

# An Open Source, Fiducial Based, Visual-Inertial Motion Capture System

Michael Neunert

Agile & Dexterous Robotics Lab  
ETH Zürich, Switzerland  
Email: neunertm@ethz.ch

Michael Bloesch

Autonomous Systems Lab  
ETH Zürich, Switzerland  
Email: bloeschm@ethz.ch

Jonas Buchli

Agile & Dexterous Robotics Lab  
ETH Zürich, Switzerland  
Email: buchlij@ethz.ch

**Abstract**—Many robotic tasks rely on the accurate localization of moving objects within a given workspace. This information about the objects' poses and velocities are used for control, motion planning, navigation, interaction with the environment or verification. Often motion capture systems are used to obtain such a state estimate. However, these systems are often costly, limited in workspace size and not suitable for outdoor usage. Therefore, we propose a lightweight and easy to use, visual-inertial Simultaneous Localization and Mapping approach that leverages cost-efficient, paper printable artificial landmarks, so called fiducials. Results show that by fusing visual and inertial data, the system provides accurate estimates and is robust against fast motions and changing lighting conditions. Tight integration of the estimation of sensor and fiducial pose as well as extrinsics ensures accuracy, map consistency and avoids the requirement for precalibration. By providing an open source implementation and various datasets, partially with ground truth information, we enable community members to run, test, modify and extend the system either using these datasets or directly running the system on their own robotic setups.

## I. INTRODUCTION

### A. Motivation

For many tasks in mobile robotics, it is important to estimate a robot's state with respect to its workspace, i.e. its pose and velocities expressed in an inertial coordinate system aligned to the robot's workspace. Such tasks include navigation, motion planning or manipulation. One way to measure the position and orientation of a robot is to use a motion capture system (such as e.g. Vicon, Optitrack or PTI Visualeyex). These systems are usually highly accurate and provide pose estimates w.r.t. to a calibrated reference system. However, these systems can be very costly and the user might be limited to a certain workspace size. Furthermore, many common systems such as Vicon and Optitrack use passive markers that reflect infrared light which limits their usage to indoor setups. Also, most systems require a tedious calibration procedure that needs to be repeated frequently to maintain accuracy. Since these systems do not have access to inertial data, they have to rely on finite differences of position measurements to estimate velocity information, which usually leads to highly quantized (compare Figure 9) or delayed data.

Another approach to state estimation is Visual Odometry (VO), which is sometimes also fused with inertial data. VO can provide very accurate local estimates of the robot motion. However, usually only a finite number of previous observations (frames) are included during the pose estimation step and

no loop closure is performed. Thus, VO is prone to drift over time and does not provide a globally consistent path. Additionally, VO only provides a pose estimate to the initial pose and not to the workspace. Compared to VO, Simultaneous Localization and Mapping (SLAM) introduces the notion of a global map and therefore can ensure consistency by performing loop closure. However, map building, storing and loop-closure detection can be computationally and memory demanding.

In this work, we propose a lightweight, cost effective motion capture system based on a monocular, visual-inertial SLAM system that tightly fuses inertial measurements and observations of artificial visual landmarks, also known as "fiducials" which constitute the map. By using artificial landmarks that provide rich information, the estimation, mapping and loop closure effort is minimized. In this implementation, we use AprilTags [1] as our fiducials. Since these tags also provide a unique identification number, they can be robustly tracked and estimated in the applied Extended Kalman Filter (EKF). Additionally, loop closure is handled implicitly and no additional loop closure detection step is required. A single observation of a tag is sufficient to estimate the relative transformation between tag and robot.

Most tag based localization systems use the relative pose estimates from the tag observations. In this work, we chose a tightly coupled approach where the *corner detections* are used as observations, forming a holistic sensor fusion algorithm. Hence, the system can work with very few tags and observations while still remaining accurate. This reduces the map size of the estimation problem and lowers computational demands. Therefore, the complexity of the proposed system is much lower than common SLAM approaches while still providing accurate, globally consistent estimates. By also including inertial measurements, robust performance during fiducial occlusion, motion blur from fast motions and changing lighting conditions is ensured. The approach requires to artificially prepare the workspace but also provides relative pose information within the workspace. Hence, instead of an alternative for SLAM solutions, we see the developed system as a lightweight tool that can be used as an inexpensive outdoor-capable motion capture system, for verification of other state estimation systems or for absolute localization in a given workspace.

## B. Related Work

Many existing fiducial-based localization systems are targeted at augmented reality or were designed to be used with cameras only. Hence, many systems use vision data only (e.g. [2], [3], [1], [4], [5]). These systems have two major drawbacks over the presented system. Firstly, they fail to provide any estimate during occlusion or motion blur. Secondly, linear velocities and body rates can only be computed based on the position and orientation estimates and are thus highly quantized. While this might be negligible for virtual reality applications, it can cause issues when closing a control loop over these estimates. To mitigate these issues, a motion model can be assumed [6]. However, this makes the approach specific to the implemented motion model.

The motion estimation and map building elements of the presented approach are closely related to monocular, visual-inertial SLAM, which has proven to be very effective [7], [8], [9], [10], [11]. The difference between our approach and fiducial-free visual-inertial SLAM solutions is that we tightly integrate artificial landmarks that result in highly robust and unique features in image space. As a result, our landmarks can be robustly re-detected and their detections are almost outlier free which increases the robustness of the approach. Additionally, each landmark has a 6-DoF pose (position and orientation) rather than only 3-DoF as in commonly used point landmarks. Furthermore, since a single measurements fully constrains the 6-DoF relative pose, a single landmark is sufficient for estimating the pose. This also allows for a simple yet accurate landmark initialization which usually is a problem in monocular SLAM approaches [12]. Furthermore, pose landmarks allow for aligning the map and the estimates to a given frame in the workspace and thus, the system can provide an absolute localization in the workspace which can be crucial for tasks that assume a prepared environment.

While both, fiducial-based localization and SLAM are well studied problems, not many approaches exist that combine both. One approach where fiducials are combined with SLAM is presented in [13]. However, the inertial measurements are not used to estimate velocities but only used for a fall-back pose estimation if all fiducials are occluded. Another similar system as the one presented in this work has been described in [14]. Since this work is part of the development of a commercial product (InterSense IS-1200) the authors remain relatively vague about their sensor fusion algorithm as well as the achievable performance of their system. Furthermore, dedicated hardware is required which poses additional costs for the user and contradicts the goals of this project to provide a cost-efficient, open source framework. A third visual-inertial, fiducial based localization system is presented in [15]. While also here inertial measurements and visual data are fused in an EKF, the approach does not include measurements from an accelerometer which can be helpful during fast linear motions and can provide a notion of gravity. Additionally, it is assumed that the poses of the tags are perfectly known in a workspace frame. Therefore, one can only place the tags in

known configuration and imperfect calibration will lead to an inconsistent map. In [16] a fiducial-based SLAM approach is presented. However, here the fiducials are only represented as point features and thus only the 3D positions of the fiducials are estimated.

## C. Contributions

We present a lightweight motion estimation system based on monocular visual-inertial EKF-SLAM using artificial landmarks. This work tightly couples two proven concepts, SLAM and fiducial-based localization, by using corner observations of 6 DoF landmarks and adapting the corresponding Kalman filter innovation term. This allows for smaller map sizes and leaner estimation. In contrast to existing approaches relying on 6 DoF fiducials, the presented framework processes visual measurements and inertial data in a single estimator which results in consistent data and avoids precalibration and recalibration. We have developed this tool out of a need for an open source, lightweight, accurate visual-inertial motion capture system. We provide the system as free to use open source software. Since it only relies on standard hardware (an IMU and a camera) and a Robot Operating System (ROS) software interface, both often available on robotic platforms, it can be deployed easily. The source code, the datasets as well as a more detailed technical manual can be found at <https://bitbucket.org/adrlab/rcars>, allowing easy integration and full reproducibility of the presented results.

## D. Notation and Conventions

In the following sections, scalars are indicated with small letters (e.g.  $f_x$ ). Vectors are indicated with small, bold letters (e.g.  $\mathbf{r}$ ). Matrices are indicated by non-bold capital letters (e.g.  $K$ ). A capital subscript leading a variable name describes the coordinate frame that the quantity is expressed in. Position vectors are denoted by  $\mathbf{r}$ . The trailing subscript describes the direction of the vector from its origin to its goal position (read from left to right), e.g.  ${}^A\mathbf{r}_{QP}$  is a position vector expressed in frame  $A$  that points from point  $Q$  to point  $P$ . Quaternions are denoted by  $\mathbf{q}$ . The trailing subscript denotes the coordinate systems involved in the passive rotation, e.g.  $\mathbf{q}_{AB}$  represents the passive rotation from the coordinate system  $B$  to the coordinate system  $A$ . Hence, to rotate a position vector expressed in  $B$  to  $A$ , we would compute  ${}^A\mathbf{r}_{QP} = \mathbf{q}_{AB}({}^B\mathbf{r}_{QP})$ .

## II. SYSTEM DESCRIPTION

The present localization system consists of two main components, a detector for the fiducials and an EKF for sensor fusion. In a first step, the image acquired by the camera is undistorted. Afterwards, the detector is run on the image which outputs the corner coordinates as well as a unique identifier number (id) associated with each detected tag. Furthermore, it estimates the relative transformation between each tag and the camera. This estimation is based on an iterative optimization minimizing the reprojection errors between the projected 3D corner points and their detections in image space. In a second step, the EKF uses the information from re-detected tag

corners to estimate the robot’s state, including pose, linear velocity and body rates. Additionally, the filter continuously estimates the position and orientation of the tags with respect to the camera coordinate frame. When a tag is seen for the first time, its pose is initialized using the relative transformation between the camera and the tag as provided by the detector. After this initialization, the tag pose will be refined within the EKF by using the reprojection errors of its corners in each subsequent re-observation. To ensure consistency, the extrinsic calibration between camera and IMU as well as the additive IMU biases are also included in the filter state.

#### A. Fiducials

Over the past years, a large variety of fiducial systems have been developed. Very popular implementations include ARToolKit [17], ARTag [3], CyberCode [18] and multiring color fiducials [19]. In our implementation, we use AprilTags [1] which are 2-dimensional, printable bar codes. The reason for this choice was the achievable high accuracy [1] and the numerous available detector implementations in C/C++. In our system, we use the detector implemented in cv2cg<sup>1</sup>. In our evaluations, this implementation has proven to be fast and providing accurate and robust tag detections.

#### B. Hardware

The proposed system requires a camera and an IMU. While the transformation between IMU and camera can be estimated online, the camera intrinsics have to be given a-priori. In our setup, we are using a Skybotix VI-Sensor [20]. This sensor consists of two cameras in a stereo configuration and an IMU. While the sensor is a stereo camera we are only relying on the left camera in this work. The sensor is set up to output images at 20 Hz and IMU data at 200 Hz.

#### C. Camera Model

In this project, we assume a pinhole camera model which is applicable to most cameras with common field-of-views. The expected input for the detector and filter is an undistorted image. Therefore, the user is free to choose a distortion model as long as an undistorted image is provided. In the case of the VI-Sensor we are using a radial tangential distortion model. The pinhole camera model is represented by the overall projection  $\pi$  which depends on the camera intrinsics, i.e., the focal lengths,  $f_x$  and  $f_y$ , and the camera’s principle point  $\mathbf{c} = (c_x, c_y)$ . It maps a 3D point  $P$  expressed in the camera coordinate frame,  ${}_{V}r_{VP}$ , to its corresponding pixel coordinates  $\mathbf{p} = \pi({}_{V}r_{VP})$ .

#### D. Filter

In order to fuse the information gained from the observed tags together with the on-board inertial measurement we implement an extended Kalman filter. Relying on appropriate sensor models, this filter uses the inertial measurements in order to propagate the robot’s state and performs an update step based on the available tag corner measurements. In the

following paragraphs we will explain the sensor models used and derive the required filter equations. For readability, this derivation is carried out for the case of a single tag, but is directly applicable to the case of multiple tags.

1) *Coordinate Systems*: In our filter setup we assume different coordinate frames. First, we assume an inertial workspace coordinate system  $W$ . We assume that gravity points in negative z-direction in this frame. Furthermore, we define the IMU coordinate system  $B$  and the camera frame  $V$ . Finally, we define a coordinate system  $T$  for each tag which coincides with the geometrical center of the tag and where z is perpendicular to the tag plane.

2) *Sensor Models*: First, we introduce the sensor model used for the IMU. It assumes Gaussian noise as well as additive bias terms for accelerometer and gyroscope measurements. This can be formulated as follows:

$$\tilde{\mathbf{f}} = \mathbf{f} + \mathbf{b}_f + \mathbf{w}_f, \quad (1)$$

$$\dot{\mathbf{b}}_f = \mathbf{w}_{bf}, \quad (2)$$

$$\tilde{\boldsymbol{\omega}} = \boldsymbol{\omega} + \mathbf{b}_\omega + \mathbf{w}_\omega, \quad (3)$$

$$\dot{\mathbf{b}}_\omega = \mathbf{w}_{b\omega}, \quad (4)$$

where  $\tilde{\mathbf{f}}$  and  $\tilde{\boldsymbol{\omega}}$  are the actual measurements of the proper acceleration and rotational rates,  $\mathbf{b}_f$  and  $\mathbf{b}_\omega$  are the additive bias terms, and all terms of the form  $\mathbf{w}_*$  represent continuous white Gaussian noise processes.

In addition to the IMU data, we will also include measurements related to the observed tags. For this measurements we propose a tight coupling by using a corner reprojection based visual model. Given the relative position and attitude of a specific tag with respect to the camera frame,  ${}_{V}r_{VT}$  and  ${}_{TV}q$ , we can compute the position of the  $i^{\text{th}}$  tag corner  ${}_{T}r_{TC_i}$  (fixed to the tag coordinate frame  $T$ ) as viewed from the camera:

$${}_{V}r_{VC_i} = {}_{V}r_{VT} + {}_{TV}q^{-1}({}_{T}r_{TC_i}). \quad (5)$$

By using the camera projection map  $\pi$ , we can project the above quantity onto the image plane and derive the corresponding pixel coordinate measurement  $\tilde{\mathbf{p}}_i$ , where we assume an additive Gaussian noise model ( $\mathbf{n}_{p,i} \sim \mathcal{N}(0, \mathbf{R}_p)$ ):

$$\tilde{\mathbf{p}}_i = \pi({}_{V}r_{VC_i}) + \mathbf{n}_{p,i}. \quad (6)$$

The advantage of the selected visual measurement model is that the noise is modelled directly on the pixel location of the detected corners. While the detector provides an estimation of the tag pose relative to the current camera frame and this could be directly used within the filter, fitting an accurate noise model to this relative pose would have been difficult since the magnitude of the noise strongly depends on the current location and orientation of the tag in the camera frame. In contrast, the noise on the reprojected tag corners is, to a large extent, indifferent with respect to the camera pose and can thus be assumed to be constant and identical for all tags and measurements.

<sup>1</sup><http://code.google.com/p/cv2cg/>

3) *Filter States and Prediction Model*: The above visual sensor model assumes the knowledge of the tag pose. Instead of using fixed values, which could quickly lead to inconsistencies, we propose to include the pose of the tag into the filter state. Therefore, the filter will be able to refine the tag pose and ensure map consistency. Employing a robocentric representation of the sensor state and the tag pose, we get the following filter state:

$$\mathbf{x} := (\mathbf{r}, \mathbf{v}, \mathbf{q}, \mathbf{b}_f, \mathbf{b}_\omega, \mathbf{r}_T, \mathbf{q}_T, \mathbf{r}_V, \mathbf{q}_V) \quad (7)$$

$$:= ({}^B\mathbf{r}_{WB}, {}^B\mathbf{v}_B, \mathbf{q}_{WB}, {}^B\mathbf{b}_f, {}^B\mathbf{b}_\omega, \\ {}^V\mathbf{r}_{VT}, \mathbf{q}_{TV}, {}^B\mathbf{r}_{BV}, \mathbf{q}_{VB}). \quad (8)$$

In the above state,  $\mathbf{r}$ ,  $\mathbf{v}$ , and  $\mathbf{q}$  are the robocentric position, velocity, and attitude of the sensor. Furthermore,  $\mathbf{r}_T$  and  $\mathbf{q}_T$  are used for parametrizing the pose of the tag, while  $\mathbf{r}_V$  and  $\mathbf{q}_V$  represent the extrinsic calibration between IMU and camera. Computing the total derivatives of the selected state and inserting the IMU model (1)-(4) yields:

$$\dot{\mathbf{r}} = -(\tilde{\omega} - \mathbf{b}_\omega - \mathbf{w}_\omega)^\times \mathbf{r} + \mathbf{v} + \mathbf{w}_r, \quad (9)$$

$$\dot{\mathbf{v}} = -(\tilde{\omega} - \mathbf{b}_\omega - \mathbf{w}_\omega)^\times \mathbf{v} \\ + \tilde{\mathbf{f}} - \mathbf{b}_f - \mathbf{w}_f + \mathbf{q}^{-1}(\mathbf{g}), \quad (10)$$

$$\dot{\mathbf{q}} = -\mathbf{q}(\tilde{\omega} - \mathbf{b}_\omega - \mathbf{w}_\omega), \quad (11)$$

$$\dot{\mathbf{b}}_f = \mathbf{w}_{bf}, \quad \dot{\mathbf{b}}_\omega = \mathbf{w}_{b\omega}, \quad (12)$$

$$\dot{\mathbf{r}}_T = -\mathbf{q}_V((\tilde{\omega} - \mathbf{b}_\omega - \mathbf{w}_\omega)^\times (\mathbf{q}_V^{-1}(\mathbf{r}_T) + \mathbf{r}_V) + \mathbf{v}) \\ + \mathbf{w}_{rt}, \quad (13)$$

$$\dot{\mathbf{q}}_T = -(\mathbf{q}_T \otimes \mathbf{q}_V)(\tilde{\omega} - \mathbf{b}_\omega - \mathbf{w}_\omega + \mathbf{w}_{qt}), \quad (14)$$

$$\dot{\mathbf{r}}_V = \mathbf{w}_{rv}, \quad \dot{\mathbf{q}}_V = \mathbf{w}_{qv}. \quad (15)$$

We include additional continuous white Gaussian noise processes  $\mathbf{w}_r$ ,  $\mathbf{w}_{rt}$ ,  $\mathbf{w}_{qt}$ ,  $\mathbf{w}_{rv}$ , and  $\mathbf{w}_{qv}$  in order to excite the full filter state and for modeling errors caused by the subsequent discretization of the states. For all white Gaussian noise processes  $\mathbf{w}_*$ , the corresponding covariance parameters  $\mathbf{R}_*$  describe the magnitude of the noise. While most covariance parameters can be chosen by considering the corresponding sensor specifications, some remain as tuning parameters. Using a simple Euler forward integration scheme a set of discrete time prediction equations can be derived. In order to achieve a minimal and consistent parametrization, the derivatives of the quaternions are expressed in a 3D local angular velocity. This has to be considered during the discretization and when implementing the filter. Also note that, during the prediction, the IMU-related states ( $\mathbf{v}$ ,  $\mathbf{b}_\omega$ ) are coupled to the estimated tag pose ( $\mathbf{r}_T$ ,  $\mathbf{q}_T$ ) based on the IMU-camera extrinsics estimates ( $\mathbf{r}_V$ ,  $\mathbf{q}_V$ ).

4) *Update Model*: The update step is performed by directly employing the reprojection error as the Kalman filter innovation term. For each tag corner  $i$  and based on equation (6) we can define an innovation term  $\mathbf{y}_i$ :

$$\mathbf{y}_i = \tilde{\mathbf{p}}_i - \pi(\mathbf{v}\mathbf{r}\mathbf{V}\mathbf{C}_i). \quad (16)$$

This results in a 8D innovation term for every tag detected in the current camera frame (2D per tag corner). For each newly

observed tag the state is augmented by an additional tag pose, i.e. position and attitude. The augmentation uses the estimated relative pose from the tag tracker in order to initialize the state with a good linearization point. The corresponding covariance matrices are initialized to large values and typically converge very quickly. Optionally, tags with known absolute location can also be fed to the filter. Especially for datasets with a large number of tags this comes in handy since the EKF does not scale well with increasing state dimension. Above around 20 tag poses in the filter state the prediction step becomes very costly for a single core implementation.

### III. RESULTS

In order to assess the performance of the proposed approach, we define different test procedures. In a first test, we verify the accuracy of the fiducial pose estimation. In a second test, we then evaluate the accuracy of the motion estimation computed by our EKF. Both tests are verified with ground truth data obtained from a high class external motion capture system. Additionally, two large scale datasets are processed for verifying the accuracy when closing larger loops. A third test evaluates the estimation of the extrinsic calibration. Lastly, we test the applicability of the presented motion estimation within an online closed loop control on a quadruped robot.

#### A. Datasets

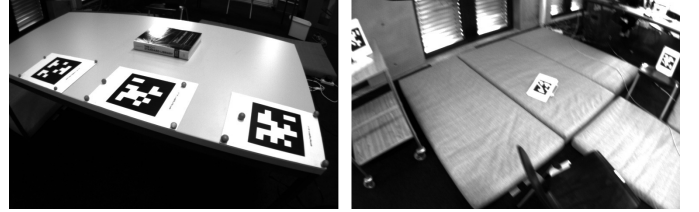


Fig. 1. Images extracted from datasets "table" (left) and "dataset\_1" (right).

In total, we are using five datasets which are all available for download with the source code. The first dataset "table" consists of three tags that are placed flat on a table at the same orientation as shown in Figure 1. The distances between the tags are chosen to be of similar magnitude. The second dataset "dataset\_1" also contains three tags. This time, we tried to create a challenging dataset, where the tags are sparsely distributed around a larger workspace of about 4x4x4m. Furthermore, the tags are intentionally oriented and located in such a way that the viewing angle for the camera is not ideal and that only a minimal amount of frames contain two neighboring tags at the same time. Additionally, the sensor is moved fast, such that motion blur occurs occasionally. Overall, this increases the level of difficulty in estimating the tags' locations. An on-board image taken by the camera, showing the challenging setup as well as the motion blur is shown in Figure 1. For evaluating the extrinsic calibration estimation, we are using a set of datasets all taken within the same workspace to ensure a consistent setup. In these datasets, the sensor is subject to extensive motion in order

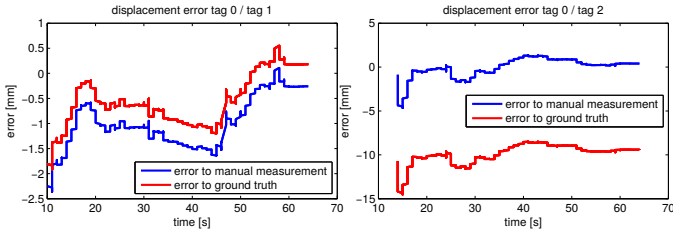


Fig. 2. Displacement error between estimated tag positions and reference from manual measurements as well as external motion capture. As can be seen, the error decreases over time, since the tags’ positions are iteratively refined by the EKF. Finally, submillimeter accuracy is achieved. The larger offset on the right plot most likely results from inaccurate marker and coordinate system placement in the external motion capture system.

to properly excite the full filter state and thereby promote the convergence of the estimated extrinsic calibration. The last two dataset ”cube“ and ”pavillon“ contain round trips on our campus. Both datasets include around 35 tags and span areas of approximately  $25 \times 25 \times 6$  m. By moving from basement to ground level or from indoors to outdoors, these datasets are subject to changing lighting conditions. To provide comparable results, no test specific parameter tuning has been performed, i.e. the same parameters are used throughout all tests.

### B. Fiducial Estimation Test

For the verification of the continuous fiducial estimation procedure, we use the ”table“ dataset. In this dataset, manually measuring the offset between the tags is simple. Thus, we can use these measurements as ground truth information and compare it to the estimates of the external motion capture system. This allows us also to evaluate the accuracy of the marker and coordinate system placement during the set up of the external motion capture system. To isolate the fiducial estimation for testing, we are disabling the extrinsic calibration in this test and use the sensor’s factory calibration.

We compare the norm of the relative translation, i.e. the distance between tag 0 and tag 1 as well as between tag 0 and tag 2 with the manual distance measurements. This error plots are shown in Figure 2. The plots shows two interesting aspects. The errors in the beginning of the sequence is quite small. This indicates that the initial guess obtained from the reprojection error optimization on the first frame is fairly accurate. Over time, our EKF then further refines the poses, reaching approximately millimeter accuracy which is of equal magnitude as manual measuring errors.

The figures also shows the error with respect to the external motion capture measurement. Here, the error is shifted by about 1cm for the translation between tag 0 and tag 2. Since a zero mean error curve would be expected, this constant offset most likely results from inaccurate marker and reference coordinate system placement. Since the tags in this dataset are placed flat on a table and aligned with the table’s edge, we can also analyze the rotation error of our tag pose estimates. To do so, we compute the relative rotation between two tags. We then convert this rotation to an axis-angle representation and use the angle as our error measurement. Due to the tag alignment, the relative rotation between two tags can be assumed to be

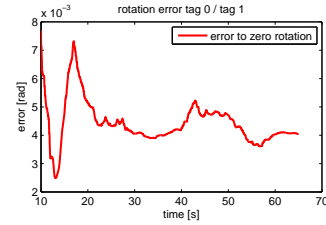


Fig. 3. Rotation error between estimated tag positions and zero rotation. The error is obtained by converting the relative rotation to angle-axis representation of which the angle is plotted. As can be seen, the error decreases over time, since the tags’ rotations are iteratively refined by the EKF. The error starts at around 0.5 degrees and reduce to about 0.2 degrees.

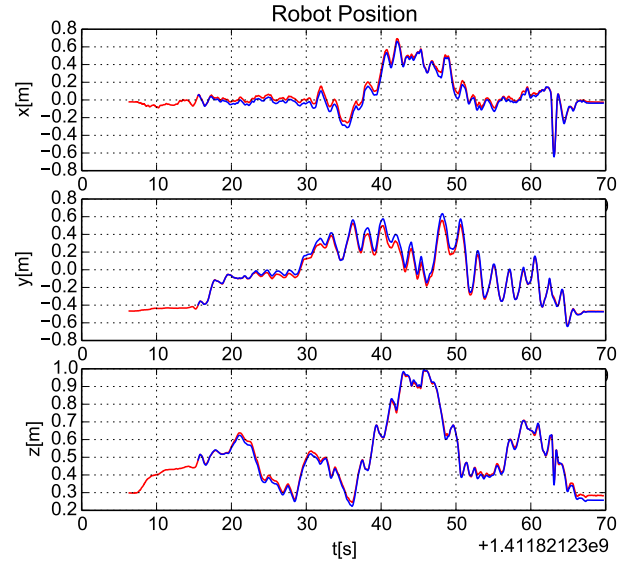


Fig. 4. Comparison between estimated robot position (blue) and ground truth position (red) for the dataset ”table“. As can be seen, the maximum position offset between both measurements lie only within a centimeter scale which is the same magnitude as the achievable measurement accuracy in this setup.

identity. This is also confirmed by the external motion capture system up to the fourth decimal of the relative rotation angle. Figure 3 shows the error between estimated rotation and the identity rotation for the relative rotation between tag 0 and tag 1. As can be seen, the error is initially around 0.5 degrees. Through continuous refinement of the tag poses within the EKF, this error reduces to around 0.2 degrees over time. This error is of same magnitude as printing and measurement accuracy.

The experiments described above show the high achievable fiducial estimation accuracy in translation and rotation. Furthermore, these results underline that tag pose refinement significantly reduces displacement and rotational errors present in the single frame pose estimate used for initialization. This will eventually improve the consistency of the relative tag poses and thus should also improve the robot’s pose estimation.

### C. Motion Estimation Test

1) *Dataset Table*: To assess the performance of the motion estimation, we use both datasets described above. The goal of our estimation framework is to localize against our workspace, where we choose tag 1 as the origin. This choice is

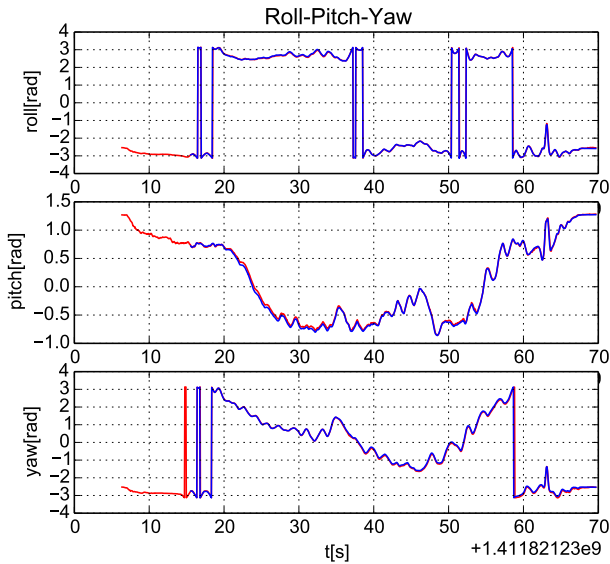


Fig. 5. Comparison between estimated robot orientation (blue) and ground truth orientation (red) for the dataset "table". Due to the wrap-around at  $\pm\pi$  the plot is discontinuous. However, since quaternions are used for the internal representation of the filter, the output of the filter is smooth. As also seen in the position data, estimated and ground truth rotations agree up to measurement uncertainty.

arbitrary and one could choose any tag as a reference defining the workspace location and orientation. Since our estimator automatically estimates the orientation of the workspace with respect to gravity, no manual alignment is required. Figure 4 shows a comparison of the position estimates of the filter and ground truth data from the external motion capture system for the *table* dataset. This plot nicely illustrates the robust tracking behavior of the system. Even though the reference tag is not detectable at every instance of the dataset, the estimated fiducials provide a stable reference for the filter to localize against, such that tracking errors remain a few centimeters. In Figure 5 the estimated orientation and ground truth orientation for the same datasets are compared. Also here, the estimator shows a robust tracking with minimal deviations. The maximum error observed in pitch direction is about 0.05 rad which corresponds to less than 3 degrees. Since the ground truth reference data is a relative pose between the sensor and the reference tag computed from the individual poses, the error magnitudes observed above lie within the measurement accuracy of the ground truth data. While this underlines the performance of the approach, it does not give any indication about its limits. Therefore, we tried to push the system to its limits using *dataset\_1* which contains several artificial challenges as described above.

2) *Dataset Dataset\_1*: In this experiment, again the estimated position is compared to ground truth data and the results are shown in Figure 6. As the plots show, the position starts to deviate from ground truth in the last third of the sequence. While results are not as good as in the *table* dataset, *dataset\_1* can be seen as a worst-case benchmark scenario. Most of the difficulties for the algorithm are artificially posed and the performed motion is faster than in many robot applications.

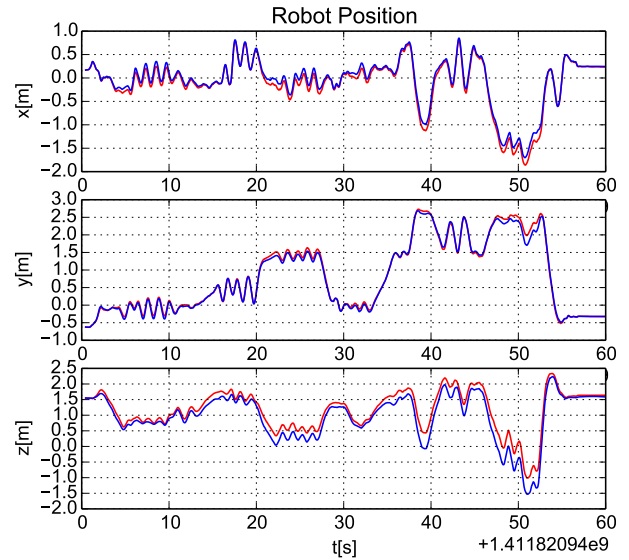


Fig. 6. Comparison between estimated robot position (blue) and ground truth position (red) for the dataset "dataset\_1". This dataset has been made artificially difficult with sparse tag coverage and fast motions to show the robustness of the filter. While the estimates diverges when only the briefly observed tag on the very left can be used for localization, it converges back to the ground truth information when localizing against the other tags again.

Due to the sparse tag placement and fast motions, the detector was unable to detect any tag in many of the images of the sequence. This is shown in Figure 8 where these instances are marked with the value 1. In total, the filter is provided with inertial measurements only for almost 20% of the sequence. Additionally, tag 0 is only seen together with another tag for in total 9 frames. Thus, little localization information is provided for this tag, leading to a high uncertainty of the tags pose. Still, it is the only visible tag for about 15% of the dataset. Thus, the filter is only provided with uncertain vision information and noisy inertial measurements during these parts. However, the filter remains stable and is able to converge close to ground truth data again when the other tags are visible again.

Also in the orientation, the effects of sparsely distributed tags combined with fast motions are visible. Figure 7 shows the difference between ground truth and estimated orientation for *dataset\_1*. As can be seen, the orientation estimate is fairly accurate throughout the dataset with a slight deviation in yaw at the beginning of the trajectory and a small deviation of pitch of about 9 degrees towards the end. When looking at the linear velocity estimates for this dataset shown in Figure 10, one can see that the estimates agree well with the velocity data obtained by using finite differences on the ground truth data. Interestingly, the estimated velocities are virtually outlier free while the finite differences show occasional peaks. This effect still occurs, even though a high quality motion capture system has been used. This underlines the limitations of using finite differences for velocity estimates and encourages the use of inertial data. This effect is even more pronounced when looking at Figure 9 which shows the rotational velocity estimates and their counterparts computed using finite differences on the ground truth orientation. The difference in noise level

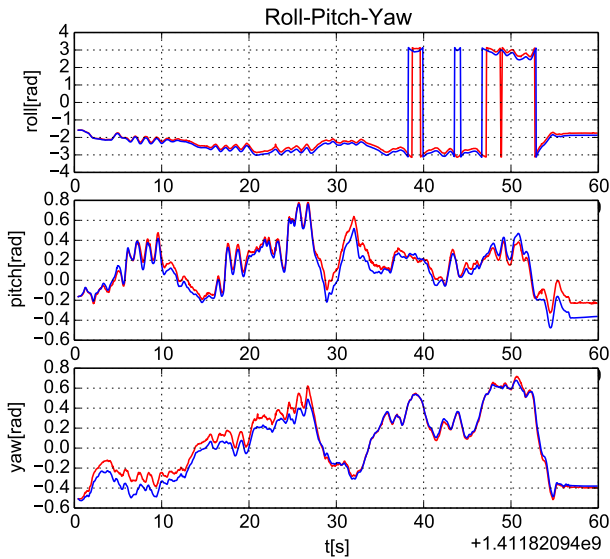


Fig. 7. Comparison between estimated robot orientation (blue) and ground truth orientation (red) for the dataset "orientation\_1". Due to the wrap-around at  $\pm\pi$  the plot is discontinuous. However, since quaternions are used for the internal representation of the filter, the output of the filter is smooth.

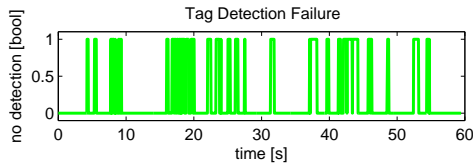


Fig. 8. Plot indicating whether one or multiple tags were detected (indicated as 0) or no tag was detected (indicated as 1) for *dataset\_1*. Overall, in almost 20% of all images no tag could be detected.

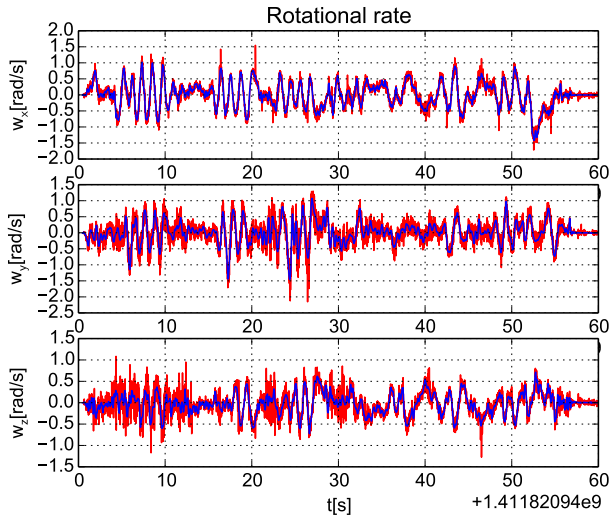


Fig. 9. Comparison of rotational velocity estimates (blue) and rotational velocities calculated by using finite differences of the ground truth orientation data (red). Like also with the linear velocities shown in Figure 10, the estimation matches the ground truth data. Here the significance of using inertial measurements for low-noise estimates over finite differences on pose information is even more prominent.

between both measurements is significant. One reason is that the IMU directly measures rotational rates using gyroscopes. Furthermore, rotations tend to be more difficult to estimate for external motion capture systems. This effect gets amplified when differentiating this noisy signal.

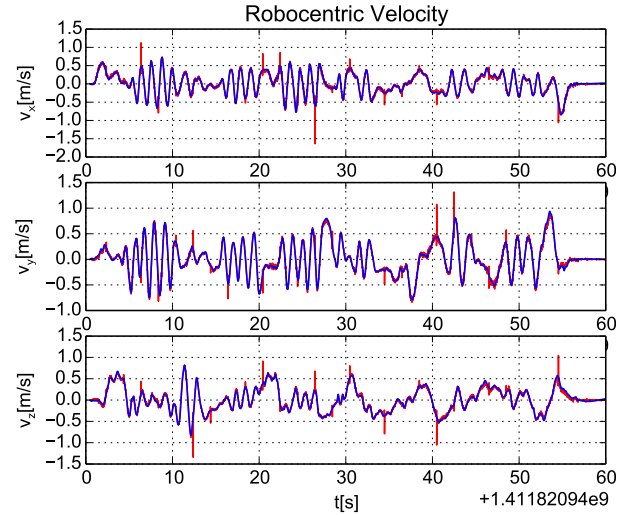


Fig. 10. Comparison of linear velocity estimates (blue) and linear velocities calculated by using finite differences of the ground truth position data (red). While the estimated velocities agree well with the velocities computed from ground truth data, they are virtually outlier free. While a high quality external motion capture system is used, this data still shows the limitations of finite differences for velocity estimates.

3) *Large scale datasets*: One advantage of the presented approach over a commercial motion capture system is the workspace size. Since our system only relies on paper printed tags, a large workspace can be covered easily. As we did not have a motion capture system available that is capable of covering such a large area, especially not outdoors, we are using loop closure to estimate the accuracy of the approach. For this test, we are using the datasets "pavillon" and "cube" which both include loop closure sequences. As a quality measure, the reprojection errors as well as the offset between the estimator's predicted tag position and the detector's instantaneous tag measurement are used. Both measures are taken at the first time that we reobserve a tag after the round trip, before updating the estimator.

For the dataset "cube" the average reprojection error at loop closure is 56.07 pixels. Taking the detector pose estimate as a reference, the position offset is 0.86 m. One round trip until loop closure is about 70 m long and follows a trail of 36 tags. Therefore, the relative position error is around 1.2 %. For the dataset "pavillon" the average reprojection error at loop closure is 51.01 pixels. Taking the detector pose estimate as a reference, the position offset is 0.38 m. One round trip until loop closure is about 80 m and follows a trail of 33 tags. Therefore, the relative position error is around 0.5 %. Please note that position errors are calculated using the detector estimate. This estimate cannot be assumed to be a ground truth measurement. Therefore, the accuracy figures above are

subject to measurement inaccuracies of the detector and based on one measurement only.

#### D. Online Extrinsic Estimation

In order to assess the quality of the estimated camera-IMU extrinsics, we evaluated the corresponding values after the system was sufficiently excited such that the values could converge. Since no real groundtruth references were available for the extrinsics, we evaluated the repeatability of the obtained estimates. For this we recorded 10 datasets within the same environment while performing similar motions with a total duration around 50-60 seconds. The obtained RMS-values were 1.5 cm for the translational part and 0.0035 rad for the rotational part of the extrinsics. Both values ranging near what can typically be obtained through a dedicated calibration routine.

#### E. Estimation in Closed Loop Control

Motion capture systems are increasingly used in closed-loop control. Since latency, noise and outliers can significantly deteriorate the closed-loop behavior of the plant, estimation in the loop is a challenging task. Therefore, we test our system in such an application. For this test, we are using our quadruped robot HyQ on a hydraulically actuated treadmill. The control task is to keep the robot in the center of the treadmill by only regulating the speed of the treadmill, i.e. the robot's walking motion is assumed to be a disturbance. The control system is a cascade of an inner velocity and an outer position control loop. The sensor input to the position control loop are the robot's position and velocity in the workspace.

As can be seen from the accompanying video<sup>2</sup>, the closed loop system is able to stabilize the robot's position on the treadmill while changing the walking speed. The estimate of the absolute position of the robot in the workspace allows us to move the robot to the treadmill center during initialization.

### IV. CONCLUSION AND FUTURE WORK

In this paper, we have presented an open-source, visual-inertial state estimation system, that tightly integrates monocular SLAM and fiducial based estimation. By relying on standard hardware already present on most robots the system can be applied cost efficiently. Experiments demonstrate its good accuracy and high robustness, which indicates that it could replace an expensive motion capture systems in applications that do not require sub-millimeter precision or very fast update rates only offered by highly expensive motion capture systems. This has been verified by using the system in a closed-loop control task. Large scale tests have demonstrated long term accuracy, map consistency and loop closure refinement. Experiments under fast motions and sparse tag coverage of the workspace underline the importance of including inertial measurements compared to fiducial only approaches. Furthermore, the inertial measurements ensure high quality translational and rotational velocity estimates which can outperform these of a commercial system. Results have shown that good coverage

of fiducials is important for a good estimation quality. In the future, we aim at supporting differently sized fiducials such that their size can be optimized for their intended location. Furthermore, we will investigate the combination with natural landmarks in order to further improve the estimation accuracy.

### ACKNOWLEDGEMENT

The authors would like to thank Sammy Omari, the Autonomous Systems Lab and Skybotix for their support with the external motion capture system and the VI-Sensor. Furthermore, the authors would like to thank Manuel Lussi for the support with the treadmill experiments. This research has been funded partially through a Swiss National Science Foundation Professorship award to Jonas Buchli.

### REFERENCES

- [1] E. Olson, "Apriltag: A robust and flexible visual fiducial system," in *IEEE International Conference on Robotics and Automation*, 2011.
- [2] S. Zickler, T. Laue, O. Birbach, M. Wongphati, and M. Veloso, "Ssl-vision: The shared vision system for the robocup small size league," in *RoboCup 2009: Robot Soccer World Cup XIII*. Springer, 2010.
- [3] M. Fiala, "Artag, a fiducial marker system using digital techniques," in *Conference on Computer Vision and Pattern Recognition*, 2005.
- [4] M. Faessler, E. Mueggler, K. Schwabe, and D. Scaramuzza, "A monocular pose estimation system based on infrared leds," in *IEEE International Conference on Robotics and Automation*, 2014.
- [5] A. Breitenmoser, L. Kneip, and R. Siegwart, "A monocular vision-based system for 6d relative robot localization," in *International Conference on Intelligent Robots and Systems*, 2011.
- [6] H. Lim and Y.-S. Lee, "Real-time single camera slam using fiducial markers," in *ICCVS-SICE, 2009*, Aug 2009.
- [7] A. Mourikis, S. Roumeliotis, *et al.*, "A multi-state constraint kalman filter for vision-aided inertial navigation," in *IEEE International Conference on Robotics and Automation*, 2007.
- [8] J. Kelly and G. S. Sukhatme, "Visual-Inertial Sensor Fusion: Localization, Mapping and Sensor-to-Sensor Self-calibration," *Int. Journal of Robotics Research*, vol. 30, 2011.
- [9] E. S. Jones and S. Soatto, "Visual-inertial navigation, mapping and localization: A scalable real-time causal approach," *The International Journal of Robotics Research*, vol. 30, no. 4, 2011.
- [10] G. Nützi, S. Weiss, D. Scaramuzza, and R. Siegwart, "Fusion of imu and vision for absolute scale estimation in monocular slam," *Journal of intelligent & robotic systems*, vol. 61, no. 1-4, 2011.
- [11] M. Li and A. I. Mourikis, "High-precision, consistent ekf-based visual-inertial odometry," *The International Journal of Robotics Research*, vol. 32, no. 6, 2013.
- [12] J. Sola, T. Vidal-Calleja, J. Civera, and J. M. M. Montiel, "Impact of landmark parametrization on monocular ekf-slam with points and lines," *International journal of computer vision*, vol. 97, no. 3, 2012.
- [13] M. Mairi, F. Ababsa, and M. Malle, "Vision-inertial tracking system for robust fiducials registration in augmented reality," in *Computational Intelligence for Multimedia Signal and Vision Processing, 2009. CIMSVP '09. IEEE Symposium on*, March 2009.
- [14] E. Foxlin and L. Naimark, "Vis-tracker: a wearable vision-inertial self-tracker," in *Virtual Reality, 2003. Proceedings. IEEE*, March 2003.
- [15] S. You and U. Neumann, "Fusion of vision and gyro tracking for robust augmented reality registration," in *IEEE Virtual Reality*, 2001.
- [16] M. Bryson and S. Sukkarieh, "Building a robust implementation of bearing-only inertial slam for a uav," *Journal of Field Robotics*, vol. 24, no. 1-2, pp. 113-143, 2007.
- [17] I. P. H. Kato, M. Billinghurst, and I. Poupyrev, "Artoolkit user manual, version 2.33," *Human Interface Technology Lab, University of Washington*, vol. 2, 2000.
- [18] J. Rekimoto and Y. Ayatsuka, "Cybercode: designing augmented reality environments with visual tags," in *Proceedings of DARE 2000 on Designing augmented reality environments*. ACM, 2000.
- [19] Y. Cho, J. Lee, and U. Neumann, "A multi-ring color fiducial system and an intensity-invariant detection method for scalable fiducial-tracking augmented reality," in *In IWAR*. Citeseer, 1998.
- [20] J. Nikolic, J. Rehder, M. Burri, P. Gohl, S. Leutenegger, P. T. Furgale, and R. Y. Siegwart, "A synchronized visual-inertial sensor system with fpga pre-processing for accurate real-time slam," in *IEEE International Conference on Robotics and Automation*, 2014.

<sup>2</sup><http://youtu.be/Ckf1QAuTKqc>