points of the original trajectory (black dotted line in Figure 4), $\hat{\mathbf{p}}'_{N,k}$ are points of the realigned duplicated trajectory (black dashed line in Figure 4) and w_k, w'_k are weights—linear functions of

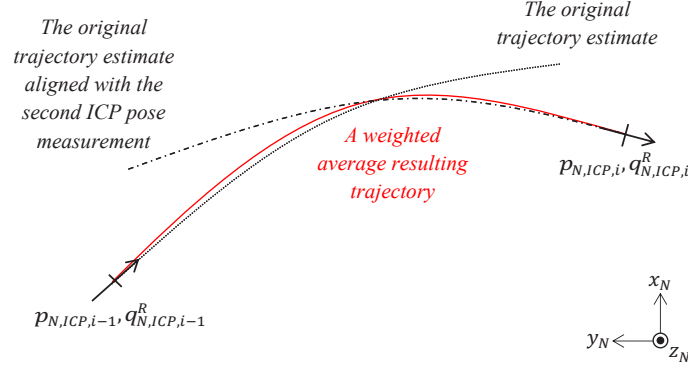


Figure 4: The principle of *trajectory approach*: when the new ICP measurement arrives (time-step i), trajectory estimate based on measurements other than ICP (black dotted line) is duplicated and aligned with the incoming ICP measurement (black dashed line) and weighted average (red solid line) of these two trajectories is computed.

time equal to 1 at time-step of associated ICP measurement and equal to 0 at time-step of the other ICP measurement. The resulting trajectory is used to generate the velocity measurements in the N-frame as follows:

$$\mathbf{v}_{N,weighted,k} = \frac{\mathbf{p}_{N,weighted,k} - \mathbf{p}_{N,weighted,k-1}}{t(k) - t(k-1)} \quad (33)$$

where $t(k)$ and $t(k-1)$ are time-steps of poses of the resulting weighted trajectory. The k denotes indexing of the fusion algorithm high-frequency samples. Velocities can be expressed in R-frame using the attitude estimates $\hat{\mathbf{q}}_{N,k}^R$:

$$\mathbf{v}_{R,weighted,k} = C_{(\hat{\mathbf{q}}_{N,k}^R)} \mathbf{v}_{N,weighted,k} \quad (34)$$

and can be used directly as measurement, whose projection onto the error state vector yields:

$$\Delta \mathbf{y}_{v,weighted} = [\emptyset_3 \quad \emptyset_3 \quad I_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3] \Delta \mathbf{x} + \mathbf{m}_{v,weighted} \quad (35)$$

The velocity expressed in R-frame can be used this way as measurement, but its values for the time period between two consecutive ICP outputs are known only *after* the second ICP measurement arrives. Thus it is necessary to recompute state estimates for this whole time period (typically in length of 300 IMU samples), including the new velocity measurements.

To process the attitude information provided as the ICP output, we use a simple incremental approach such that the drift of the ICP world frame with respect to the N-frame is suppressed. To achieve this, we extract only the increment in attitude between two consecutive ICP poses:

$$\mathbf{q}_{N,ICP,i}^R = \mathbf{q}_{R',ICP}^R \otimes \mathbf{q}_{N,ICP,i-1}^{R'} \quad (36)$$

$$\mathbf{q}_{R',ICP}^R = \mathbf{q}_{N,ICP,i}^R \otimes \left(\mathbf{q}_{N,ICP,i-1}^{R'} \right)^{-1} \quad (37)$$

where $\mathbf{q}_{R',ICP}^R$ is rotation that occurred between two consecutive ICP measurements $\mathbf{q}_{N,ICP,i-1}^{R'}$ and $\mathbf{q}_{N,ICP,i}^R$. We apply this rotation to the attitude state estimated at time-step $k' \equiv i-1$:

$$\mathbf{y}_{q,ICP} = \mathbf{q}_{R',ICP}^R \otimes \hat{\mathbf{q}}_{N,k'}^R \quad (38)$$

To express the measurement residual, we define the following error quaternion:

$$\delta \mathbf{q}_{ICP,i} = \hat{\mathbf{q}}_{N,k}^R \otimes (\mathbf{y}_{q,ICP})^{-1} \quad (39)$$

where $\hat{\mathbf{q}}_{N,k}^R$ is the attitude estimated at time-step $k \equiv i$. We express this residual rotation by means of rotation vector $\delta\theta_{ICP,i}$

$$\delta\theta_{ICP,i} = 2\delta\vec{\mathbf{q}}_{ICP,i} \quad (40)$$

which can be projected onto the error state as

$$\Delta\mathbf{y}_{\delta\theta,ICP} = [\emptyset_3 \quad I_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3] \Delta\mathbf{x} + \mathbf{m}_{\delta\theta,ICP} \quad (41)$$

Although the ICP is very accurate in measuring translation between consecutive measurements, the attitude measurement is not as precise. Noise introduced in the pitch angle can cause wrong velocity estimates expressed in R-frame, resulting in problem described as *climbing robot*—the system tends to slowly drift in the vertical axis. Since the output of the *trajectory approach* is velocity $\mathbf{v}_{R,weighted,i}$, applying a constraint assuming only planar motion in the R frame is fully justified, easy to implement and resolves this issue.

4.3.4 Visual odometry measurement model

As explained in Section 3.5, the visual odometry (VO) is an algorithm for estimating translation and rotation of a camera body based on images recorded by the camera. The current implementation of the data fusion utilizes only the rotation part of the motion estimated by the VO, since it is not affected by the scale. The set of 3D landmarks maintained by the VO is not in any way processed by the fusion algorithm—it is used by the VO to improve its attitude estimates internally. Similarly, the bundle adjustment ensures more consistent measurements, yet still, it does not enter the data fusion models.³ The way we incorporate the VO measurements is equivalent to the ICP *trajectory approach*, however, reduced only to the incremental processing of the attitude measurements. This way, the whole VO processing block can easily be replaced by an alternative (for example by stereo vision based VO), provided the output—the estimated rotation—is available in the same way. The motivation is to have the VO measurement model independent on the VO internal implementation details. The implementation of the VO attitude aiding is identical to the ICP attitude aiding; the attitude increment is extracted and used to construct a new measurement $\mathbf{y}_{q,VO}$:

$$\mathbf{q}_{N,VO,i}^R = \mathbf{q}_{R',VO}^R \otimes \mathbf{q}_{N,VO,i-1}^{R'} \quad (42)$$

$$\mathbf{q}_{R',VO}^R = \mathbf{q}_{N,VO,i}^R \otimes \left(\mathbf{q}_{N,VO,i-1}^{R'} \right)^{-1} \quad (43)$$

where $\mathbf{q}_{R',VO}^R$ is rotation that occurred between two consecutive VO measurements $\mathbf{q}_{N,VO,i-1}^{R'}$ and $\mathbf{q}_{N,VO,i}^R$. We apply this rotation to the attitude state estimated at time-step $k' \equiv i - 1$:

$$\mathbf{y}_{q,VO} = \mathbf{q}_{R',VO}^R \otimes \hat{\mathbf{q}}_{N,k'}^R \quad (44)$$

Then, the measurement residual is expressed as error quaternion:

$$\delta\mathbf{q}_{VO,i} = \hat{\mathbf{q}}_{N,k}^R \otimes (\mathbf{y}_{q,VO})^{-1} \quad (45)$$

where $\hat{\mathbf{q}}_{N,k}^R$ is the attitude estimated at time-step $k \equiv i$. We express this residual rotation by means of rotation vector $\delta\theta_{VO,i}$

$$\delta\theta_{VO,i} = 2\delta\vec{\mathbf{q}}_{VO,i} \quad (46)$$

which can be projected onto the error state as

$$\Delta\mathbf{y}_{\delta\theta,VO} = [\emptyset_3 \quad I_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3 \quad \emptyset_3] \Delta\mathbf{x} + \mathbf{m}_{\delta\theta,VO} \quad (47)$$

where $\mathbf{m}_{\delta\theta,VO}$ is the VO attitude measurement noise.

³The same idea applies for the ICP-based localization: although it builds an internal map, this map is independent from our localization estimates. This would not be the case in a SLAM approach with integrated loop closures.

5 Experimental evaluation

Our evaluation procedure involves several different tests. First, we describe our evaluation methodology in Section 5.1. It covers obtaining ground truth positioning measurements for both in- and outdoors. Then we present and discuss our field experiments with the global behavior of our state estimation (Section 5.2). We also show two examples of typical behavior of the filter in order to give more insight on its general characteristics (Section 5.3). We take advantage of them to explain the importance of the *trajectory approach*, compared to more standard measurement models. Finally, we analyze the behavior of the filter under failure case scenarios involving partial or full outage of each sensory modality (Section 5.4).

5.1 Evaluation metrics

In order to validate the results of our fusion system, we need accurate measurements of part of our system states to confront with the proposed filter. For indoor measurements, we use a Vicon motion capture system with nine cameras covering more than 20 m² and giving a few millimeter accuracy at 100 Hz.

For external tracking, we use a theodolite from Leica Geosystems, namely the Total Station TS15; see Figure 5 (left). It can track a reflective prism to measure its position continuously at an average frequency of 7.5 Hz. The position precision of the theodolite is 3 mm in continuous mode. However, this system cannot measure the orientation of the robot. Moreover, the position measured is that of the prism and not directly of the robot, therefore we calibrated the position of the prism with respect to the robot body using the theodolite and precise blueprints. However, the position of the robot cannot be recovered from the position of the prism without the information about orientation. That explains why, in the validations below, we do not compare the position of the robot but the position of the prism from the theodolite and reconstructed from the states of our filter. With these ground-truth measurements, we use different metrics for evaluation. First, we simply plot the error as a function of time. More precisely, we consider *position error*, *velocity error*, and *attitude error* and we compute it by taking the norm of the difference between the prediction made by our filter and the reference value.

Since this metric shows how the errors evolve over time, a more condensed measure is needed to summarize and compare the results of different versions of the filter. Therefore, we use the *final position error* expressed as a percentage of the total trajectory length:

$$e_{rel} = \frac{\|\mathbf{p}_l - \mathbf{p}_{ref,l}\|}{\text{distance travelled}} \quad (48)$$

where l is the index of the last position sample \mathbf{p}_l with the corresponding reference position $\mathbf{p}_{ref,l}$.

While this metric is convenient and widely used in the literature, it is however representative only of the end point error regardless of the intermediary results. This can be misleading for long trajectories in confined environment as the end-point might be close to the ground truth by chance. This is why we introduce, as a complement, the *average position error*:

$$e_{avg}(l) = \frac{\sum_{i=1}^l \|\mathbf{p}_i - \mathbf{p}_{ref,i}\|}{l} \quad (49)$$

where $1 \leq l \leq \text{total number of samples}$. To improve legibility of this metric in plots, we express the e_{avg} as a function of time

$$e'_{avg}(t) = e_{avg}(l(t)) \quad (50)$$

where $l(t)$ simply maps time t to the corresponding sample l .

5.2 Performance overview of the proposed data fusion

With these metrics, we can actually evaluate the performance of our system in a quantitative way. We divided the tests into indoor and outdoor experiments.

5.2.1 Indoor performance

For the indoor tests, we replicated semi-structured environment found in USAR environments, including ramps, boxes, a catwalk, a small passage, etc. Figure 5 (right) shows a picture of part of the environment. Due to limitations of our motion capture set-up, this testing environment is not as large as typical indoor USAR environments. Nevertheless, it features most of the complex characteristics that make state estimation challenging in such an environment.



Figure 5: The experimental setup with the Leica reference theodolite for obtaining ground truth trajectory (left). Part of the 3D semi-structured environment for indoor test with motion capture ground truth (right).

For this evaluation, we recorded approximately 2.4 km of indoor data with ground truth; 28 runs represent standard conditions (765 m in total), 36 runs represent failure cases of different sensory modalities induced artificially (1613 m in total). Table 2 presents the results of each combination of sensory modalities for the 28 standard conditions runs; the failure scenarios are analyzed in Section 5.4 separately.

The sensory modalities combinations can be divided into two groups by including or excluding the ICP modality; these two groups differ by the magnitude of the final position error. From this fact, we conclude that the main source of error is slippage of the caterpillar tracks—the VO modality in our fusion system corrects only the attitude of the robot. Also, the results confirmed sensitivity to erroneous attitude measurements originating from the sensory modalities. In this instance, VO has slightly worsen the median of the final position error—the indoor experiments are not long enough to make the difference between drift rates of the bare IMU+OD combination and possible VO errors that originate from incorrect pairing of image features. Nevertheless, the results are not significantly different.⁴ A significant improvement is brought with the ICP modality, which compensates the tracks slippage and reduces the resulting median of the final position errors by 50% (approximately). As expected during the filter design, fusing all sensory modalities yields the best result (not significantly different that without VO), with a median of 1.2% final position error; the occasional VO attitude measurement errors are diminished by the ICP modality attitude measurement (and vice versa).

5.2.2 Outdoor performance

We ran outdoor tests in various different environments; namely a street canyon and a urban park with trees and stairs in Zurich. Figure 6 shows pictures of the environments.

⁴All statistical significance results are assessed using the Wilcoxon signed-rank test with $p < 0.05$ testing whether the median

Exp.	Distance traveled [m]	Exp. duration [s]	Final position error in % of the distance travelled			
			OD, IMU	OD, IMU, VO	OD, IMU, ICP	OD, IMU, ICP, VO
1	47.42	254	2.17	2.30	1.71	0.79
2	36.52	186	1.99	2.21	0.36	0.14
3	48.74	244	3.15	2.63	0.50	0.18
4	29.40	237	2.22	2.06	0.42	0.45
5	82.10	585	2.51	2.24	0.90	0.71
6	74.64	452	2.05	3.64	0.98	1.24
7	74.65	387	1.70	1.72	2.28	0.58
8	30.57	194	1.98	3.42	1.59	2.29
9	26.58	287	2.67	2.23	1.90	1.19
10	26.57	236	1.53	3.94	0.77	2.11
11	26.96	208	1.25	1.20	0.95	0.66
12	29.13	211	1.27	1.29	0.88	0.87
13	26.35	180	1.37	1.25	0.94	0.77
14	40.23	240	6.58	6.70	0.88	0.99
15	21.01	167	5.26	5.27	0.61	0.57
16	19.04	209	5.94	5.95	0.55	0.60
17	10.95	405	3.44	2.89	2.15	2.05
18	8.65	238	2.87	2.77	1.36	1.38
19	9.36	284	4.14	3.91	1.83	1.85
20	9.02	282	2.90	3.36	2.73	2.65
21	10.82	308	3.79	3.23	1.43	1.41
22	9.45	237	5.36	5.45	2.66	2.68
23	12.75	204	2.65	2.84	2.66	1.79
24	7.81	179	1.58	1.83	2.82	3.06
25	10.85	165	3.85	4.14	3.25	2.17
26	10.83	163	2.36	1.84	0.62	0.68
27	12.79	237	15.42	14.95	2.48	2.53
28	12.07	239	28.42	27.07	2.89	2.98
Lower quartile Median Upper quartile			2.0 2.7 4.0	2.1 2.9 4.0	0.8 1.4 2.4	0.7 1.2 2.1

Table 2: Comparison of combinations of different modalities evaluated on indoor experiments performed under standard conditions with the Vicon system providing ground truth in position and attitude. Final position error expressed in percents of the total distances traveled was chosen as metric for each experiment; the total distance of the 28 experiments was 765 m, including traversing obstacles.

In those environments we recorded in total approximately 2 km, with ground truth available for 1.6 km, the rest were returns from the experimental areas. These 1.6 km are split into 10 runs and, as for the indoor experiments, Table 3) presents the results of each combination of sensory modalities for each run.

Contrary to the indoor experiments, combining all four modalities does not improve precision of localization compared to ICP, IMU and odometry fusion (the fusion of all is significantly worse than ICP, IMU and odometry only). Although some runs show improvement while combining all the sensory modalities (runs 7, 9 and 10) or are at least comparable with the best result 0.4|**0.6**|1.2 (runs 4, 5 and 6), there are several experiments, where VO failed due to the specificities of the environments. Such failures result in erroneous attitude estimates significantly exceeding expected VO measurement noise and compromising localization accuracy of the fusion algorithm. The reasons for failures are described in the Section 5.4 together with other failure cases. Since we did not artificially induce these VO failures as we did in the case of the indoor experiments, we do not exclude these runs from the performance evaluation in Table 3—we consider such environments standard for USAR. Moreover, we treat them as another proof of the fusion algorithm sensitivity to erroneous attitude measurement originating both from VO and ICP modalities and address

of correlated samples is different.



Figure 6: Pictures of the outdoor environments in Zurich. Left: street canyon, right: urban park.

Experiment	Distance traveled [m]	Exp. duration [s]	Final position error in % of the distance travelled			
			OD, IMU	OD, IMU, VO	OD, IMU, ICP	OD, IMU, ICP, VO
1: basement 1	120.62	825	2.08	26.61	1.83	17.84
2: basement 2	175.67	853	1.37	12.53	2.42	5.91
3: hallway straight	159.42	738	1.10	20.48	0.43	12.22
4: street 1	135.18	584	2.78	0.72	0.24	0.62
5: street 2	259.86	992	9.74	0.80	0.26	0.80
6: park big loop	145.31	918	2.65	2.66	1.03	1.76
7: park small loop	88.20	601	1.94	1.60	1.25	0.97
8: park straight	99.29	560	1.20	20.18	0.62	11.50
9: 2 floors	238.28	1010	9.10	0.62	0.58	0.43
10: 2 floors opposite	203.23	1107	3.23	6.79	0.51	0.42
Lower quartile Median Upper quartile			1.4 2.4 3.2	0.8 4.7 20.2	0.4 0.6 1.2	0.6 1.4 11.5

Table 3: Comparison of combinations of different modalities evaluated on outdoor experiments performed under standard conditions with the Leica system providing ground truth in position.

them in the conclusions and future work.

5.3 In-depth analysis of the examples of performance

In order to have more insight on the characteristics of the filter, we selected some trajectories and show more information than just the final position error metric.

5.3.1 Example of data fusion performance in indoor environment

In this example we address the caterpillar tracks slippage when traversing an obstacle (Figure 7). Since we are looking forward to USAR missions, such environment with conditions inducing high slippage can be expected, e.g. collapsed buildings full of debris and dust that impairs traction on smooth surfaces such as exposed concrete walls or floors, mass traffic accidents with oil spills, etc. The Vicon system was used to obtain precise position and orientation ground truth for computing the *average position error* development in time.

When traversing a slippery surface, any track odometry inevitably fails with the tracks moving with significantly diminishing traction. For this reason, trajectory and state estimates resulting from the IMU+OD fusion showed unacceptable error growth; see Figure 8. The robot was operated to attempt climbing up the

yellow slippery board (Figure 7), which deteriorated the traction to the point the robot was sliding back down with each attempt to steer. Because of the slippage, it failed to reach the top. Then, it was driven around the structure and up, to further slowly slip down the slope backwards, with the tracks moving forward to spoil traction. The effect of the slippage on the OD is apparent from the purple line in Figure 8. The corresponding average position error of the bare combination of IMU+OD starts to build up as soon as the robot enters the slippery slope. At 75 seconds, the IMU+OD has already an error of 0.5 m and finishes at 200 sec at an error of 4.4 m (outside Figure 8). Without exteroceptive modalities this problem is unsolvable and as expected, including these modalities significantly improves the localization accuracy; the final average position error is only 0.14 m for the IMU+OD+VO+ICP combination. The resulting state estimates for combination of all modalities are shown in Figure 9 and 10. Figure 9 depicts position estimates (the upper left quarter) with the reference values. The difference between the estimate and the reference is plotted in the bottom left quarter; similarly, the right half of the figure displays the velocity estimate. In the left part of the Figure 10, the attitude estimate expressed in Euler angles is shown with its error compared to the Vicon reference. The right part of this figure demonstrates estimation of the sensor biases, which are part of the system state. Note, that the biases in angular rates are initialized to values obtained as the mean of angular rate samples measured when the robot stays stationary before each experiment—short self-calibration. Concluded, adding the exteroceptive sensor modalities—as proposed in our filter design—compensated the effect caused by high slippage shown in this example as shown by the shape of trajectories and the average position error.

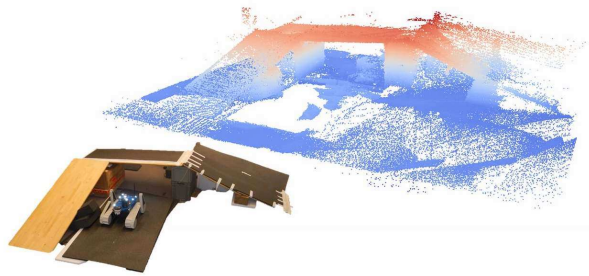


Figure 7: The 3D structure for testing of obstacle traversability shown a metric map created by ICP.

5.3.2 Example of data fusion performance in outdoor environment

This outdoor experiment took place on the Clausiusstrasse street (nearby ETH in Zurich) (Figure 11) and the purpose was testing the exteroceptive modalities (the ICP and the VO) in open urban space. In this standard setting, both the ICP and the VO are expected to perform reasonably well, though the ICP—compared to a closed room—is missing a significant amount of spatial information (laser range is limited approximately to 50 meters, no ceiling etc.). The Leica theodolite was used to obtain the ground truth position during this experiment (Figure 5).

The results are shown in Figures 12 and 13, demonstrating the improvement of performance when including more modalities up to the full setup. The basic dead-reckoning combination (IMU+OD) showed clearly drift in the yaw angle caused by accumulating error due to angular rate sensor noise integration (see the purple trajectory in the left part of Figure 12). By including the VO attitude measurements (resulting in IMU+OD+VO) the drift was compensated. Though the VO is not in fact completely drift-free, the performance is clearly better than the angular rate integration—it is rather the scale of the trajectory that matters. The IMU+OD+VO modality combination suffered from inaccurate track odometry velocity measurements (the green line in Figure 12), but this problem was resolved by including the ICP modality into the fusion scheme. The IMU+OD+ICP+VO combination proved to provide the best results; see the average position error plot in Figure 12 (right). The attitude estimates and estimates of the sensor biases

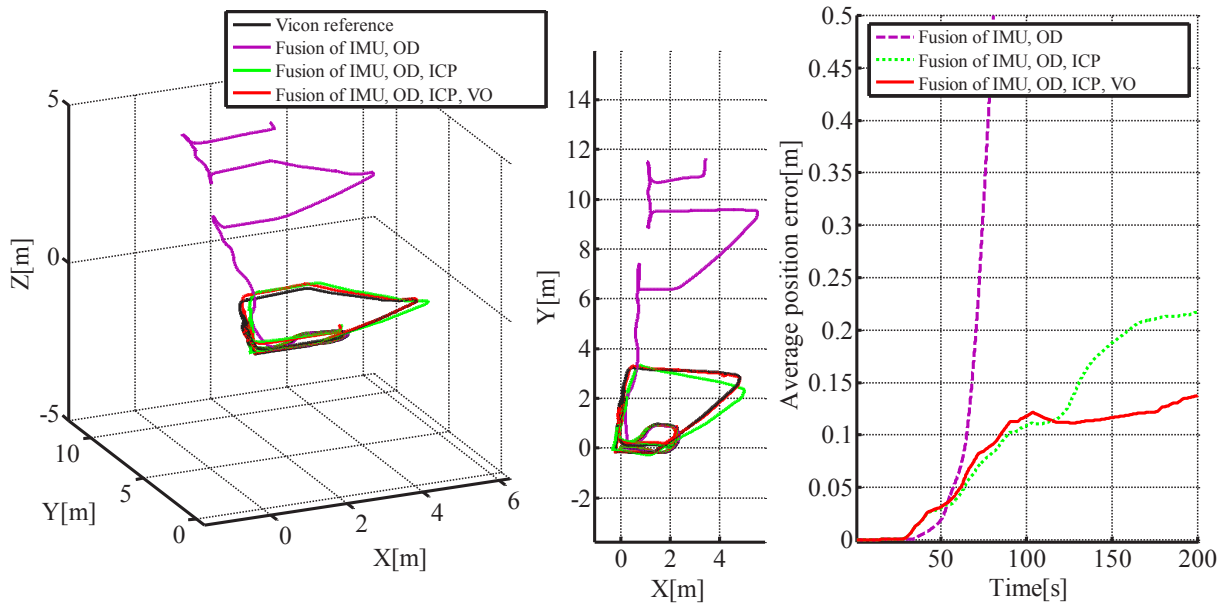


Figure 8: Trajectories obtained by fusing different combinations of modalities during the indoor experiment testing obstacle (depicted in Figure 7) traversability under high slippage (left, middle); development of the average position error (right).

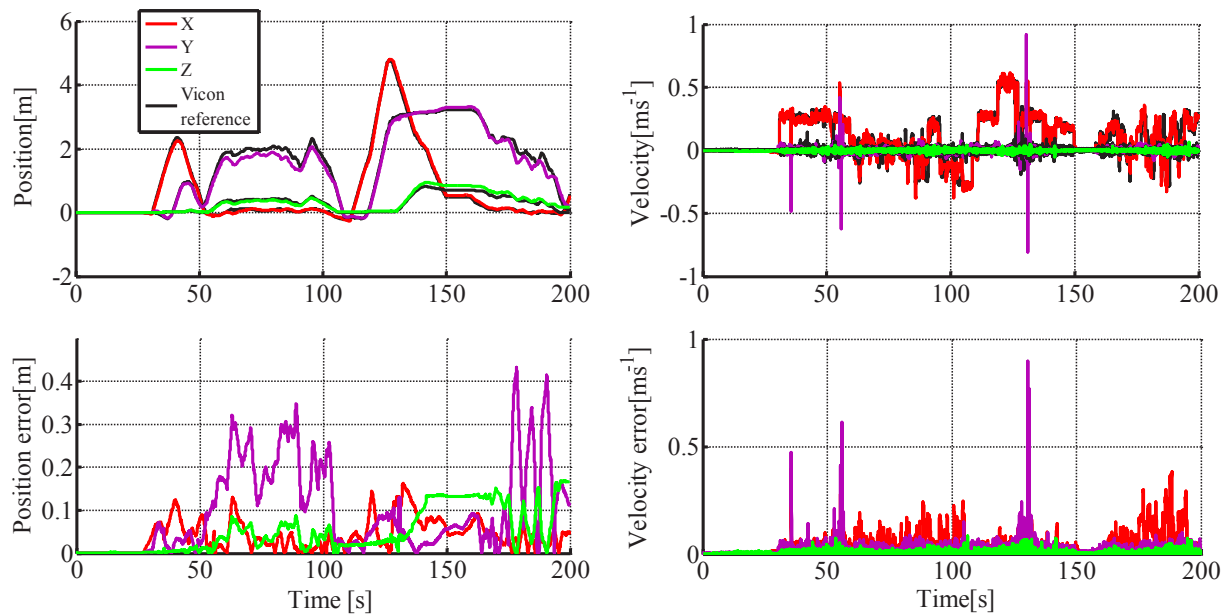


Figure 9: The corrected position (top left) and velocity estimates (top right) for the IMU+OD+ICP+VO combination corresponding to the trajectory in Figure 8 (testing obstacle traversability). Errors in position and velocity are obtained as norm of difference between the Vicon reference and the corresponding state at each time-step (bottom left, bottom right). The Vicon reference for both position and velocity is shown in black.

are shown in Figure 14.

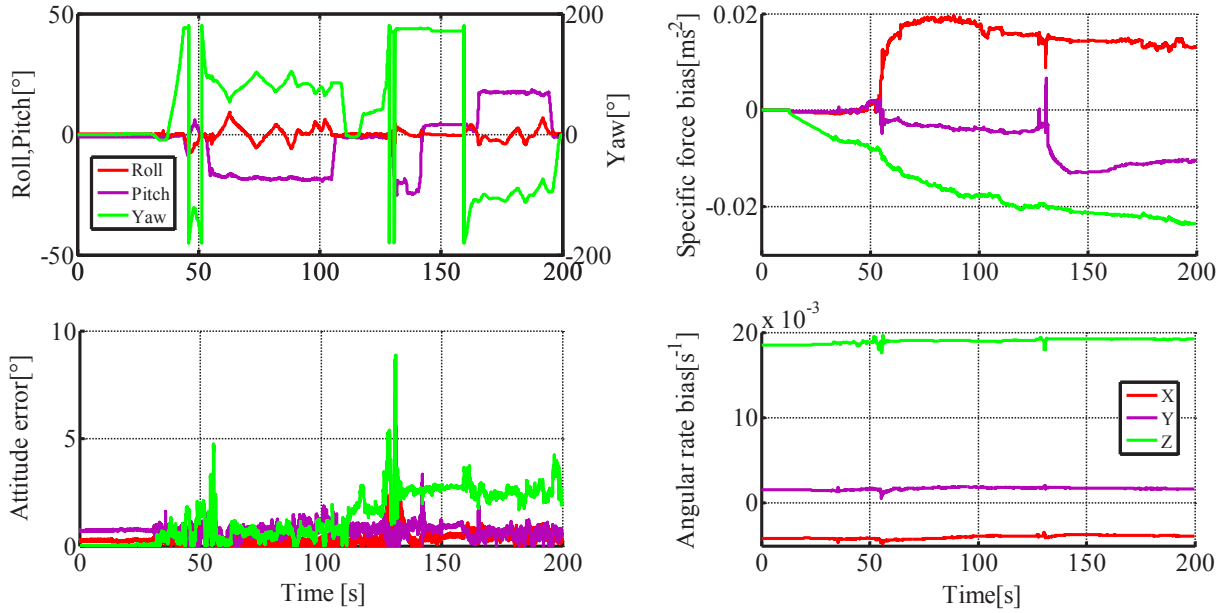


Figure 10: The corrected attitude estimates (top left) for the full multi-modal combination IMU+OD+ICP+VO corresponding to the trajectory shown in Figure 8 (testing obstacle traversability). Errors in attitude are obtained as the difference between the Vicon reference and the corresponding state at each time-step (bottom left). Estimated biases for the specific forces (top right) and angular rates (bottom right).



Figure 11: An example of trajectory driven by the robot over the Clausiusstrasse street.

5.3.3 Evaluation of the measurement model

We claim that a standard measurement model—as usually used for measurements coming at comparable frequency—is not well suitable for measurements with significant differences in sampling frequencies as well as in values, which correspond to the same state observed. This is crucial, when the difference in states obtained from the IMU or the OD at high frequency is very large compared to the measurements provided by the ICP or the VO sensory modalities at relatively low frequency—such as in case of high slippage.

Table 4) shows the overall comparison of the three measurement models we evaluated for fusing the ICP and the VO sensory modalities in the filter. Figure 15 presents a typical example of trajectory reconstructed by all the three measurement approaches we introduced in Section 4.3.3. The *velocity approach*—the state-of-the-art practice—that considers those information as relative measurements, is the least precise, with the

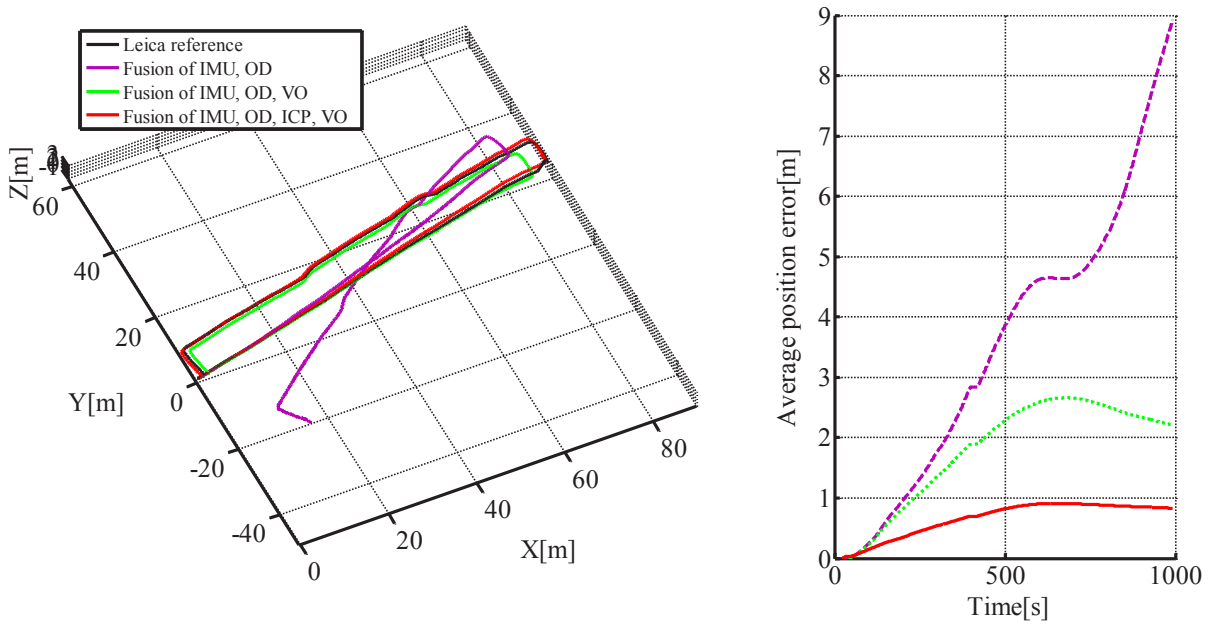


Figure 12: Trajectories obtained by fusing different combinations of modalities during the outdoor experiment with Leica reference system (left) and the corresponding average position error in time (right).

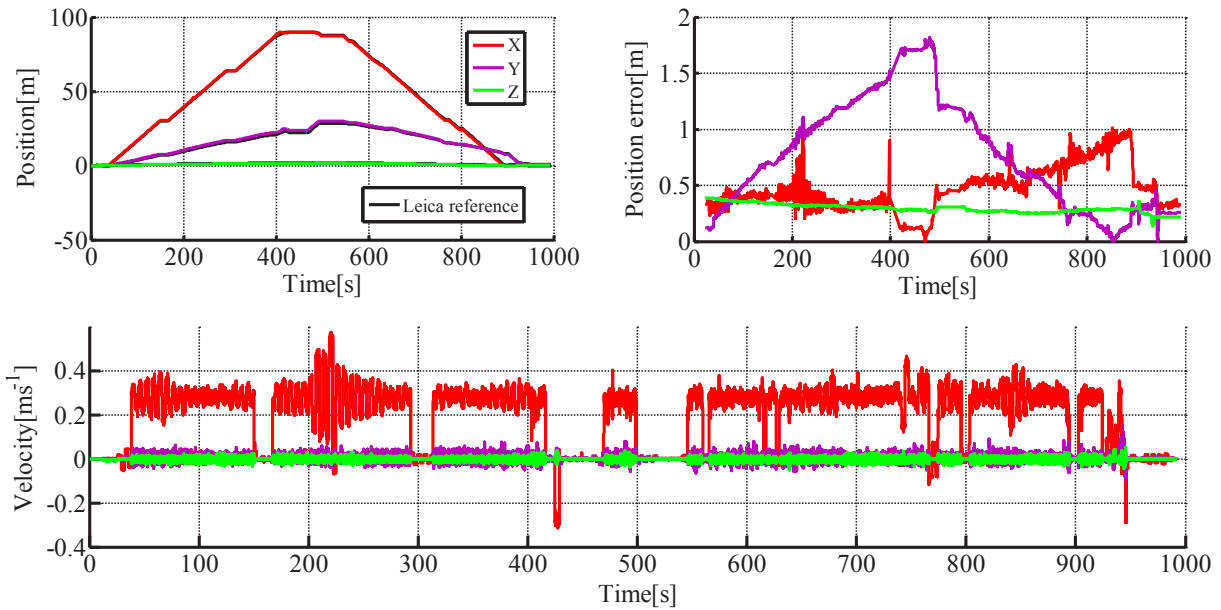


Figure 13: The position and velocity estimates (top left and bottom respectively) for the IMU+OD+ICP+VO combination corresponding to the outdoor trajectory in Figure 12; errors in position obtained as norm of differences between the Leica reference and the corresponding state at each time-step (top right).

highest average position error; see Figure 15 (right). This is due to the *corner cutting* behavior emphasized in Figure 15 (middle). The *incremental position approach* performs reasonably well in indoor environments, which are well -conditioned for the ICP and the VO sensory modalities—especially the ICP algorithm is really precise as there are enough features to unambiguously fix all degrees of freedom. On the other hand,

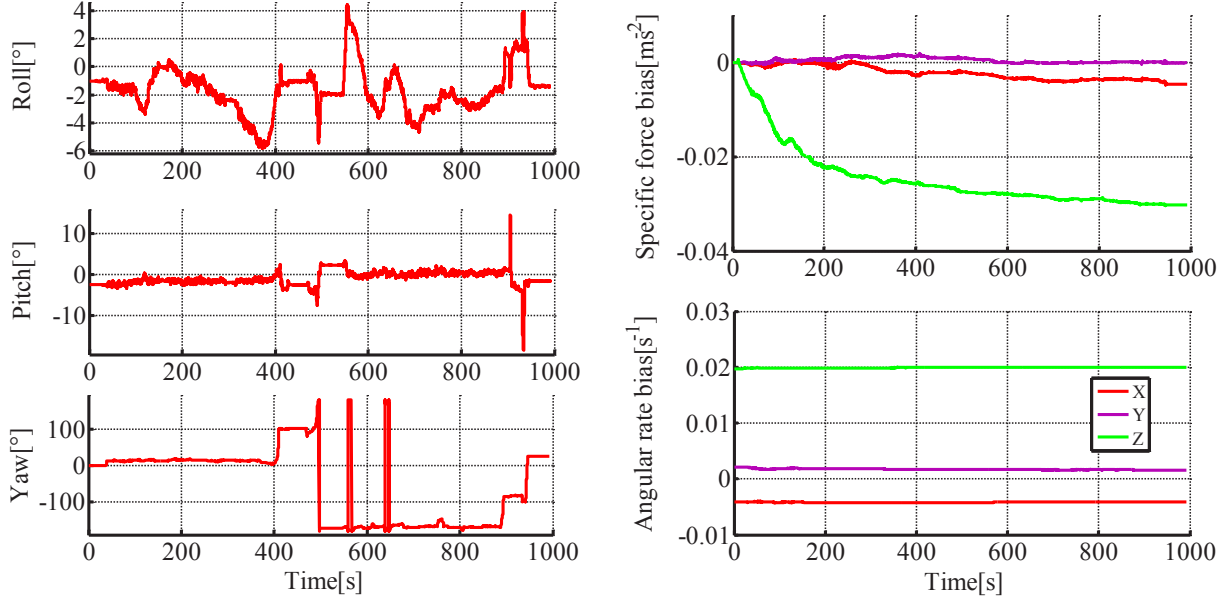


Figure 14: The attitude estimates (left) for the IMU+OD+ICP+VO combination corresponding to the outdoor trajectory in Figure 12; biases estimated for the specific forces (top right) and angular rates (bottom right).

Model	Indoor		Outdoor	
	e_{rel}	e_{avg}	e_{rel}	e_{avg}
incremental position	0.4 0.7 1.2	0.1 0.1 0.2	0.8 1.5 11.0	0.7 2.4 6.1
velocity	1.0 1.3 2.3	0.1 0.1 0.3	0.9 1.8 12.2	0.8 2.5 6.1
trajectory	0.7 1.2 2.1	0.0 0.1 0.2	0.6 1.4 11.5	0.6 2.2 6.1

Table 4: Comparison of the different measurement models; for each model we show the lower|**median**|higher quartile statistics of the relative and average metrics. The average metric e_{avg} is evaluated for the last sample of each experiment, see (49). We distinguish the in- and outdoor environments.

in larger environments with less constraints (expected for USAR), the *trajectory approach* allows the IMU and the OD information to better correct the drift of the ICP and the VO sensory modalities.

5.4 Failure case analysis

As seen in the previous sections, there are plenty of occasions in USAR environments for which the generic assumptions of the EKF are not valid. The most frequent example is track slippage that violates the assumption of Gaussian observation centered on the actual value.

Our failure case analysis reviews each sensory modality involved in the filter to see how the resulting estimate degrades with partial outage of the modality. IMUs are not subject to much partial failure other than bias and noise, that are already accounted for in our filter.

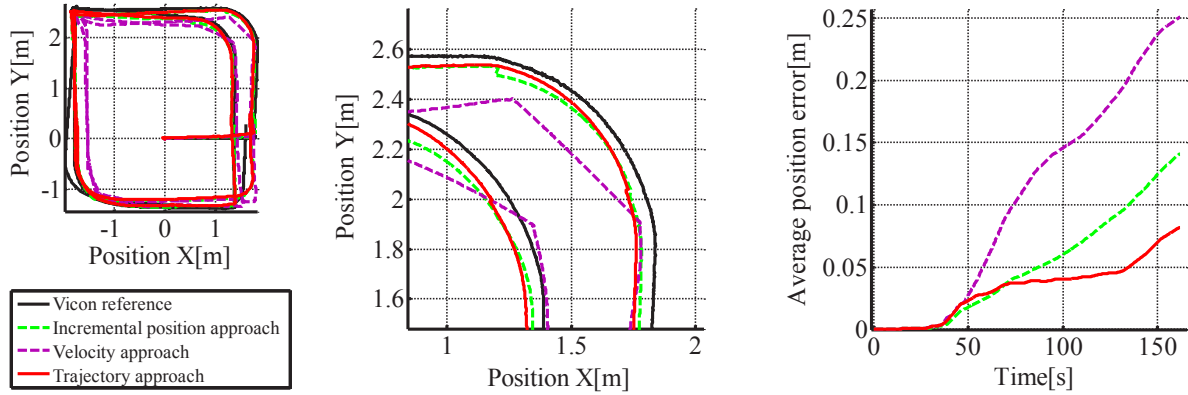


Figure 15: Comparison of effects of the three different ICP aiding approaches on the estimated trajectory (left, middle) and on the average position error (right). Note the *corner cutting* effect of the *velocity approach*.

5.4.1 Robot slippage and sliding

A typical failure case of the odometry modality is significant slippage. Small slippage occurs routinely when turning skid-steer robots and is usually accounted for by the uncertainty in the odometry model. However, on surfaces like ice, or inclined wet or smooth surfaces, stronger slippage can occur. Stronger slippages or sliding are outliers of the odometry observation model. IMU, ICP and VO sensory modalities are not affected in such a case. In order to simulate such a situation, we placed the robot on a trolley and moved it manually.

Figure 16 shows both the trajectory from the top (top-left plot) and the comparison between the fusion of all four sensory modalities and the fusion of only IMU+OD. We can see that the latter wrongly estimates no motion whereas the fusion of all modalities correctly estimates the trajectory. The failure of the partial filter can be explained by the low acceleration of the platform during the test. As the IMU acceleration signal is quite noisy, confidence on the IMU cannot compensate for the odometry modality asserting an absence of motion.

It should be noted that such a failure of the odometry modality does not lead to a failure of our complete filter.

5.4.2 Partial occlusion of visual field of view

Partial occlusion, overexposure, or projections of dirt on the camera could lead to faulty estimation of the motion by the VO. In order to test this situation, we occluded one of the cameras of the omnicaamera (see Figure 17). Reduction of the field of view of the omnicaamera causes in the vast majority of cases a reduction in the number of visual features being robustly detected by the VO. The insufficient number of features can then cause the VO to incorrectly estimate the attitude. This information then propagates into the state estimate and can make the fusion algorithm fail.

Figure 18 shows the result of the filter in such a case. We can see that during a first loop of the trajectory, the state estimation is correct. Then, lacking a sufficient number of features, the VO computes an erroneous estimate and the final state estimate degenerates. On the contrary, by leaving out the visual odometry, the state estimation would continue to perform satisfactorily.

It is to be noted that more than the portion of the field of view occluded, it is the number, quality and distribution of features that matters. One typical way to prevent this issue is to monitor the number of

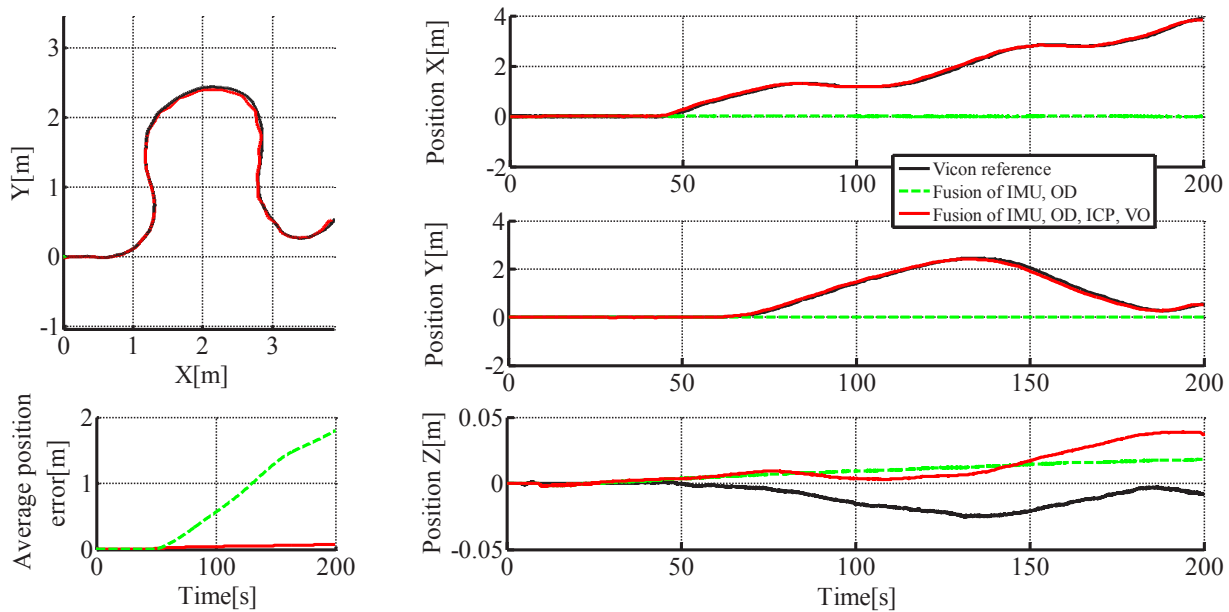


Figure 16: Test trajectory for robot slippage. Black line: ground-truth; red solid line: state estimate with all four modalities; green dashed line: IMU and odometry fusion. Top left: top view of the trajectory; bottom left: average error as a function of time; top, middle, bottom right: evolution of x , y , z coordinate.

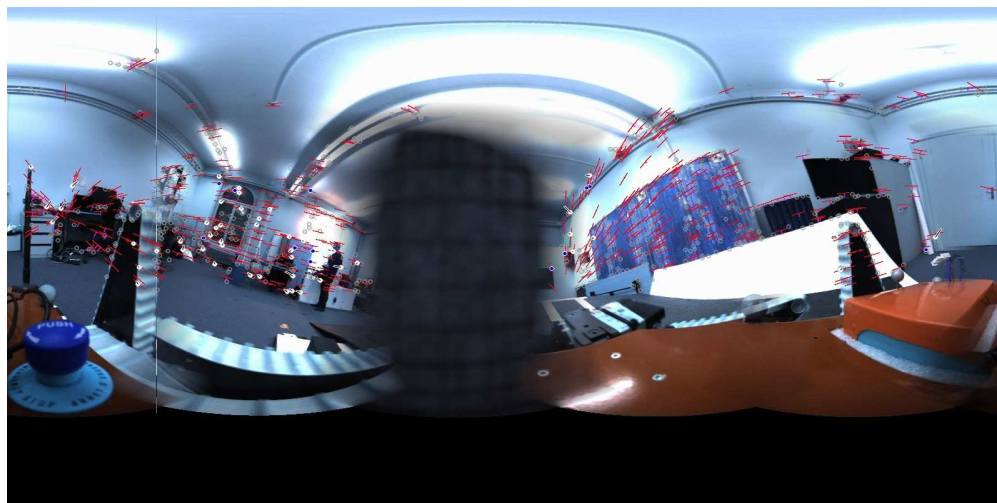


Figure 17: Picture from the partially occluded omniscamera. Notice the dark rectangle in the middle

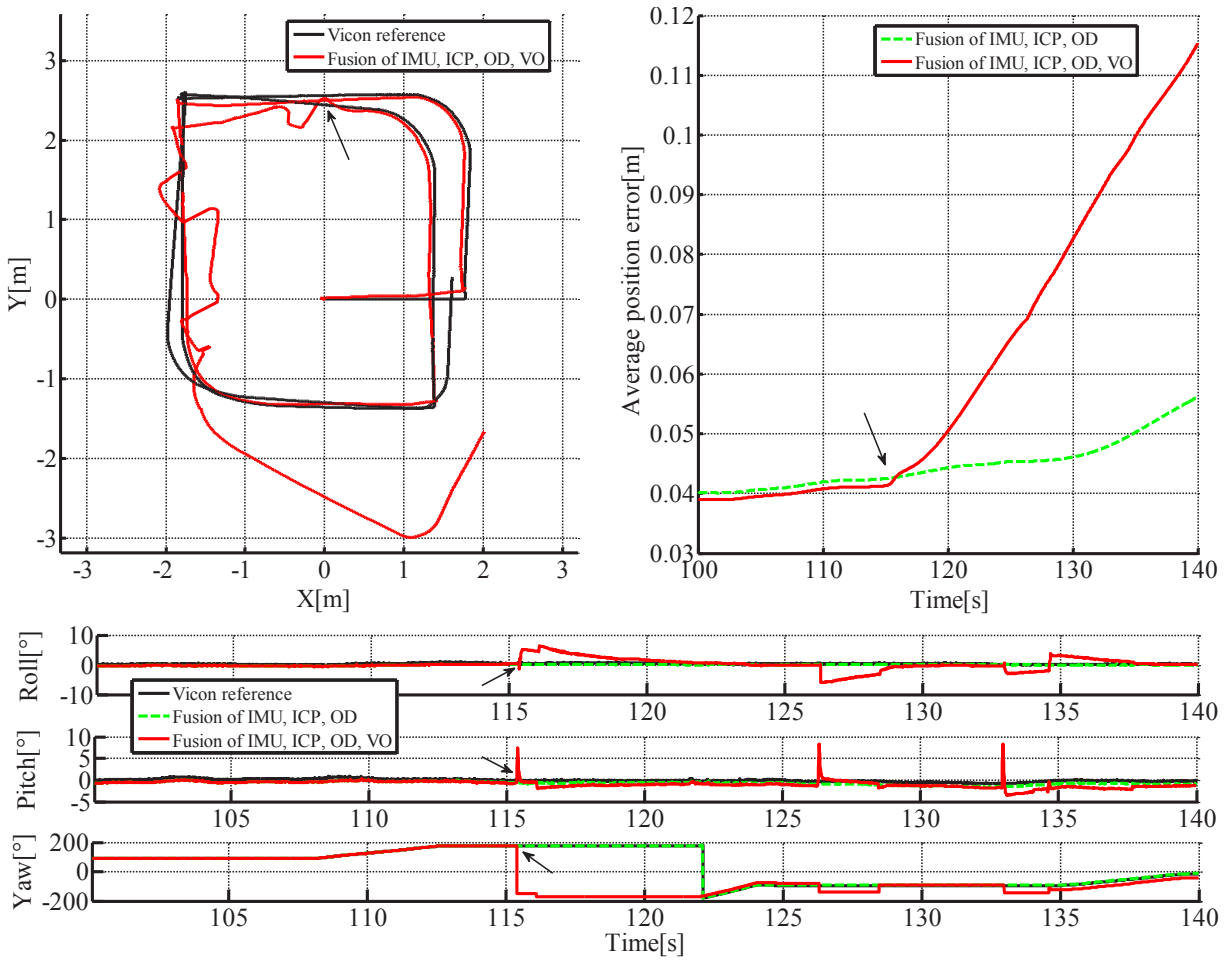


Figure 18: Trajectory reconstruction with partial omnicaamera occlusion. Black line: ground truth; Solid red line: state estimate with all four modalities; Dashed green line: state estimate excluding visual odometry; Black arrow: visual odometry failure. Top left: top view of the trajectory; Top right: average position error around visual odometry failure; Bottom: attitude estimated along the trajectory.

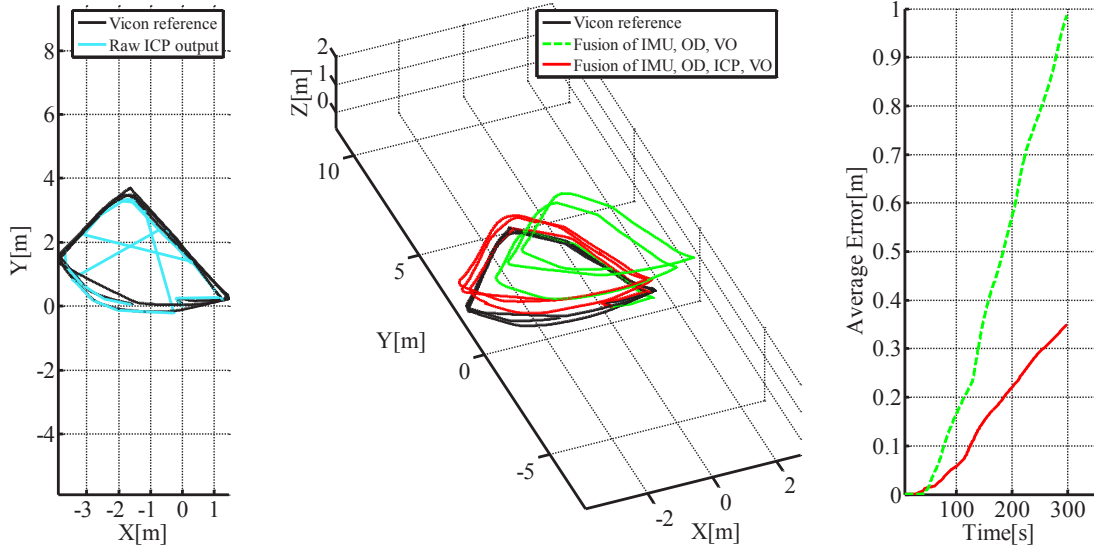


Figure 19: Trajectory estimates in case of low ICP frequency. Black line: ground truth; cyan line: positions estimated by ICP alone; red line: state estimate with all four modalities; green dashed line: state estimate excluding ICP measurements. Left: top view; middle: 3D view; right: average position error.

features and eventually their distribution in the field of view—our VO tries to have corresponding features spread over the whole image. As Figure 18 shows, even with the partial occlusion of the field of view, the VO performed correctly during most of the trial. This observation holds also for too dark or over-exposed images.

5.4.3 Temporary laser scanner outage

As demonstrated above, our trajectory approach to fusion of ICP measurements is able to cope with the relatively low frequency of laser scanning. As the laser is moving, it can be blocked in case of collision or high vibration of the platform (safety precaution at the level of the motor controller). When this happens, it is necessary to initiate a recalibration procedure that can take around 30 s.

We simulated this situation by throttling the laser point clouds, which resulted in ICP measurement outages of up to 40 s. Figure 19 shows the trajectory estimates for this test. On the left the cyan polygon shows the position estimates of ICP linked by straight lines (no filtering). It is to be noted that in this case, the positions are accurate compared to the ground-truth but of very low sampling rate. We can see in the middle and right graphs that the filter estimates degrade gracefully. There is some drift, mostly along elevation due to slippage, but even with this low frequency, the ICP measurements help correcting the state estimates over just the IMU, odometry and visual odometry.

5.4.4 Moving obstacle and limited laser range

Unlike the cameras, laser range sensors are not sensitive to illumination conditions. On the other hand, they have a limited sensor range which can induce a lack of points in large environments. Close range obstacles might then be the dominant cluster of points and hence the ICP registration might converge to a wrong local minimum, following the motion of the obstacles.

In order to test this situation, we artificially limited the range of the laser range sensor to 2 m. This is similar to heavy smoke or dust scenarios that can arise in USAR conditions. This prevents the laser to observe the

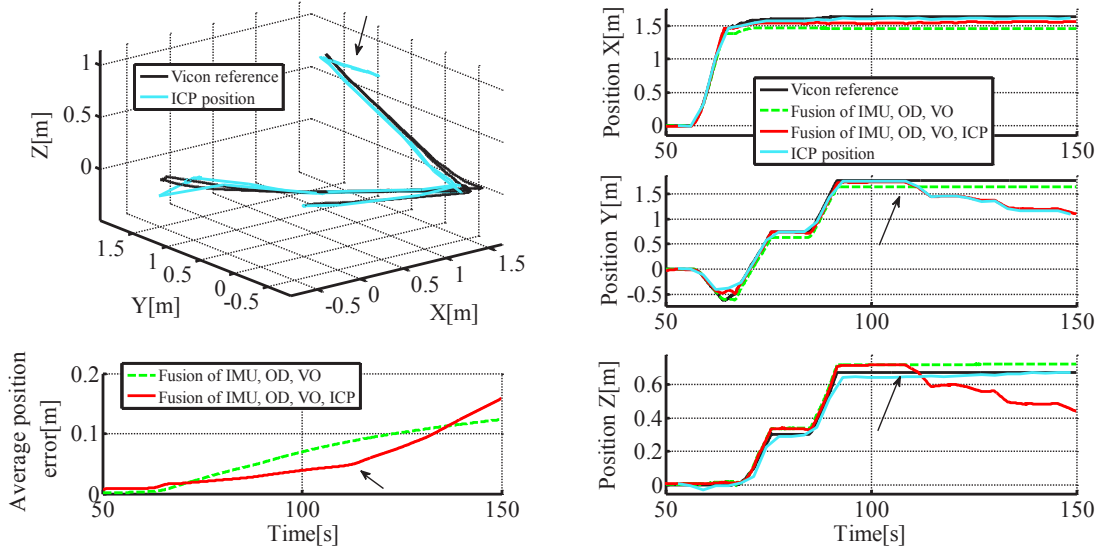


Figure 20: Trajectory estimates in case of moving obstacle in reduced field of view. Solid black line: ground truth; solid red line: state estimate with all four modalities; dashed green line: state estimate excluding ICP measurements; cyan line: position estimated by ICP alone; black arrow: start of moving obstacle. Top left: 3D view; bottom left: average error as a function of time; right: x , y , and z coordinates as a function of time.

walls and the ceiling that are, indoors, usually the strongest cues for correct point cloud registration.

Additionally, we used a large board to simulate a moving obstacle of significant size. This caused the ICP to drift, following the motion of the board.

Figure 20 shows the result of the filter compared to ground truth. We can see that when the large obstacle starts to move, the estimate of the ICP drifts with it. As a consequence the whole filter drifts as well. This is analogous to the slippage situation, in which the ICP modality compensates the combined estimate of the other three modalities. Using the omnicaamera information not only as visual compass but also as a complete visual odometry modality would probably allow to make the difference between those two situations.

5.4.5 Map deformation

As explained above, the ICP map is not globally optimized. This means that the map might have some large scale deformations due to the accumulation of small errors. We were able to observe this particularly in a long corridor that we used to assess the impact of map deformation on the state estimate.

Figure 21 shows an instance of the deformed map. We drove along two superposed corridors over two floors. We can see that both ends of the corridor are not aligned: the ground plane of the blue end has a roll angle of several degrees compared to the red end. We used the theodolite system to acquire ground truth on the upper floor.

Figure 22 shows the impact of map deformation on the state estimate. The top graph shows that even if the ICP estimate is erroneous, the full filter maintains a correct, drift-free estimate. The bottom graph compares the estimate of the roll angle between ICP only and the fusion. It clearly shows the drift in roll of the ICP estimate and the lack of impact it has on the fusion. The difference with previous failure case lies on the kind of drift. The drift of the roll angle can be compensated for by the IMU, especially the accelerometer. On the other hand, the drift in position of previous failure case is not observable by the other modalities.

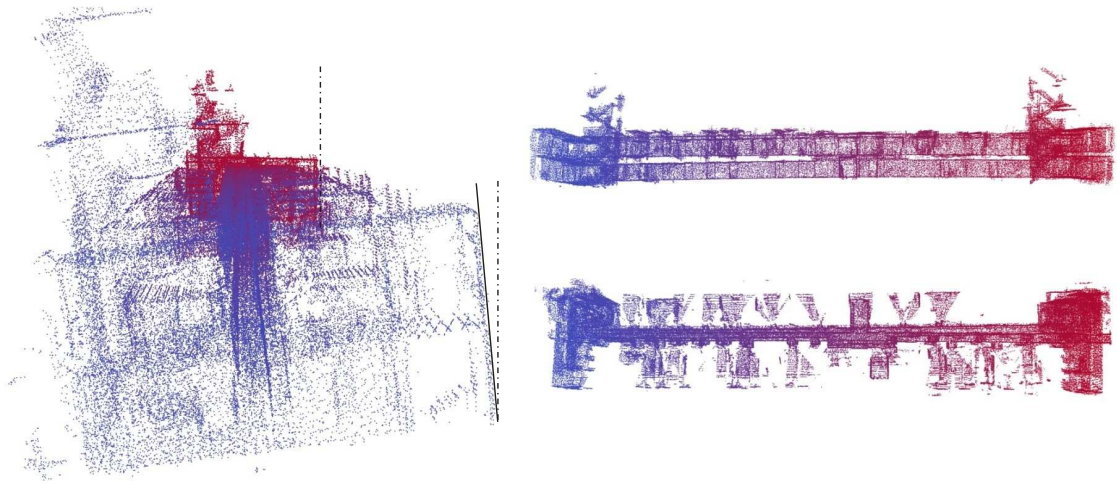


Figure 21: Deformed point cloud map created by ICP. The points are colored alongside the corridor from red (initial position) to blue. Left: front view; top right: side view; bottom right: top view.

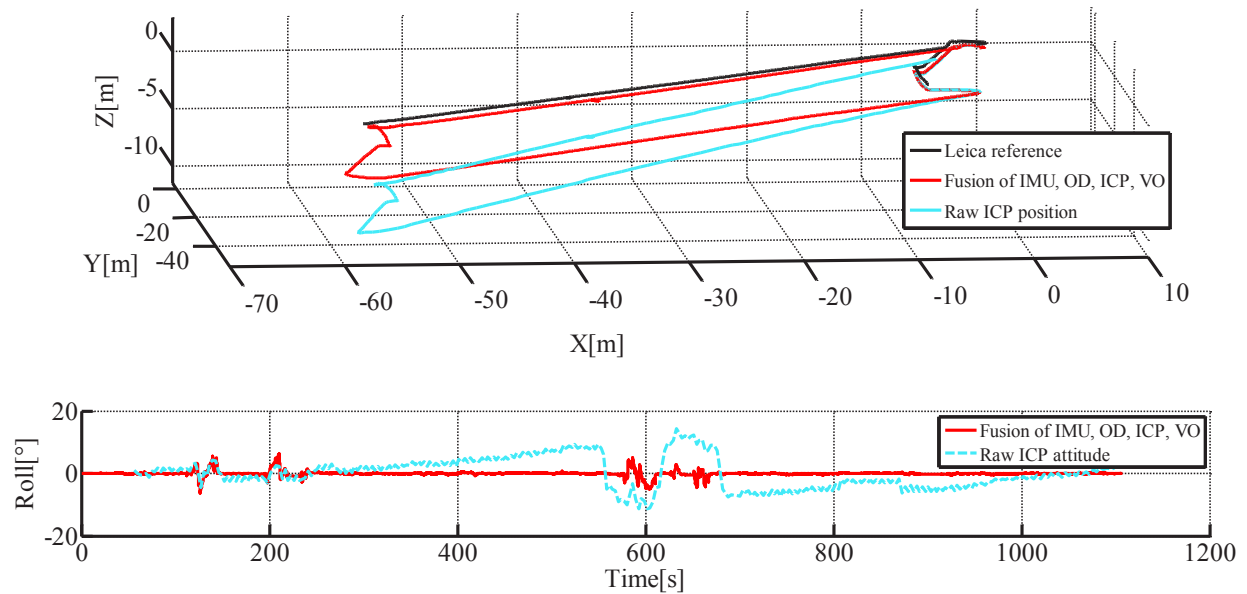


Figure 22: Trajectory estimates in case of map deformation. Solid black line: ground truth; solid red line: state estimate with all four modalities; cyan line: position estimated by ICP alone. Top: side view; bottom: roll angle along the trajectory.

6 Conclusion

We designed and evaluated a multi-modal data fusion system for state estimation of a mobile skid-steer robot intended for Urban Search and Rescue missions. USAR missions often involve in- and outdoor environments with challenging conditions like slippage, moving obstacles, bad or changing light conditions, etc. In order to cope with such environment, our robot is equipped with both proprioceptive (IMU, tracks odometry) and exteroceptive (laser rangefinder, omnidirectional camera) sensors. We designed such a data fusion scheme in order to adequately include measurements from all these four modalities with order of magnitude difference in update frequency from 90 Hz to $\frac{1}{3}$ Hz.

We tested our algorithm on approximately 4.4 km of field tests (over more than 9 hours of data) both in- and outdoors. In order to have precise quantitative analysis, we recorded ground truth using either a Vicon motion capture system (indoors) or a Leica theodolite tracker (outdoors). This way we proved that our scheme is a significant improvement upon standard approaches. Combining all four modalities: IMU, tracks odometry, visual odometry and ICP-based localization, we achieved precision in the total distance driven of 1.2% error in the indoor environment and 1.4% error in the outdoor environment. Moreover, we characterized the reliability of our data fusion scheme against sensor failures. We designed failure case scenarios according to potential failures of each sensory modality that are likely to occur during real USAR missions. In course of this testing, we evaluated robustness with respect to: heavy slippage (odometry failure case), reduction of field of view of the omnicaamera (visual odometry failure case), and reduction of the laser rangefinder together with large moving obstacles spoiling the created metric map (ICP-based localization failure case).

While our filter demonstrates good accuracy during our field tests and is robust against some of the failures expected in USAR, there is still a space for improvement—the need for an automatic failure detection and resolution. Exploring different methods of detecting anomalous measurements and rejecting them in order to improve the overall performance is one of the ways, but it is currently left for future work. Furthermore, developing a visual odometry solution capable of providing also estimates of scaled translation is another topic for the future.

It is not surprising that combining more modalities yields more precision. But we were able to show that if such a rich multi-modal system is well designed, it will perform reasonably well even in cases, where other systems exploiting fewer modalities fail completely. We describe how to design such system using the commonly used EKF. In this way we contribute by proposing and comparing three different approaches to treat the ICP measurements; out of which the *trajectory approach* proved to perform best.

To contribute to the robotics community, we release our datasets used in this paper, including the ground truth measurements from Vicon and Leica systems.

Acknowledgments

The research presented here was supported by the European Union FP7 Programme under the NIFTi project (No. 247870; <http://www.nifti.eu>) and the TRADR project (No. 609763; <http://www.tradr-project.eu>). François Pomerleau was supported by a fellowship from the Fonds québécois de recherche sur la nature et les technologies (FQRNT). Vladimir Kubelka was supported by the Czech Science Foundation (Project Registration No. 14-13876S). We would like to thank the anonymous reviewers for their constructive suggestions that greatly improved the manuscript.

References

- Almeida, J. and Santos, V. M. (2013). Real time egomotion of a nonholonomic vehicle using lidar measurements. *Journal of Field Robotics*, 30(1):129–141.
- Anousaki, G. and Kyriakopoulos, K. J. (2004). A dead-reckoning scheme for skid-steered vehicles in outdoor environments. In *Robotics and Automation (ICRA), 2004. Proceedings of the IEEE International Conference on*, pages 580–585.
- Bachrach, A., Prentice, S., He, R., and Roy, N. (2011). RANGE—Robust autonomous navigation in GPS-denied environments. *Journal of Field Robotics*, 28(5):644–666.
- Barfoot, T., Stenning, B., Furgale, P., and McManus, C. (2012). Exploiting reusable paths in mobile robotics: Benefits and challenges for long-term autonomy. In *Computer and Robot Vision (CRV), 2012 Ninth Conference on*, pages 388–395.
- Besl, P. and McKay, H. (1992). A method for registration of 3-D shapes. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 14(2):239–256.
- Breckenridge, W. G. (1999). Quaternions - proposed standard conventions. Technical report, JPL.
- Brodsky, T., Fermueller, C., and Aloimonos, Y. (1998). Directions of motion fields are hardly ever ambiguous. *International Journal of Computer Vision*, 26(1):5–24.
- Chen, Y. and Medioni, G. (1991). Object modeling by registration of multiple range images. In *Robotics and Automation (ICRA), 1991. Proceedings of the IEEE International Conference on*, pages 2724–2729.
- Chetverikov, D., Svirko, D., Stepanov, D., and Krsek, P. (2002). The Trimmed Iterative Closest Point algorithm. In *Pattern Recognition, 2002. Proceedings of the 16th International Conference on*, pages 545–548.
- Chiu, H.-P., Williams, S., Dellaert, F., Samarasekera, S., and Kumar, R. (2013). Robust vision-aided navigation using sliding-window factor graphs. In *Robotics and Automation (ICRA), 2013 IEEE International Conference on*, pages 46–53.
- Chowdhary, G., Johnson, E. N., Magree, D., Wu, A., and Shein, A. (2013). GPS-denied Indoor and Outdoor Monocular Vision Aided Navigation and Control of Unmanned Aircraft. *Journal of Field Robotics*, 30(3):415–438.
- Civera, J., Grasa, O. G., Davison, A. J., and Montiel, J. M. M. (2010). 1-Point RANSAC for extended Kalman filtering: Application to real-time structure from motion and visual odometry. *Journal of Field Robotics*, 27(5):609–631.
- Dissanayake, G., Sukkariéh, S., Nebot, E., and Durrant-Whyte, H. (2001). The aiding of a low-cost strap-down inertial measurement unit using vehicle model constraints for land vehicle applications. *IEEE Transactions on Robotics and Automation*, 17(5):731–747.
- Ellekilde, L.-P., Huang, S., Miro, J. V., and Dissanayake, G. (2007). Dense 3D Map Construction for Indoor Search and Rescue. *J. Field Robotics*, 24(1-2):71–89.
- Endo, D., Okada, Y., Nagatani, K., and , K. (2007). Path following control for tracked vehicles based on slip-compensating odometry. In *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems IROS 2007*, pages 2871–2876.
- Fraundorfer, F. and Scaramuzza, D. (2012). Visual Odometry : Part II: Matching, Robustness, Optimization, and Applications. *Robotics Automation Magazine, IEEE*, 19(2):78–90.
- Galben, G. (2011). New Three-Dimensional Velocity Motion Model and Composite Odometry–Inertial Motion Model for Local Autonomous Navigation. *IEEE Transactions on Vehicular Technology*, 60(3):771–781.

- Jesus, F. and Ventura, R. (2012). Combining monocular and stereo vision in 6d-slam for the localization of a tracked wheel robot. In *Safety, Security, and Rescue Robotics (SSRR), 2012 IEEE International Symposium on*, pages 1–6.
- Kalman, R. E. (1960). A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):34–45.
- Kelly, J., Sibley, G., Barfoot, T., and Newman, P. (2012). Taking the Long View: A Report on Two Recent Workshops on Long-Term Autonomy. *Robotics & Automation Magazine, IEEE*, 19(1):109–111.
- Kohlbrecher, S., Stryk, O. V., Meyer, J., and Klingauf, U. (2011). A flexible and scalable slam system with full 3d motion estimation. In *Safety, Security, and Rescue Robotics (SSRR), 2011 IEEE International Symposium on*, pages 155–160.
- Konolige, K., Agrawal, M., and Sola, J. (2011). Large-scale visual odometry for rough terrain. In Kaneko, M. and Nakamura, Y., editors, *Robotics Research*, volume 66 of *Springer Tracts in Advanced Robotics*, pages 201–212. Springer.
- Kruijff, G. J. M., Janicek, M., Keshavdas, S., Larochelle, B., Zender, H., Smets, N. J. J. M., Mioch, T., Neerincx, M. A., van Diggelen, J., Colas, F., Liu, M., Pomerleau, F., Siegwart, R., Hlavac, V., Svoboda, T., Petricek, T., Reinstein, M., Zimmerman, K., Pirri, F., Gianni, M., Papadakis, P., Sinha, A., Balmer, P., Tomatis, N., Worst, R., Linder, T., Surmann, H., Tretyakov, V., Surmann, H., Corrao, S., Pratzler-Wanczura, S., and Sulk, M. (2012). Experience in System Design for Human-Robot Teaming in Urban Search and Rescue. In *Field and Service Robotics*, pages 1–14, Matsushima, Japan.
- Kubelka, V. and Reinstein, M. (2012). Complementary filtering approach to orientation estimation using inertial sensors only. In *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, pages 599–605.
- Kummerle, R., Grisetti, G., Strasdat, H., Konolige, K., and Burgard, W. (2011). g2o: A general framework for graph optimization. In *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pages 3607–3613. IEEE.
- Lamon, P. and Siegwart, R. (2004). Inertial and 3D-odometry fusion in rough terrain - towards real 3D navigation. In *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems (IROS 2004)*, pages 1716–1721.
- Li, H. and Hartley, R. (2006). Five-point motion estimation made easy. In *Pattern Recognition, 2006. ICPR 2006. 18th International Conference on*, volume 1, pages 630–633.
- Loan, C. V. (1978). Computing integrals involving the matrix exponential. *Automatic Control, IEEE Transactions on*, 23(3):395–404.
- Ma, J., Susca, S., Bajracharya, M., Matthies, L., Malchano, M., and Wooden, D. (2012). Robust multi-sensor, day/night 6-dof pose estimation for a dynamic legged vehicle in gps-denied environments. In *Robotics and Automation (ICRA), 2012 IEEE International Conference on*, pages 619–626.
- McElhoe, B. A. (1966). An Assessment of the Navigation and Course Corrections for a Manned Flyby of Mars or Venus. *Aerospace and Electronic Systems, IEEE Transactions on*, 2(4):613–623.
- Morales, Y., Carballo, A., Takeuchi, E., Aburadani, A., and Tsubouchi, T. (2009). Autonomous robot navigation in outdoor cluttered pedestrian walkways. *Journal of Field Robotics*, 26(8):609–635.
- Nagatani, K., Okada, Y., Tokunaga, N., Kiribayashi, S., Yoshida, K., Ohno, K., Takeuchi, E., Tadokoro, S., Akiyama, H., Noda, I., Yoshida, T., and Koyanagi, E. (2011). Multirobot exploration for search and rescue missions: A report on map building in robocuprescue 2009. *Journal of Field Robotics*, 28(3):373–387.
- Nemra, A. and Aouf, N. (2010). Robust INS/GPS Sensor Fusion for UAV Localization Using SDRE Nonlinear Filtering. *IEEE Sensors Journal*, 10(4):789–798.

- Nuchter, A., Lingemann, K., Hertzberg, J., and Surmann, H. (2007). 6D SLAM - 3D mapping outdoor environments. *Journal of Field Robotics*, 24(8-9):699–722.
- Oskiper, T., Chiu, H.-P., Zhu, Z., Samarasekera, S., and Kumar, R. (2010). Multi-modal sensor fusion algorithm for ubiquitous infrastructure-free localization in vision-impaired environments. In *Intelligent Robots and Systems (IROS), 2010 IEEE/RSJ International Conference on*, pages 1513–1519.
- Pomerleau, F., Colas, F., Siegwart, R., and Magnenat, S. (2013). Comparing ICP Variants on Real-World Data Sets. *Autonomous Robots*, 34(3):133–148.
- Reinstein, M. and Hoffmann, M. (2013). Dead reckoning in a dynamic quadruped robot based on multimodal proprioceptive sensory information. *IEEE Transactions on Robotics*, 29(2):563–571.
- Reinstein, M., Kubelka, V., and Zimmermann, K. (2013). Terrain adaptive odometry for mobile skid-steer robots. In *Proc. IEEE Int Robotics and Automation (ICRA) Conf*, pages 4706–4711.
- Rodriguez F, S. A., Fremont, V., and Bonnifait, P. (2009). An experiment of a 3D real-time robust visual odometry for intelligent vehicles. In *Proc. 12th Int. IEEE Conf. Intelligent Transportation Systems ITSC '09*, pages 1–6.
- Rublee, E., Rabaud, V., Konolige, K., and Bradski, G. (2011). ORB: An efficient alternative to SIFT or SURF. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2564 –2571.
- Sakai, A., Tamura, Y., and Kuroda, Y. (2009). An efficient solution to 6DOF localization using Unscented Kalman Filter for planetary rovers. In *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems IROS 2009*, pages 4154–4159.
- Savage, P. G. (1998). Strapdown inertial navigation integration algorithm design part 2: Velocity and position algorithms. *Journal of Guidance, Control, and Dynamics*, 21(No. 2):208 – 221.
- Scaramuzza, D. and Fraundorfer, F. (2011). Visual odometry [tutorial]. *IEEE Robot Autom Mag*, 18(4):80–92.
- Shen, J., Tick, D., and Gans, N. (2011). Localization through fusion of discrete and continuous epipolar geometry with wheel and IMU odometry. In *Proc. American Control Conf. (ACC)*, pages 1292–1298.
- Smith, G. L., Schmidt, S. F., and McGee, L. A. (1962). Optimal filtering and linear prediction applied to a midcourse navigation system for the circumlunar mission. Technical report, U.S. Government Printing Office.
- Sukumar, S. R., Bozdogan, H., Page, D. L., Koschan, A. F., and Abidi, M. A. (2007). Sensor selection using information complexity for multi-sensor mobile robot localization. In *Robotics and Automation, 2007 IEEE International Conference on*, pages 4158–4163.
- Suzuki, T., Kitamura, M., Amano, Y., and Hashizume, T. (2010). 6-DOF localization for a mobile robot using outdoor 3D voxel maps. In *Proc. IEEE/RSJ Int Intelligent Robots and Systems (IROS) Conf*, pages 5737–5743.
- Svoboda, T., Pajdla, T., and Hlaváč, V. (1998). Motion estimation using central panoramic cameras. In *IEEE International Conference on Intelligent Vehicles*, pages 335–340.
- Tardif, J., Pavlidis, Y., and Daniilidis, K. (2008). Monocular visual odometry in urban environments using an omnidirectional camera. In *Intelligent Robots and Systems, 2008. IROS 2008. IEEE/RSJ International Conference on*, pages 2531–2538.
- Titterton, D. H. and Weston, J. L. (1997). *Strapdown Inertial Navigation Technology*. The Lavenham Press ltd, Lavenham, UK.
- Trawny, N. and Roumeliotis, S. I. (2005). Indirect Kalman Filter for 3D Attitude Estimation - A Tutorial for Quaternion Algebra. Technical report, University of Minnesota.

- Weiss, S. M. (2012). Vision based navigation for micro helicopters. Dissertation, ETH Zurich.
- Yi, J., Zhang, J., Song, D., and Jayasuriya, S. (2007). Imu-based localization and slip estimation for skid-steered mobile robots. In *Proc. IEEE/RSJ Int. Conf. Intelligent Robots and Systems IROS 2007*, pages 2845–2850.
- Yoshida, T., Irie, K., Koyanagi, E., and Tomono, M. (2010). A sensor platform for outdoor navigation using gyro-assisted odometry and roundly-swinging 3D laser scanner. In *Proc. IEEE/RSJ Int Intelligent Robots and Systems (IROS) Conf*, pages 1414–1420.