


# Iterated extended Kalman filter based visual-inertial odometry using direct photometric feedback

**Journal Article****Author(s):**

Bloesch, Michael; Burri, Michael; Omari, Sammy; [Hutter, Marco](#) ; Siegwart, Roland

**Publication date:**

2017-09

**Permanent link:**

<https://doi.org/10.3929/ethz-b-000187364>

**Rights / license:**

[In Copyright - Non-Commercial Use Permitted](#)

**Originally published in:**

The International Journal of Robotics Research 36(10), <https://doi.org/10.1177/0278364917728574>

# IEKF-based Visual-Inertial Odometry using Direct Photometric Feedback

Journal Title  
XX(X):1–18  
©The Author(s) 2017  
Reprints and permission:  
sagepub.co.uk/journalsPermissions.nav  
DOI: 10.1177/ToBeAssigned  
www.sagepub.com/



Michael Bloesch<sup>1</sup>, Michael Burri<sup>1</sup>, Sammy Omari<sup>1</sup>, Marco Hutter<sup>1</sup>, Roland Siegwart<sup>1</sup>

## Abstract

This paper presents a visual-inertial odometry framework which *tightly* fuses inertial measurements with visual data from one or more cameras, by means of an iterated extended Kalman filter (IEKF). By employing image patches as landmark descriptors, a photometric error is derived, which is directly integrated as an innovation term in the filter update step. Consequently, the data association is an inherent part of the estimation process and no additional feature extraction or matching processes are required. Furthermore, it enables the tracking of non-corner shaped features, such as lines, and thereby increases the set of possible landmarks. The filter state is formulated in a fully robocentric fashion, which reduces errors related to nonlinearities. This also includes partitioning of a landmark's location estimate into a bearing vector and distance and thereby allows an *undelayed* initialization of landmarks. Overall, this results in a compact approach which exhibits a high level of robustness with respect to low scene texture and motion blur. Furthermore, there is no time-consuming initialization procedure and pose estimates are available starting at the second image frame. We test the filter on different real datasets and compare it to other state-of-the-art visual-inertial frameworks. The experimental results show that robust localization with high accuracy can be achieved with this filter-based framework.

## Keywords

Visual-Inertial Odometry, Iterated Extended Kalman Filter, Photometric Error, Tight Information Fusion, Multiple Cameras

## 1 Introduction

Robust and high-bandwidth estimation of ego-motion is a key factor to enable the operation of autonomous robots. For dynamically controlled robots, such as aerial vehicles or legged robots, a reliable state estimate is essential: Failures of the state estimator can quickly lead to damage of the hardware and its surroundings. Thus, as autonomous robots become more capable and extend their range of applications, it is essential that the corresponding ego-motion estimation can perform well in increasingly difficult environments. The corresponding selection of sensors should be kept as lightweight and low-cost as possible in order to employ them on a wide range of robotic systems. Furthermore, in the context of vision-based estimation, extreme conditions such as strongly varying lighting, missing texture, fast motion, or dynamic objects may need to be accounted for.

Past research has shown that combining the complementary information from an Inertial Measurement Unit (IMU) and visual sensors can be a very capable approach in terms of accuracy and reliability. Consequently this approach has been successfully applied to robotic systems such as unmanned aerial robots (Weiss et al. (2013); Shen et al. (2014)) or legged robots (Stelzer et al. (2012); Ma et al. (2015)). Since assessing the precision of an algorithm is often simpler than evaluating its robustness, many researchers have focused on optimizing the accuracy of their approaches. The evaluation is typically done by measuring the accumulated position error over given traveled distances. Depending on the experimental setup, state-of-the-art algorithms reduce position errors to 0.1% of the traveled distance

(Leutenegger et al. (2015); Forster et al. (2016); Usenko et al. (2016)). Such a demonstration of high accuracy can serve as surrogate for the well-functioning of an approach. However, all odometry frameworks inherently suffer from drift and, if the primary goal is localization accuracy, a back-end framework doing global mapping, re-localization and loop closure will be indispensable (e.g. Lynen et al. (2015)). Furthermore, if the ego-motion estimation is employed within a feedback loop on an autonomous robot, other aspects like reliability and estimation time-delay become more central.

The well-established Kalman Filtering techniques represent sensor fusion frameworks that allow computationally efficient and high-bandwidth state estimation. Due to the inherent marginalization, the filter states at each timestep can refer to different physical quantities, e.g., a landmark's position can be estimated w.r.t. the moving sensor frame (and thereby represent a varying quantity over time). This enables the use of a fully robocentric formulation of the state and thereby reduces observability/nonlinearity related issues (Castellanos et al. (2004)). To mitigate the problem of intrinsic unobservability of a landmark's initial distance from the observer, the landmark position can be parameterized by its bearing vector and distance (Montiel et al. (2006)).

<sup>1</sup>ETH Zürich, Switzerland

### Corresponding author:

Michael Bloesch, Autonomous Systems Lab, ETH Zürich, LEE J 225, Leonhardstrasse 21, 8092 Zurich, Switzerland.

Email: bloeschm@ethz.ch

Consequently a landmark's distance can be initialized with a high uncertainty without affecting its bearing vector estimate (which can be initialized with a low uncertainty). Especially for scenarios with fast motions and short feature tracks, this becomes invaluable as it allows a seamless initialization of landmarks and thereby the extraction of visual information out of a landmark's second observation onwards. In this context, a sound representation of the filter state is crucial for the applicability of sensor fusion algorithms. An approach that has become increasingly popular for 3D orientations is based on a manifold encapsulation technique (Hertzberg et al. (2011)), which in the present work is also applied to bearing vectors.

The proposed approach combines an iterated extended Kalman filter (IEKF), a fully robocentric formulation of visual-inertial odometry, and a photometric error model. This is achieved by associating every landmark with a multilevel patch feature, where the innovation term is derived by projecting the patch into the current image and computing the photometric error for every patch pixel. To keep the computational effort tractable, a QR-decomposition based reduction is applied for obtaining an *equivalent* 2D innovation term per observed landmark. This method takes into account the local texture of a landmark and thereby gains more information along the directions where the patch gradients are stronger. In addition, this offers the possibility to track non-corner shaped features, such as lines, increasing the set of possible image features which is beneficial in scenarios with missing texture.

The main contributions of this work lie in the *fully robocentric* formulation of a visual-inertial odometry as well as in the *tight* integration of the photometric error. The application of manifold encapsulation to bearing vectors can be seen as enabling technique and allows for a sound robocentric representation of 3D landmarks. Furthermore, in contrast to our previous work (Bloesch et al. (2015)), the present approach inherently takes care of landmark tracking by employing an IEKF. This allows per-landmark iterative updates and thus provides simultaneous landmark tracking and full state refinement while considering inertial and visual data. To the best of our knowledge, this *tight* integration of data association and estimation process has no precedent in visual-inertial odometry.

All in all, this paper describes a *fully robocentric* and *direct* visual-inertial odometry framework which runs in real-time on computationally constrained platforms. To increase robustness and usability, we implement multi-camera support (with or without overlapping field of view) and enable online calibration of camera-IMU extrinsics. An in-depth derivation and evaluation of the framework is provided, including experiments on publicly available datasets (Burri et al. (2016)). Our framework, which we refer to as Rovio (RObust Visual-Inertial Odometry), is implemented in C++ and is available as open-source software <sup>\*</sup>.

## 2 Related Work

Within the field of computer vision, Davison (2003) proposed one of the first real-time 3D monocular localization and mapping frameworks. Similarly to the work in this paper,

the author made use of an EKF framework where he co-estimates the absolute position of 3D landmarks. Since then, various research groups have contributed improvements and proposed further approaches. A key issue is to improve the consistency of the estimation framework that is affected by its inherent nonlinearity (Julier and Uhlmann (2001); Castellanos et al. (2004)). One approach is to make use of a robocentric representation for the tracked landmarks and thereby significantly reduce the effect of nonlinearities (Castellanos et al. (2004); Civera et al. (2009)). As an alternative, Huang et al. (2008) propose the use of a so-called observability constrained extended Kalman filter, whereby the inconsistencies can be avoided by using special linearization points while evaluating the system Jacobians.

A somewhat related problem is the choice of the specific representation of a landmark's location. Since the depth of a newly detected landmark is unknown for monocular setups, the initial 3D location estimate exhibits a high uncertainty along the view axis. To integrate this landmark from the beginning into the estimation process, Montiel et al. (2006) proposed the use of an inverse-depth parametrization (IDP). They parametrize each landmark location by the camera position where the landmark was initially detected, by a bearing vector (parametrized with azimuth and elevation angles), as well as the inverse depth of the landmark. The increase in consistency for the IDP and other parametrization methods was further analyzed and confirmed by Solà et al. (2012).

While most standard visual odometry approaches are based on detected and tracked point landmarks as source of visual information, so-called *direct* approaches directly use the image intensities in their estimation framework. Jin et al. (2003) propose to model the environment as a collection of planar patches and to derive a corresponding photometric error between camera frames. Their work is similar to ours in that they also embed the photometric error directly into a filtering framework (but they do not use any inertial data which limits them to slow motions). Molton et al. (2004) also track locally planar image patches in a filter-based SLAM framework. By employing gradient-based image alignment, they also co-estimate surface normals but keep data association separated from the subsequent EKF-based information fusion. Silveira et al. (2008) also use planar regions and minimize the photometric error with respect to a reference frame in order to estimate the relative motion as well as other parameters like illumination parameters and patch normals. They subsequently merge the output in an EKF. More recently, by employing highly optimized SIMD (Single Instruction Multiple Data) implementations, first real-time, CPU-based approaches for semi-dense motion estimation using a monocular camera (Engel et al. (2014); Forster et al. (2014)) have recently been proposed.

Incorporating inertial measurements in the estimation can significantly improve the robustness of the system, provides the estimation process with the notion of gravity, and allows for a more accurate and high bandwidth estimation of the velocities and rotational rates. By adapting the original EKF proposed by Davison (2003), additional IMU

<sup>\*</sup><https://github.com/ethz-asl/rovio>

measurements can be relatively simply integrated into the ego-motion estimation, whereby calibration parameters can be co-estimated online (Kelly and Sukhatme (2011); Jones and Soatto (2011)). Leutenegger et al. (2015) describe a *tightly* coupled approach in which the robot trajectory and sparse 3D landmarks are estimated in a joint optimization problem using inertial error terms as well as the reprojection error of the tracked landmarks in the camera images. This is done in a windowed bundle adjustment approach over a set of keyframe images and a temporal inertial measurement window. Similarly, Mourikis and Roumeliotis (2007) estimate the trajectory in an IMU-driven filtering framework using the reprojection error of 3D landmarks as measurement updates. Instead of adding the landmarks to the filter state, they marginalize them out using a nullspace decomposition, thus leading to a small filter state size. Since inertial measurements are often obtained at a higher rate than image data, methods for combining multiple inertial measurements are desirable to reduce the computational costs. Forster et al. (2016) have presented a concise IMU measurements pre-integration method such that they can be efficiently included in a factor graph framework. Recently, Usenko et al. (2016) have extended their previous work on semi-dense visual odometry (Engel et al. (2014)) in order to integrate inertial measurements. They minimize a joint energy term composed of visual and inertial error terms in order to estimate the ego-motion of their sensor.

Probably the most comparable work to ours was developed by Tanskanen et al. (2015), who implemented an EKF-based framework for merging patch-based photometric errors with IMU measurements. They parameterize their landmarks by the pose of the camera when the landmark was detected as well as the corresponding bearing vector and inverse depth (analogously to Montiel et al. (2006)). Our work differs in that it uses a fully robocentric formulation of the current state, which has various implications on the filtering and visual processing framework. We also integrate a QR-decomposition based measurement space reduction and perform per-landmark update iterations, which are both key to the efficiency and accuracy of our system.

### 3 Prerequisites on Rotations and Unit Vectors

#### 3.1 Notation

For better readability and comprehensibility, we give a brief overview of the employed notations and the algebra of 3D rotations and unit vectors. Three different coordinate frames are used throughout the paper: the inertial world coordinate frame,  $\mathcal{I}$ , the IMU fixed coordinate frame,  $\mathcal{B}$ , as well as the camera fixed coordinate frame,  $\mathcal{C}$ . Only in section 6, where multi-camera setups are discussed, the distinction between the different camera frames will be made. The origin associated with a specific coordinate frame is denoted by the same symbol. In this context, a term of the form  $\mathcal{I}r_{BC}$  denotes the coordinates of a vector from the origin of  $\mathcal{B}$  to the origin of  $\mathcal{C}$ , expressed in the coordinate frame  $\mathcal{I}$ . Furthermore,  $q_{BI}$  is employed in an abstract manner for representing the rotation between a frame  $\mathcal{I}$  and  $\mathcal{B}$ . A good way to think of a rotation is

as a mapping  $q_{BI} : \mathbb{R}^3 \rightarrow \mathbb{R}^3$  between the two associated coordinate frames: Given a physical vector  $r_{BC}$ , a rotation maps the corresponding coordinates from the right index frame to the left index frame, e.g.,  ${}_{\mathcal{B}}r_{BC} = q_{BI}(\mathcal{I}r_{BC})$ . We also employ the mapping  $C(q) : SO(3) \rightarrow \mathbb{R}^{3 \times 3}$  which is defined such that  $q(r) \triangleq C(q)r$  and basically returns the  $3 \times 3$  rotation matrix.

As further abbreviations, we use  $v_B$  for denoting the absolute velocity of  $\mathcal{B}$ , and  $\omega_{IB}$  for the vector describing the relative rotational velocity of the coordinate frame  $\mathcal{B}$  w.r.t. the coordinate frame  $\mathcal{I}$ . In some cases we use further denotations like tildes (measurements) or hats (estimates) if we want to highlight a specific aspect of a quantity. The superscript  $\times$  is used to denote the skew symmetric matrix  $v^\times \in \mathbb{R}^{3 \times 3}$  of a vector  $v \in \mathbb{R}^3$ .

#### 3.2 Representation of 3D Rotations

The set of 3D rotations, the special orthogonal group  $SO(3)$  with group operation  $\otimes$ , is not a vector space and thus adaptations are required in order to enable traditional optimization based methods (e.g. filtering). A mathematically sound and increasingly popular method is to map the region around a selected linearization point to a proper vector space and thereby introduce a local parametrization. There are slight variation in how this concept is formalized and we follow the approach of Hertzberg et al. (2011) as we found it to provide a useful level of abstraction for modeling.

$SO(3)$  is a Lie group and has a logarithmic and an exponential map which map to and from a corresponding Lie algebra  $\mathbb{R}^3$ :

$$\log : SO(3) \rightarrow \mathbb{R}^3, \quad (1)$$

$$q_{BI} \mapsto \log(q_{BI}) = \theta_{BI},$$

$$\exp : \mathbb{R}^3 \rightarrow SO(3), \quad (2)$$

$$\theta_{BI} \mapsto \exp(\theta_{BI}) = q_{BI}.$$

There is a certain amount of freedom in selecting these maps. Here, we select the exponential and logarithmic maps such that  $\theta_{BI}$  in the above equations coincides with the passive rotation vector of the rotation  $q_{BI}$ . We can write the following identities (the last identity is known as Rodrigues' formula):

$$\exp(-\theta) = \exp(\theta)^{-1}, \quad (3)$$

$$\exp(q(\theta)) = q \otimes \exp(\theta) \otimes q^{-1}, \quad (4)$$

$$C(\theta) = I - \frac{\sin(\|\theta\|)\theta^\times}{\|\theta\|} + \frac{(1 - \cos(\|\theta\|))\theta^\times{}^2}{\|\theta\|^2}. \quad (5)$$

The exponential and logarithmic maps can be used to introduce a boxplus ( $\boxplus$ ) and a boxminus ( $\boxminus$ ) operator, which adopt the role of addition and subtraction operators for rotations. Using a slightly different notation than Hertzberg et al. (2011), we define:

$$\boxplus : SO(3) \times \mathbb{R}^3 \rightarrow SO(3), \quad (6)$$

$$q, \theta \mapsto \exp(\theta) \otimes q,$$

$$\boxminus : SO(3) \times SO(3) \rightarrow \mathbb{R}^3, \quad (7)$$

$$q, p \mapsto \log(q \otimes p^{-1}).$$

Similarly to regular addition and subtraction, both operators fulfill the following axioms:

$$\mathbf{q} \boxplus \mathbf{0} = \mathbf{q}, \quad (8)$$

$$(\mathbf{q} \boxplus \boldsymbol{\theta}) \boxminus \mathbf{q} = \boldsymbol{\theta}, \quad (9)$$

$$\mathbf{q} \boxplus (\mathbf{p} \boxminus \mathbf{q}) = \mathbf{p}. \quad (10)$$

This approach distinguishes between actual rotations which are on  $SO(3)$  (Lie group) and differences of rotations which lie on  $\mathbb{R}^3$  (Lie algebra). The above operators take care of appropriately transforming the elements into their respective spaces and allow a smooth embedding of rotational quantities in filtering and optimization frameworks.

The definition of differentials involving rotation can be adapted by replacing the regular plus and minus operators by the above boxplus and boxminus operators. For instance the differential of a mapping  $\mathbf{q}(x) : \mathbb{R} \rightarrow SO(3)$  can be defined as:

$$\frac{\partial}{\partial x} \mathbf{q}(x) := \lim_{\epsilon \rightarrow 0} \frac{\mathbf{q}(x + \epsilon) \boxminus \mathbf{q}(x)}{\epsilon}. \quad (11)$$

The same can be done for the other way round where we have a mapping  $x(\mathbf{q}) : SO(3) \rightarrow \mathbb{R}$ :

$$\frac{\partial}{\partial \mathbf{q}} x(\mathbf{q}) := \lim_{\epsilon \rightarrow 0} \begin{bmatrix} \frac{x(\mathbf{q} \boxplus (\mathbf{e}_1 \epsilon)) - x(\mathbf{q})}{\epsilon} \\ \frac{x(\mathbf{q} \boxplus (\mathbf{e}_2 \epsilon)) - x(\mathbf{q})}{\epsilon} \\ \frac{x(\mathbf{q} \boxplus (\mathbf{e}_3 \epsilon)) - x(\mathbf{q})}{\epsilon} \end{bmatrix}^T \quad (12)$$

where  $\mathbf{e}_{1/2/3}$  are orthonormal basis vectors. This results in the following frequently-used derivatives (these may vary depending on conventions):

$$\partial / \partial t (\mathbf{q}_{BI}(t)) = \mathbf{B} \boldsymbol{\omega}_{IB}(t), \quad (13)$$

$$\partial / \partial \mathbf{q} (\mathbf{q}(\mathbf{r})) = (\mathbf{q}(\mathbf{r}))^\times, \quad (14)$$

$$\partial / \partial \mathbf{q} (\mathbf{q}^{-1}) = -\mathbf{C}(\mathbf{q})^T, \quad (15)$$

$$\partial / \partial \mathbf{q} (\mathbf{q} \otimes \mathbf{p}) = \mathbf{I}, \quad (16)$$

$$\partial / \partial \mathbf{q} (\mathbf{p} \otimes \mathbf{q}) = \mathbf{C}(\mathbf{p}), \quad (17)$$

$$\partial / \partial \boldsymbol{\theta} (\exp(\boldsymbol{\theta})) = \boldsymbol{\Gamma}(\boldsymbol{\theta}), \quad (18)$$

$$\partial / \partial \mathbf{q} (\log(\mathbf{q})) = \boldsymbol{\Gamma}^{-1}(\log(\mathbf{q})). \quad (19)$$

The derivative of the exponential map is given by the Jacobian  $\boldsymbol{\Gamma}(\boldsymbol{\theta}) \in \mathbb{R}^{3 \times 3}$  which has the following analytical expression:

$$\boldsymbol{\Gamma}(\boldsymbol{\theta}) = \mathbf{I} - \frac{(1 - \cos(\|\boldsymbol{\theta}\|)) \boldsymbol{\theta}^\times}{\|\boldsymbol{\theta}\|^2} + \frac{(\|\boldsymbol{\theta}\| - \sin(\|\boldsymbol{\theta}\|)) \boldsymbol{\theta}^\times{}^2}{\|\boldsymbol{\theta}\|^3}. \quad (20)$$

The above derivations are independent of the actual numerical parametrization of 3D rotations (e.g. quaternions or rotation matrices) as long as it is lossless and the associated operations adhere to certain rules. A more detailed discussion and derivations can be found in Bloesch et al. (2016). The employed parametrization in the implementation is based on unit quaternions using the Hamilton convention. For a unit quaternion  $\mathbf{q}$  with real part  $q_0$  and imaginary part  $\tilde{\mathbf{q}}$  we employ the following exponential and logarithmic maps:

$$\exp(\boldsymbol{\theta}) = (q_0, \tilde{\mathbf{q}}) = \left( \cos(\|\boldsymbol{\theta}\|/2), \sin(\|\boldsymbol{\theta}\|/2) \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|} \right) \quad (21)$$

$$\log(\mathbf{q}) = 2 \operatorname{atan2}(\|\tilde{\mathbf{q}}\|, q_0) \frac{\tilde{\mathbf{q}}}{\|\tilde{\mathbf{q}}\|} \quad (22)$$

where  $\boldsymbol{\theta} \in \mathbb{R}^3$  can be interpreted as the corresponding rotation vector. Both maps together with the quaternion multiplication are the only parametrization specific operations that are required.

### 3.3 Representation of 3D Unit Vectors

While the above handling of rotations has been used similarly in previous filtering frameworks (e.g. Li and Mourikis (2013); Bloesch et al. (2013)), we extend the methodology to 3D unit vectors on the 2-sphere  $S^2$ . This is done analogously to Hertzberg et al. (2011), whereas we employ a parametrization yielding simple analytical derivatives and guarantee second order differentiability. A main issue with 3D unit vectors is to select orthonormal vectors for spanning the tangent space such that a suitable difference operator can be defined. Assigning orthonormal vectors to every point on the 2-sphere creates a vector field and as stated by the ‘‘hairy ball theorem’’, there is no continuous way of doing so over the full 2-sphere. To solve this issue we employ a rotation,  $\boldsymbol{\mu} \in SO(3)$ , as underlying representation for unit vectors and define the following quantities:

$$\mathbf{n}(\boldsymbol{\mu}) := \boldsymbol{\mu}(\mathbf{e}_z) \in S^2 \subset \mathbb{R}^3, \quad (23)$$

$$\mathbf{N}(\boldsymbol{\mu}) := [\boldsymbol{\mu}(\mathbf{e}_x), \boldsymbol{\mu}(\mathbf{e}_y)] \in \mathbb{R}^{3 \times 2}, \quad (24)$$

where  $\mathbf{e}_{x/y/z} \in \mathbb{R}^3$  are the basis vectors of an arbitrary orthonormal coordinate system. The actual unit vector is given by  $\mathbf{n}(\boldsymbol{\mu})$  which results when rotating  $\mathbf{e}_z$  by  $\boldsymbol{\mu}$  (if the context is clear we directly refer to the unit vector using  $\boldsymbol{\mu}$ ). The matrix  $\mathbf{N}(\boldsymbol{\mu})$  is composed of the rotated  $\mathbf{e}_x$  and  $\mathbf{e}_y$  and spans the tangent space. While such a construction of the tangent space is not deterministic since infinitely many rotations  $\boldsymbol{\mu}$  provide the same unit vector  $\mathbf{n}(\boldsymbol{\mu})$ , we have the advantage that smooth transformations of the rotation  $\boldsymbol{\mu}$  induce smooth transformations of the associated tangent space.

The tangent space can be used to define the following boxplus and boxminus operators:

$$\boxplus : SO(3) \times \mathbb{R}^2 \rightarrow SO(3), \quad (25)$$

$$\boldsymbol{\mu}, \mathbf{u} \mapsto \exp(\mathbf{N}(\boldsymbol{\mu})\mathbf{u}) \otimes \boldsymbol{\mu},$$

$$\boxminus : SO(3) \times SO(3) \rightarrow \mathbb{R}^2, \quad (26)$$

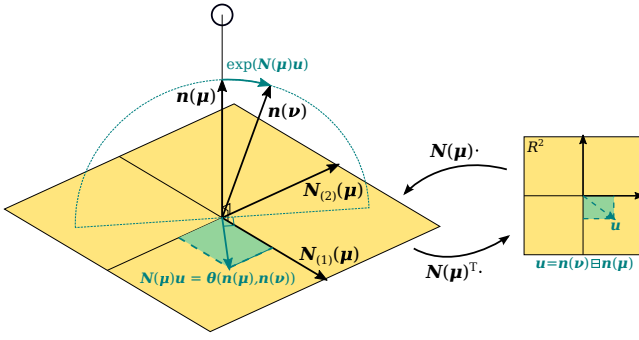
$$\boldsymbol{\nu}, \boldsymbol{\mu} \mapsto \mathbf{N}(\boldsymbol{\mu})^T \boldsymbol{\theta}(\boldsymbol{\mu}, \boldsymbol{\nu}),$$

where  $\boldsymbol{\theta}$  maps two unit vectors to the minimal rotation vector between them:

$$\boldsymbol{\theta}(\mathbf{n}(\boldsymbol{\mu}), \mathbf{n}(\boldsymbol{\nu})) = \frac{\arccos(\mathbf{n}(\boldsymbol{\nu})^T \mathbf{n}(\boldsymbol{\mu}))}{\|\mathbf{n}(\boldsymbol{\nu}) \times \mathbf{n}(\boldsymbol{\mu})\|} \mathbf{n}(\boldsymbol{\nu}) \times \mathbf{n}(\boldsymbol{\mu}). \quad (27)$$

A visualization of the 2-sphere and the tangent space for a specific  $\boldsymbol{\mu}$  is given in Figure 1.

The concept is slightly more complicated than in the case of 3D rotations since we truly over-parameterize a 3D unit vector (no constraint is imposed on the underlying rotation). To overcome this, we use a different notion of equivalence where we define that two unit vector parametrizations  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  are equivalent ( $\boldsymbol{\mu} \sim \boldsymbol{\nu}$ ) iff  $\mathbf{n}(\boldsymbol{\mu}) = \mathbf{n}(\boldsymbol{\nu})$ . With this, the



**Figure 1.** Representation of 3D unit vectors: The 3D unit vector  $\mathbf{n}(\boldsymbol{\mu})$  is represented as the result of applying the rotation  $\boldsymbol{\mu}$  onto the z-axis of an arbitrary inertial coordinate system. The images of the x- and y-axis are used to define an orthonormal plane to the unit vector. This plane then represents the tangent space used for the construction of the boxplus and boxminus operations. The boxminus operator takes two 3D unit vectors and represents their difference in  $\mathbb{R}^2$ . Conversely, the boxplus operator takes an element from  $\mathbb{R}^2$  and applies it on a 3D unit vector.

axioms proposed by Hertzberg et al. (2011) are fulfilled:

$$\boldsymbol{\mu} \boxplus \mathbf{0} = \boldsymbol{\mu}, \quad (28)$$

$$(\boldsymbol{\mu} \boxplus \mathbf{u}) \boxminus \boldsymbol{\mu} = \mathbf{u}, \quad (29)$$

$$\boldsymbol{\mu} \boxplus (\boldsymbol{\nu} \boxminus \boldsymbol{\mu}) \sim \boldsymbol{\nu}. \quad (30)$$

A technical detail with this parametrization is that whenever representing a difference,  $\mathbf{u} \in \mathbb{R}^2$ , we have to keep track of the corresponding tangent space. Mathematically, if we have two rotations  $\boldsymbol{\mu}$  and  $\boldsymbol{\nu}$  with  $\boldsymbol{\mu} \sim \boldsymbol{\nu}$ , it does not follow that  $(\boldsymbol{\mu} \boxplus \mathbf{u}) \sim (\boldsymbol{\nu} \boxplus \mathbf{u})$ .

Similarly to the derivatives given in section 3.2, the most commonly used derivatives for terms involving 3D unit vectors are given by:

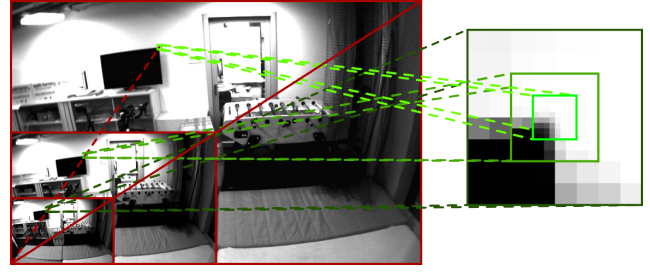
$$\frac{\partial}{\partial t} (\boldsymbol{\mu}(t)) = -\mathbf{N}(\boldsymbol{\mu}(t))^T \mathbf{n}(\boldsymbol{\mu}(t))^\times \cdot \frac{\partial}{\partial t} (\mathbf{n}(\boldsymbol{\mu}(t))), \quad (31)$$

$$\frac{\partial}{\partial \boldsymbol{\mu}} (\mathbf{n}(\boldsymbol{\mu})) = \mathbf{n}(\boldsymbol{\mu})^\times \mathbf{N}(\boldsymbol{\mu}), \quad (32)$$

$$\frac{\partial}{\partial \boldsymbol{\mu}} (\mathbf{N}(\boldsymbol{\mu})^T \mathbf{r}) = -\mathbf{N}(\boldsymbol{\mu})^T \mathbf{r}^\times \mathbf{N}(\boldsymbol{\mu}). \quad (33)$$

The first identity relates the time derivative of a 3D unit vector on its manifold to its time derivative in the 3D vector space. The second expression is the derivative of the unit vector in 3D w.r.t. to its minimal 2D representation. Those identities can be very useful when computing Jacobians, whereby the chain rule can be applied for computing the derivatives of more complex terms. An example will be provided when discussing the process model of the bearing vector state of 3D landmarks (see section 5.3 and Appendix A).

All in all, the proposed unit vector parametrization yields analogous advantages as obtained when employing the well established minimal 3D rotation parametrization. This includes a singularity-free parametrization which comes with relatively simple differentials. Furthermore the parametrization of the tangent space is orthogonal and the direction of the boxminus operation is in accordance with the shortest path between two given unit vectors (taking a step along  $\boldsymbol{\nu} \boxminus \boldsymbol{\mu}$  is optimal for going from  $\boldsymbol{\mu}$  to  $\boldsymbol{\nu}$ ,



**Figure 2.** The construction of a multilevel patch out of an image pyramid. Here each single patch is composed of  $8 \times 8$  pixels and 3 pyramid levels are depicted. These settings may vary in the actual implementation.

see Figure 1). Other parametrizations, such as azimuth and elevation angles, do not meet these properties and often exhibit singular configurations.

## 4 Multilevel Patches and Photometric Error

### 4.1 Multilevel Patch Features

Along the lines of other landmark-based visual odometry approaches (Davison (2003)) we model landmarks as distinguished stationary 3D locations in the environment. Each landmark is associated with a multilevel patch feature  $P = \{P_0, \dots, P_L\}$ , which is composed of multiple  $n \times n$  image patches,  $P_l$ , extracted at the projected landmark location on image level  $l$ . In the current default implementation we extract  $6 \times 6$  image patches on the second and third pyramid level (down-sampling factor of 2). These parameters can and should be adapted to the actual hardware setup and the scenario. An example is given in Figure 2. The simultaneous use of multiple pyramid levels leads to cross-correlations between the pixel intensities. These are not explicitly modeled but can be handled to a certain extent by tuning the corresponding error weighting.

In contrast to a standard feature descriptor, a patch-based descriptor allows to compute a photometric error and thereby to avoid the use of reprojection errors. Taking the information of every pixel gives much richer information about the environment, which not only helps improving the robustness in bad lighting conditions, but also inherently takes into account the texture of the tracked image patch. For instance, it enables the integration of edge-shaped features, whereby the gained information is along the perpendicular direction of the edge. In comparison, reprojection error based approaches typically attempt to minimize the distance between the predicted and detected feature location. This ignores the local texture around the landmark and, if no additional measures are taken, all landmarks are weighted equally.

### 4.2 Projection Model and Linear Warping

Given the bearing vector  $\boldsymbol{\mu}$  of a landmark, the pixel coordinates in the camera frame can be retrieved by using the camera model  $\boldsymbol{\pi}$ . Assuming a known intrinsic calibration, the pixel coordinates  $\mathbf{p}$  can directly be expressed by  $\mathbf{p} = \boldsymbol{\pi}(\boldsymbol{\mu})$ . If the camera is moving, the feature moves through the image and is seen from a different perspective. To account for a certain patch distortion effect, a linear warping matrix is

tracked with each feature. This is done by concatenating all Jacobians when transforming a landmark location. For instance, if we detect a feature in some frame at pixel  $\mathbf{p}_1$ , transform the corresponding bearing vector  $\boldsymbol{\mu}_1 = \boldsymbol{\pi}^{-1}(\mathbf{p}_1)$  with a process model  $\boldsymbol{\mu}_2 = \mathbf{f}(\boldsymbol{\mu}_1)$ , and then re-project the bearing vector in a subsequent frame  $\mathbf{p}_2 = \boldsymbol{\pi}(\boldsymbol{\mu}_2)$ , we obtain the following linear warping matrix:

$$\mathbf{D} = \frac{\partial \boldsymbol{\pi}(\boldsymbol{\mu}_2)}{\partial \boldsymbol{\mu}_2} \frac{\partial \mathbf{f}(\boldsymbol{\mu}_1)}{\partial \boldsymbol{\mu}_1} \frac{\partial \boldsymbol{\pi}^{-1}(\mathbf{p}_1)}{\partial \mathbf{p}_1} \in \mathbb{R}^{2 \times 2}. \quad (34)$$

In essence, this maps the two patch axes from the point of patch extraction (which were aligned with the image axes) to the two distorted patch axes in the projection image. This approach tracks the distortion locally around the patch and ignores any larger scale information like the geometric shape of a patch. To avoid large distortions and the accumulation of errors, the patches are re-extracted regularly and the warping matrix is reset to identity.

### 4.3 Photometric Error and Patch Alignment

The photometric error between a given multilevel patch feature and a specific image is computed by extracting a warped patch at the estimated location and evaluating the pixel-wise intensity error. For a given multilevel patch feature (with coordinates  $\mathbf{p}$  and multilevel patch  $P = \{P_0, \dots, P_L\}$ ) at a specific image level  $l$  and patch pixel  $\mathbf{p}_j$ , the photometric error can be formalized as follows:

$$e_{l,j}(\mathbf{p}, P, I, \mathbf{D}) = P_l(\mathbf{p}_j) - a I_l(\mathbf{p} s_l + \mathbf{D} \mathbf{p}_j) - b, \quad (35)$$

where  $I_l$  is the image at the pyramid level  $l$  and  $s_l = 0.5^l$  is a scaling factor to account for the down-sampling. The linear warping matrix  $\mathbf{D}$  is used to map patch pixel coordinates to image coordinates. Furthermore, inter-frame illumination changes are taken into account by employing an affine intensity model composed of the scalars  $a$  and  $b$  (both get marginalized out). Figure 3 depicts the photometric error between a patch and its measurement in an image at a predicted location  $\mathbf{p}$ .

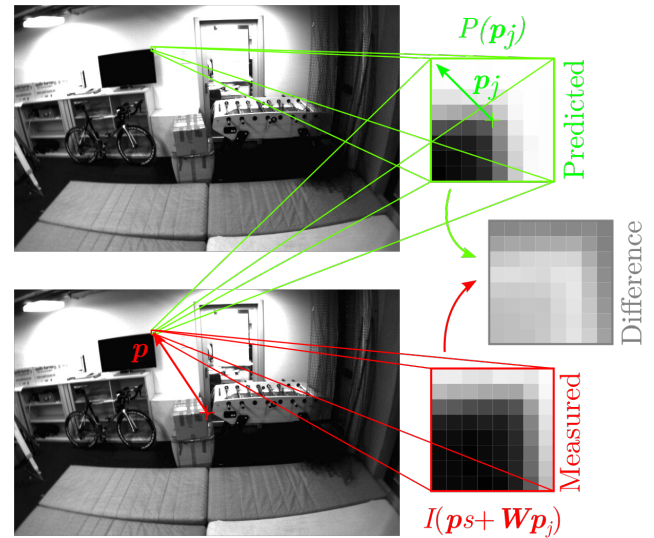
If we minimize the squared error terms for a multilevel patch, we obtain a patch alignment algorithm which is very similar to the well-known Kanade-Lucas-Tomasi (KLT) feature tracker (Lucas and Kanade (1981); Shi and Tomasi (1994)). A slight difference is given by the fact that we optimize over multiple image levels at once. The minimization can be solved by a Gauss-Newton method which iteratively linearizes the optimization problem around an estimated patch location  $\hat{\mathbf{p}}$ :

$$\mathbf{b}(\hat{\mathbf{p}} + \delta \mathbf{p}, P, I, \mathbf{D}) = \mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D}) \delta \mathbf{p} + \mathbf{b}(\hat{\mathbf{p}}, P, I, \mathbf{D}), \quad (36)$$

where  $\mathbf{b}(\hat{\mathbf{p}}, P, I, \mathbf{D})$  represents the stacked error terms from eq. (35) and  $\mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D})$  the corresponding Jacobian. The corresponding normal equations are then given by:

$$\mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D})^T \mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D}) \delta \mathbf{p} = -\mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D})^T \mathbf{b}(\hat{\mathbf{p}}, P, I, \mathbf{D}), \quad (37)$$

which can be solved for the correction  $\delta \mathbf{p}$ . This is analogous to one iteration step of the KLT feature tracker (but is not used as such in Rovio). In section 5.4 we will demonstrate



**Figure 3.** Illustration of the (signed) photometric error between a previously extracted patch (green) and its projection into an image (measured, red) at a predicted location  $\mathbf{p}$ . The bottom left grey tone of the difference patch represents 0. Only a single image level is depicted. This photometric error is directly used as the innovation term in an IEKF.

how eq. (36) is leveraged into the innovation term of the employed IEKF.

Note that due to the scaling factor  $s_l$  in eq. (35), error terms for higher image levels will have a weaker corrective influence on the filter state or the patch alignment. On the other hand, they exhibit increased robustness w.r.t. image blur or bad initial alignment and thus strongly increase the robustness of the overall alignment method.

### 4.4 Detection and Scoring

The detection of new landmarks is based on the FAST corner detector (Rosten and Drummond (2006)) which provides a large amount of candidate feature locations. After removing candidates which are close to currently tracked features, we compute a patch gradient based score for selecting new features which are added to the state. This basically represents an adaptation of the Shi-Tomasi score (Shi and Tomasi (1994)) by considering the combined Hessian on multiple image levels, instead of only a single level. The combined Hessian can be directly retrieved from the normal equations (37):

$$\mathbf{H} = \mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D})^T \mathbf{A}(\hat{\mathbf{p}}, I, \mathbf{D}), \quad (38)$$

where the minimal eigenvalue of  $\mathbf{H}$  corresponds to the adapted Shi-Tomasi score.

The advantage is that a high score is correlated with the alignment accuracy of the corresponding multilevel patch feature. Instead of returning the minimal eigenvalue, the method can return other eigenvalue based scores like the 1- or 2-norm. This is useful in environments with scarce corner data, whereby also edge-shaped features can be considered. Finally, the detection process is also coupled with a bucketing technique to achieve a good distribution of the features within the camera frame.

## 5 Filter Framework

### 5.1 Iterated Extended Kalman Filtering

The regular Kalman filter can be interpreted as the recursive optimal solution to the maximum likelihood estimation problem formulated over two subsequent time steps (Bell and Cathey (1993)). Analogously, the EKF can be associated with a nonlinear maximum likelihood estimation and can be shown to yield the same result as the first iteration step of a corresponding Gauss-Newton optimization. However, in contrast to its linear counterpart, the EKF cannot guarantee to retrieve the optimal solution, whereby linearization errors tend to become larger if the linearization point is further away from the real solution. A possibility to improve this aspect is to make use of an IEKF which is basically the recursive form of the Gauss-Newton optimization (Bell and Cathey (1993)).

A nonlinear discrete time system with state  $\mathbf{x}$ , innovation term  $\mathbf{y}$ , process noise  $\mathbf{w} \sim \mathcal{N}(0, \mathbf{W})$ , and update noise  $\mathbf{n} \sim \mathcal{N}(0, \mathbf{R})$  can be written as:

$$\mathbf{x}_k = f(\mathbf{x}_{k-1}, \mathbf{w}_{k-1}), \quad (39)$$

$$\mathbf{y}_k = h(\mathbf{x}_k, \mathbf{n}_k). \quad (40)$$

In eq. (40) we made use of an implicit formulation of the measurement model which directly yields the Kalman innovation  $\mathbf{y}_k$ . This provides more flexibility in the design by allowing the direct integration of residuals. Given an a-posteriori estimate  $\mathbf{x}_{k-1}^+$  with covariance  $\mathbf{P}_{k-1}^+$ , the prediction step of the IEKF is analogous to the EKF and yields the a-priori estimate at the next time step:

$$\mathbf{x}_k^- = f(\mathbf{x}_{k-1}^+, \mathbf{0}), \quad (41)$$

$$\mathbf{P}_k^- = \mathbf{F}_{k-1} \mathbf{P}_{k-1}^+ \mathbf{F}_{k-1}^T + \mathbf{G}_{k-1} \mathbf{W}_{k-1} \mathbf{G}_{k-1}^T, \quad (42)$$

with the Jacobians

$$\mathbf{F}_{k-1} = \frac{\partial f}{\partial \mathbf{x}_{k-1}}(\mathbf{x}_{k-1}^+, \mathbf{0}), \quad (43)$$

$$\mathbf{G}_{k-1} = \frac{\partial f}{\partial \mathbf{w}_{k-1}}(\mathbf{x}_{k-1}^+, \mathbf{0}). \quad (44)$$

Analogously to the EKF, the update step of the IEKF can be linked to an optimization problem considering the deviation from the prior  $\mathbf{x}_k^-$  and the innovation term  $h(\mathbf{x}_k^+, \mathbf{0})$ :

$$\min_{\mathbf{x}_k^+} \|\mathbf{x}_k^+ \boxminus \mathbf{x}_k^-\|_{\mathbf{P}_k^{-1}} + \|h(\mathbf{x}_k^+, \mathbf{0})\|_{(\mathbf{J}_k \mathbf{R}_k \mathbf{J}_k^T)^{-1}}. \quad (45)$$

However, in contrast to the EKF, an iterative scheme is employed where the problem is linearized around continuously refined linearization points  $\mathbf{x}_{k,j}^+$  starting with  $\mathbf{x}_{k,0}^+ = \mathbf{x}_k^-$ :

$$\min_{\Delta \mathbf{x}_{k,j}} \|\mathbf{x}_{k,j}^+ \boxminus \mathbf{x}_k^- + \mathbf{L}_{k,j}^{-1} \Delta \mathbf{x}_{k,j}\|_{\mathbf{P}_k^{-1}} + \|h(\mathbf{x}_{k,j}^+, \mathbf{0}) + \mathbf{H}_{k,j} \Delta \mathbf{x}_{k,j}\|_{(\mathbf{J}_k \mathbf{R}_k \mathbf{J}_k^T)^{-1}} \quad (46)$$

where the Jacobians are updated every iteration step:

$$\mathbf{H}_{k,j} = \frac{\partial h}{\partial \mathbf{x}_k}(\mathbf{x}_{k,j}^+, \mathbf{0}), \quad (47)$$

$$\mathbf{J}_{k,j} = \frac{\partial h}{\partial \mathbf{n}_k}(\mathbf{x}_{k,j}^+, \mathbf{0}), \quad (48)$$

$$\mathbf{L}_{k,j} = \frac{\partial(\mathbf{x}_k^- \boxplus \Delta \mathbf{x})}{\partial \Delta \mathbf{x}}(\mathbf{x}_{k,j}^+ \boxminus \mathbf{x}_k^-). \quad (49)$$

The Jacobian  $\mathbf{L}_{k,j}$  of the boxplus operator is required to account for the special linearization of certain states such as rotations or bearing vectors. Its inverse  $\mathbf{L}_{k,j}^{-1}$  is the corresponding Jacobian of the boxminus operation w.r.t. to the left operand and is required to linearize the deviation of the prior in (46). Please note that due to the special notion of differentials on manifolds the Jacobian  $\mathbf{L}_{k,j}$  is a square matrix (see eq. (11)). Also, in the case of vector spaces this Jacobian will be the identity matrix.

Setting the derivative of the cost function (46) w.r.t. the incremental update  $\Delta \mathbf{x}_{k,j}$  to zero and employing some matrix calculus yields the following recursive solution:

$$\mathbf{S}_{k,j} = \mathbf{H}_{k,j} \mathbf{L}_{k,j}^T \mathbf{P}_k^- \mathbf{L}_{k,j} \mathbf{H}_{k,j}^T + \mathbf{J}_{k,j} \mathbf{R}_k \mathbf{J}_{k,j}^T, \quad (50)$$

$$\mathbf{K}_{k,j} = \mathbf{L}_{k,j}^T \mathbf{P}_k^- \mathbf{L}_{k,j} \mathbf{H}_{k,j}^T \mathbf{S}_{k,j}^{-1}, \quad (51)$$

$$\Delta \mathbf{x}_{k,j} = \mathbf{K}_{k,j} \left( \mathbf{H}_{k,j} \mathbf{L}_{k,j} (\mathbf{x}_{k,j}^+ \boxminus \mathbf{x}_k^-) \right. \quad (52)$$

$$\left. - h(\mathbf{x}_{k,j}^+, \mathbf{0}) \right) - \mathbf{L}_{k,j} (\mathbf{x}_{k,j}^+ \boxminus \mathbf{x}_k^-),$$

$$\mathbf{x}_{k,j+1}^+ = \mathbf{x}_{k,j}^+ \boxplus \Delta \mathbf{x}_{k,j}, \quad (53)$$

whereby the iteration is terminated when the update  $\Delta \mathbf{x}_{k,j}$  is below a certain threshold. Finally, the covariance matrix is only updated once the process has converged after  $n$  iterations:

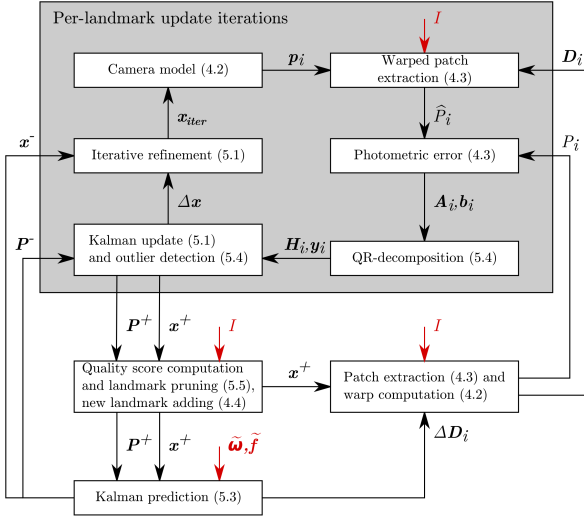
$$\mathbf{P}_k^+ = (\mathbf{I} - \mathbf{K}_{k,n} \mathbf{H}_{k,n}) \mathbf{L}_{k,n}^T \mathbf{P}_k^- \mathbf{L}_{k,n}. \quad (54)$$

Especially in setups with large initial uncertainties, the IEKF can provide a significant gain in convergence and accuracy. Using a termination criteria based on the magnitude of the performed correction, the computational overhead can be limited to cases with large update corrections (e.g. the initial measurements of a new landmark). Once the state has properly converged, the number of iterations can be kept to a minimum and the computational costs remain similar to the ones of the regular EKF.

### 5.2 Filter Setup and State Definition

Similar to other visual-inertial filtering frameworks (Kelly and Sukhatme (2011); Jones and Soatto (2011)), the inertial measurements are employed to propagate the filter state, while the visual measurements are processed and integrated during the filter update step (see Figure 4). The proposed filter setup differs in that it makes use of a fully robocentric formulation of the filter state, which has previously been tested in vision-only approaches Civera et al. (2009). The advantage of this formulation is that the position and yaw states, which are unobservable, can be fully decoupled from the rest of the filter states. This decreases the noise magnitude and improves the consistency of relevant states such as





**Figure 4.** Flow graph of the filter framework. The processing of the inertial ( $\tilde{\omega}$ ,  $\tilde{f}$ ) and visual ( $I$ ) measurements (red terms) is tightly embedded in the IEKF framework. The filter state  $x$  and covariance  $P$  are alternatively processed by an IMU-based prediction step and a vision-based update step. Most of the visual information is processed as part of the iterative update step (grey block). The order and dependencies of the image processing sub-steps are also illustrated and the numbers represent the references to the corresponding sections.

velocity or inclination angles. On the other hand, noise from the gyroscope affects all states that need to be rotated during the state propagation (see section 5.3). However, since the gyroscope noise is often relatively small and because most states are observable, this does not represent a significant issue.

The state of the filter is composed of the following elements (including  $N$  visual landmarks):

$$x := (\mathbf{r}, \mathbf{v}, \mathbf{q}, \mathbf{b}_f, \mathbf{b}_\omega, \mathbf{c}, \mathbf{z}, \boldsymbol{\mu}_0, \dots, \boldsymbol{\mu}_N, \rho_0, \dots, \rho_N) \quad (55)$$

with:

- $\mathbf{r} := {}_{\mathcal{B}}\mathbf{r}_{\mathcal{I}\mathcal{B}}$ : robocentric position of IMU,
- $\mathbf{v} := {}_{\mathcal{B}}\mathbf{v}_{\mathcal{B}}$ : robocentric velocity of IMU,
- $\mathbf{q} := \mathbf{q}_{\mathcal{I}\mathcal{B}}$ : attitude of IMU (map from  $\mathcal{B}$  to  $\mathcal{I}$ ),
- $\mathbf{b}_f$ : additive bias on accelerometer (expressed in  $\mathcal{B}$ ),
- $\mathbf{b}_\omega$ : additive bias on gyroscope (expressed in  $\mathcal{B}$ ),
- $\mathbf{c} := {}_{\mathcal{B}}\mathbf{r}_{\mathcal{B}\mathcal{C}}$ : linear part of IMU-camera extrinsics,
- $\mathbf{z} := \mathbf{q}_{\mathcal{C}\mathcal{B}}$ : rotational part of IMU-camera extrinsics,
- $\boldsymbol{\mu}_i$ : bearing vector to landmark  $i$  (expressed in  $\mathcal{C}$ ),
- $\rho_i$ : distance parameter of landmark  $i$ .

The generic parametrization for the distance  $d_i$  of a landmark  $i$  is given by the mapping  $d_i = d(\rho_i)$  with derivative  $d'(\rho_i)$ . In the context of this work we mainly tested the inverse distance parametrization,  $d(\rho_i) = 1/\rho_i$ . A brief comparison with the regular distance parametrization is provided in section 7.3.

Rotations ( $\mathbf{q}$ ,  $\mathbf{z}$ ) and bearing vectors ( $\boldsymbol{\mu}_i$ ) are parametrized as detailed in section 3.2 and section 3.3. This means that quantities like differences, uncertainties, or errors are represented as elements of a vector space with minimal dimension, i.e., 3D for rotations and 2D for

bearing vectors. By using the combined bearing vector and distance parametrization, landmarks can be initialized in an *undelayed* manner and can be integrated into the filter at detection time. The distance of a landmark is initialized with a fixed value or, if sufficiently converged, with an estimate of the current average scene distance. The corresponding covariance is set to a very large value. In comparison to other inverse-depth parametrizations, we do not over-parametrize the 3D landmark location estimates, whereby each landmark corresponds to 3 columns in the covariance matrix of the state (2 for the bearing vector and 1 for the distance parameter). This also avoids the need for re-parametrization (Solà et al. (2012)).

A singularity-free parametrization of bearing vectors on the full unit sphere is essential here. It enables the proper representation of bearing vectors and their uncertainty estimates even if outside the field of view of the camera. Furthermore, limiting the validity of the parametrization to a certain region would render online camera-IMU extrinsics calibration and multi-camera support more difficult.

### 5.3 State Propagation

The state propagation is driven by the proper acceleration measurement,  $\tilde{f} = {}_{\mathcal{B}}\tilde{f}_{\mathcal{B}}$ , and the rotational rate measurement,  $\tilde{\omega} = {}_{\mathcal{B}}\tilde{\omega}_{\mathcal{I}\mathcal{B}}$ . Both measurements are modeled as noise and bias affected leading to the following bias corrected but noise affected estimates:

$$\hat{f} = \tilde{f} - \mathbf{b}_f - \mathbf{w}_f, \quad (56)$$

$$\hat{\omega} = \tilde{\omega} - \mathbf{b}_\omega - \mathbf{w}_\omega. \quad (57)$$

Together with the estimated camera linear velocity and rotational rate

$$\hat{v}_c = \mathbf{z}(\mathbf{v} + \hat{\omega}^\times \mathbf{c}), \quad (58)$$

$$\hat{\omega}_c = \mathbf{z}(\hat{\omega}), \quad (59)$$

this yields the following set of continuous differential equations:

$$\dot{\mathbf{r}} = -\hat{\omega}_c^\times \mathbf{r} + \hat{v}_c + \mathbf{w}_r, \quad (60)$$

$$\dot{\mathbf{v}} = -\hat{\omega}_c^\times \mathbf{v} + \hat{f} + \mathbf{q}^{-1}(\mathbf{g}), \quad (61)$$

$$\dot{\mathbf{q}} = -\mathbf{q}(\hat{\omega}), \quad (62)$$

$$\dot{\mathbf{b}}_f = \mathbf{w}_{b_f}, \quad (63)$$

$$\dot{\mathbf{b}}_\omega = \mathbf{w}_{b_\omega}, \quad (64)$$

$$\dot{\mathbf{c}} = \mathbf{w}_c, \quad (65)$$

$$\dot{\mathbf{z}} = \mathbf{w}_z, \quad (66)$$

$$\dot{\boldsymbol{\mu}}_i = \mathbf{N}(\boldsymbol{\mu}_i)^T \left( \hat{\omega}_c + \mathbf{n}(\boldsymbol{\mu}_i)^\times \frac{\hat{v}_c}{d(\rho_i)} \right) + \mathbf{w}_{\boldsymbol{\mu}_i}, \quad (67)$$

$$\dot{\rho}_i = -\mathbf{n}(\boldsymbol{\mu}_i)^T \hat{v}_c / d'(\rho_i) + \mathbf{w}_{\rho_i}. \quad (68)$$

The term  $\mathbf{N}(\boldsymbol{\mu}_i)^T$  projects a 3D vector onto the 2D tangent space at the bearing vector  $\boldsymbol{\mu}_i$  (see Figure 1). Furthermore,  $\mathbf{g}$  is the gravity vector expressed in the world coordinate frame, and the terms of the form  $\mathbf{w}_*$  are white Gaussian noise processes. The corresponding covariance parameters can either be derived from the IMU specifications or can be tuned manually.

Note that the derivatives of bearing vectors and rotations lie within 2D and 3D vector spaces, respectively. This is required for achieving a minimal and consistent representation of the filter state and covariance. While most of the above derivatives are relatively well known, the dynamics of the bearing vector and distance parameter is a novelty in this work. We give a sketch of the corresponding derivation in Appendix A. It relies on the assumption that a 3D point landmark  $\mathcal{F}$  with bearing vector  $\boldsymbol{\mu}$  and distance parameter  $\rho$  remains stationary with respect to an inertial frame  $\mathcal{I}$ :

$$\mathcal{I}\mathbf{r}_{\mathcal{I}\mathcal{F}} = \mathcal{I}\mathbf{r}_{\mathcal{I}\mathcal{C}} + \mathbf{q}_{\mathcal{C}\mathcal{I}}^{-1}(\boldsymbol{\mu}d(\rho)). \quad (69)$$

In eq. (68) we can observe that the derivative of the distance parameter only depends on the velocity in direction of the bearing vector. On the other hand, the derivative of the bearing vector, eq. (67), is the sum of a velocity and rotational rate effect, whereby the magnitude of the velocity effect is proportional to the inverse distance of the landmark.

Using an appropriate Euler forward integration scheme, i.e., using the boxplus operator where appropriate, the above time continuous equation can be transformed into a set of discrete prediction equations which are used during the prediction of the filter state. For the attitude, the rotational IMU-camera extrinsics and the bearing vectors the discretization yields:

$$\begin{aligned} \mathbf{q}_{k+1} &= \mathbf{q}_k \boxplus (-\Delta t \mathbf{q}_k(\hat{\boldsymbol{\omega}}_k)), \\ &= \exp(-\Delta t \mathbf{q}_k(\hat{\boldsymbol{\omega}}_k)) \otimes \mathbf{q}_k, \\ &= \mathbf{q}_k \otimes \exp(-\Delta t \hat{\boldsymbol{\omega}}_k) \otimes \mathbf{q}_k^{-1} \otimes \mathbf{q}_k, \\ &= \mathbf{q}_k \otimes \exp(\Delta t(\mathbf{b}_{\omega,k} + \mathbf{w}_{\omega,k} - \hat{\boldsymbol{\omega}}_k)), \end{aligned} \quad (70)$$

$$\begin{aligned} \mathbf{z}_{k+1} &= \mathbf{z}_k \boxplus \Delta t \mathbf{w}_{z,k}, \\ &= \exp(\Delta t \mathbf{w}_{z,k}) \otimes \mathbf{z}_k, \end{aligned} \quad (71)$$

$$\begin{aligned} \boldsymbol{\mu}_{i,k+1} &= \exp\left(\Delta t\left((\mathbf{I} - \mathbf{n}(\boldsymbol{\mu}_{i,k})\mathbf{n}(\boldsymbol{\mu}_{i,k})^T)\hat{\boldsymbol{\omega}}_{\mathcal{C}}\right.\right. \\ &\quad \left.\left.+ \mathbf{n}(\boldsymbol{\mu}_{i,k}) \times \frac{\hat{\mathbf{v}}_{\mathcal{C}}}{d(\rho_{i,k})}\right.\right. \\ &\quad \left.\left.+ \mathbf{N}(\boldsymbol{\mu}_{i,k})\mathbf{w}_{\mu,i}\right)\right) \otimes \boldsymbol{\mu}_{i,k}. \end{aligned} \quad (72)$$

The derivation of the bearing vector discretization is given in Appendix A.

In typical visual-inertial sensor setups, the IMU measurements are often obtained at a higher rate than the images. As the proposed propagation step is driven by the IMU measurements this can result in a high computational burden. A classical approach to mitigate this issue is to make use of IMU pre-integration techniques (Forster et al. (2016)) in order to merge multiple IMU measurements into a single prediction step. However, since the duration between two consecutive images remains relatively small we employ a simpler pre-integration approach where the Jacobian is evaluated based on the mean of the IMU measurement. Thus, even if multiple IMU measurements are available between two consecutive images, eq. (42) is evaluated only once. Compared to the regular solution no notable performance loss could be observed.

#### 5.4 Direct Innovation Term and Update

In section 4.3, we discussed how to construct a photometric error term based on the pixel-wise intensity difference

between a previously extracted patch and its predicted location in a given image frame. Within an IEKF this can be directly used as innovation term. However, for the multilevel patch format that we use, this would lead to a  $6 \times 6 \times 2 = 72$  dimensional error term per patch inducing very high computational cost. Fortunately, looking at eq. (36), one can observe that the entire error term corresponding to a patch  $P_i$  and an image  $I$  is only dependent on the estimated pixel coordinates  $\mathbf{p}_i = \boldsymbol{\pi}(\boldsymbol{\mu}_i)$ . Thus, the only direct filter state dependency of this error term is given by the bearing vector and an equivalent reduced 2D error term can be derived. This can be achieved by means of a QR-decomposition of the gradient matrix in eq. (36):

$$\begin{aligned} \mathbf{A}(\mathbf{p}_i, I, \mathbf{D}_i) &= \mathbf{Q}(\mathbf{p}_i, I, \mathbf{D}_i)\mathbf{R}(\mathbf{p}_i, I, \mathbf{D}_i), \\ &= [\mathbf{Q}_1(\mathbf{p}_i, I, \mathbf{D}_i) \quad \mathbf{Q}_2(\mathbf{p}_i, I, \mathbf{D}_i)] \begin{bmatrix} \mathbf{R}_1(\mathbf{p}_i, I, \mathbf{D}_i) \\ \mathbf{0} \end{bmatrix} \end{aligned} \quad (73)$$

where the upper-triangular matrix  $\mathbf{R}_1(\mathbf{p}_i, I, \mathbf{D}_i)$  has full row-rank 2 for regular features, row-rank 1 for line features, and goes towards  $\mathbf{0}$  for uniform patches.

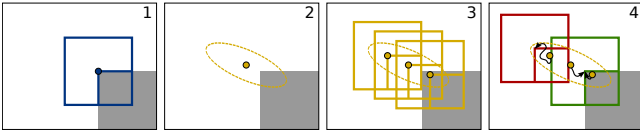
Considering the above decomposition, the innovation term for the  $j^{\text{th}}$  iteration step for a patch  $i$  yields:

$$\mathbf{y}_{i,j} = \mathbf{Q}_1(\boldsymbol{\pi}(\boldsymbol{\mu}_{i,j}^+), I, \mathbf{D}_i)^T \mathbf{b}(\boldsymbol{\pi}(\boldsymbol{\mu}_{i,j}^+), P_i, I, \mathbf{D}_i). \quad (74)$$

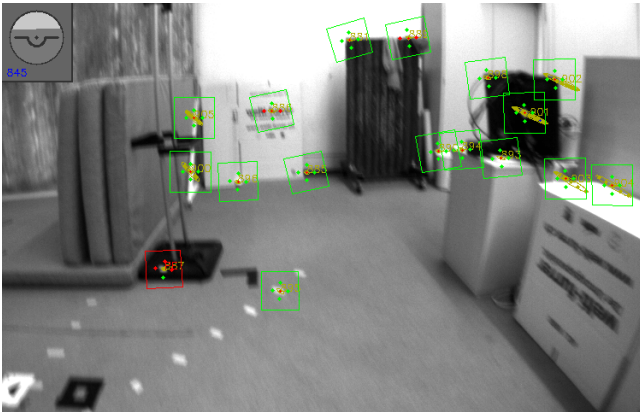
This has a maximal dimension of 2 and loses dimensions for degenerate cases like line-shaped or uniform patches. Since this represents a left-multiplication with an orthonormal matrix, the noise characteristics are assumed to be of the same magnitude on every photometric error term. To account for potentially different noise properties of the intensity errors, a weighting based scheme could be introduced. The Jacobian for the innovation term is given by

$$\mathbf{H}_{i,j} = \mathbf{R}_1(\boldsymbol{\pi}(\boldsymbol{\mu}_{i,j}^+), I, \mathbf{D}_i) \frac{d\boldsymbol{\pi}}{d\boldsymbol{\mu}}(\boldsymbol{\mu}_{i,j}^+). \quad (75)$$

Within the IEKF, the tracked landmarks are updated one after another, each undergoing a certain number of iterations. While the robocentric state formulation moved parts of the nonlinearities from the update into the propagation step, significant nonlinearities remain with the pixel intensity generation. The update iterations are taking care of aligning the patches in the current image and, simultaneously, to spread the gained information throughout the filter state. Thus, all the landmark tracking functionality is intrinsically contained in the filter. In the case where a landmark's predicted image coordinates exhibit a large uncertainty (e.g. for newly initialized landmarks), multiple hypothesis are select within the uncertainty bound. Figure 5 provides a simplified sketch of the tracking concept. An advantage of this is that non-corner features can be properly tracked by considering the prior provided by the IMU-driven process model. In the case of line-shaped features, for instance, a corrective update only applies along the perpendicular direction to the line, while the other direction remains unaffected. In the degenerate case of uniformly textured patch features, the iteration finishes after one step without changes to the filter state (since no information is contained in the patch). Figure 6 shows the tracked landmarks in a frame. Each iteration for a landmark update is depicted by



**Figure 5.** Overview of the landmark tracking concept. Step 1: a patch feature (blue square) is extracted for the landmark (blue dot). Step 2: the estimated landmark image coordinates (yellow dot) and the corresponding covariance (yellow ellipse) are provided by the filter's IMU-driven process model. Step 3: depending on the magnitude of the uncertainty multiple candidates (yellow dots) are initialized. Step 4: for each candidate an iterative update (black arrow) is performed which integrates patch intensity errors together with the motion prior. Outlier detection and quality checks are performed to select the best valid tracking (green vs red squares). Steps 3 and 4 are completely integrated in the iterative filter updates.



**Figure 6.** Live screenshot of the tracked landmarks. The projected patches (final iteration) are depicted by squares (green if successful, red if rejected). The predicted uncertainty of the landmark location are represented by yellow ellipses and each update iteration candidates is marked by a yellow dot. The final location is highlighted with a small red dot surrounded by four green or red dots. The surrounding locations are checked for higher innovation residuals (green). If more than two surrounding locations exhibit no increased innovation residuals (red) then the match is rejected (e.g. the bottom left landmark).

a yellow dot. Especially for the newly added landmarks (the four most right ones), the initial uncertainty (yellow) ellipse and the number of iterations are increased.

To account for moving objects or other disturbances, a simple Mahalanobis based outlier detection is implemented within the update step. It compares the obtained innovation residual with the predicted innovation covariance and rejects the measurement whenever the weighted norm exceeds a certain threshold. This method inherently takes into account the covariance of the state and measurements. It also considers the image gradients and thereby tends to reject gradient-less image patches more readily. In addition, a threshold on the total intensity error of a patch is introduced, whereas a patch measurement is rejected if the threshold is exceeded. Also, a landmark quality check is performed by sampling 4 nearby locations and evaluating the corresponding innovation residual. Tracking tends to be bad if not at least two locations exhibit a significantly higher residual than the matched landmark (see the bottom left landmark in Figure 6).

The computation of the photometric error relies on an image patch from a previous frame. If parts of this image patch have influenced the filter state in the past, the resulting photometric error will exhibit a correlation with the current filter state. This correlation is not modeled in the current framework and doing so would significantly increase the computational burden (one possible approach would be to co-estimate the patch pixel intensities). This is an issue which is also commonly encountered in dense approaches where cross-correlations between localization and mapping are often neglected. In our case however, the cross-correlation with the environment geometry is tracked and accounted for and the problem is limited to the texture of the environment. A refinement step on the patch intensities could reduce the pixel intensity noise and thereby reduce this effect. Investigations in this direction will be part of future work.

## 5.5 Landmark Management

The IEKF does not exhibit good scalability in terms of the size of the filter state. Consequently, only a limited number of landmarks can be tracked and they have to be selected and managed carefully in order to obtain good results. In section 4.4, we outlined an intensity based scoring which describes how informative the content of a patch can be. This is mainly used to decide what landmarks are added to the filter state. In addition to this, we maintain tracking and visibility information of a landmark, and a combined heuristic quality score is computed for each landmark which is being tracked. The quality score is composed of three sub-scores:

- The global quality: how often has a landmark been tracked since initialization
- The local quality: how often has a landmark been tracked when expected to be in the field of view (limited to recent frames)
- The local visibility: how often was the landmark in the field of view (limited to recent frames)

If a landmark exhibits a high global quality, i.e. it has often been tracked since initialization, the pruning threshold on the two local sub-scores is kept more conservative. Using an adaptive thresholding, we can control the total amount of landmarks which are currently in the frame. E.g. if we reach the maximal number of landmarks in the filter state and only a minor part is properly tracked, we make the landmark pruning stricter to get space for new landmarks.

## 6 Multi-Camera Setup

One issue with monocular visual-inertial setups is that they require sufficient motion in order to properly estimate the complete filter state. Also, a particular camera can be blind at times, because of fast lightning changes or very bad texture. Adding an additional camera can therefore improve the robustness of the overall system. In the case where a multi-camera setup has overlapping fields of view, multiple measurements of the same landmark are received at a given time. This provides information about the landmark's distance and the extrinsic calibration of the corresponding camera frames. Still, some excitation of the states is

necessary to estimate the calibration parameters. Once the calibration estimates have converged, the distance of landmarks in overlapping fields of view becomes observable, even if the sensor remains stationary.

New landmarks are still detected in single camera frames only and the parametrization of the corresponding bearing vector and distance parameter is kept with respect to the detection frame. In the case where the newly detected landmark can be seen in more than one camera frame, the initial distance estimate can be computed by triangulation.

For all subsequent time steps, the landmarks get projected in every camera frame. If the predicted pixel coordinates lie within a camera frame, an iterative update is performed (see section 5.4). If the measurement camera frame  $\mathcal{C}_m$  (where the landmark is observed) is not the same as the detection camera frame  $\mathcal{C}_d$  (where the landmark is parametrized), the corresponding bearing vector must be transformed into the measurement frame first. This can be done by

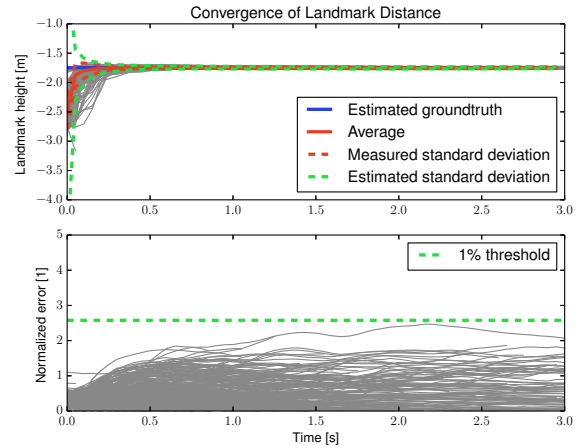
$$c_m \mu_i = z_m(c_d + z_d^{-1}(c_d \mu_i d(\rho_i)) - c_m) \quad (76)$$

where the terms of the form  $c_*$  and  $z_*$  represent the linear and rotational extrinsic IMU-camera calibration of the corresponding camera. Together with the parametrization of the landmark location  $(c_d \mu_i, \rho_i)$ , they are contained in the filter state. This represents the main difference to the monocular setup, whereas the innovation Jacobian in eq. (75) has to be right-multiplied by the Jacobian of eq. (76).

## 7 Experimental Results

The evaluation of the presented approach is split into multiple parts. These include convergence evaluation (section 7.1), parameter exploration (section 7.2), evaluation of the photometric feedback (section 7.3), comparison with other approaches (section 7.4), and tests on a real MAV (section 7.5). Whenever possible, the publicly available EuRoC datasets are employed (Burri et al. (2016)). All experiments are performed with the same selection of parameters except where explicitly mentioned. Our baseline implementation only employs the second and third image pyramid level. While taking into account the first image level can increase accuracy, it is not really useful for highly dynamic and difficult cases. As default, we use a patch size of  $6 \times 6$  together with 25 filter landmarks.

Accuracy remains an important criteria for visual-inertial odometry and can be evaluated quantitatively. To a certain extent, it can also serve as a surrogate measure for the well-functioning of an approach. To evaluate accuracy we contemplate the root mean square estimation error per traveled distance (Geiger et al. (2012)). For instance, if we want to evaluate the accuracy after 10 m of traveled distance and have a dataset which is 80 m long, we split the obtained estimation results into 8 chunks of 10 m. The chunks are then aligned with the corresponding bit of groundtruth data and the accumulated error after 10 m is evaluated. Box-plots are employed to depict the corresponding median and quartiles. Assuming that the odometry output exhibits random walk drift with increasing traveled distance (which is often a good approximation as long as the yaw error remains small), the observed errors should increase as square root of the traveled distance. We select the spacing of the traveled distance



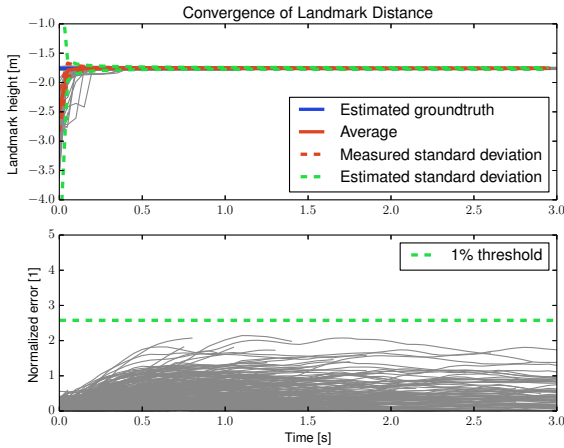
**Figure 7.** Dataset observing a horizontal plane with a *monocular* setup. Estimated landmark heights (grey) together with their average (red) and empirical standard deviation (red dashes,  $1-\sigma$ ). The groundtruth (blue) is estimated by averaging the height estimates once converged. In the top plot, the average of the estimated standard deviation over multiple landmark heights is provided (green dashes,  $1-\sigma$ ). The bottom plot depicts the normalized error (height error divided by estimated standard deviation) together with the 1% confidence threshold.

samples quadratically, and should therefore observe a linear error increase in the plots. The following results can vary depending on the setup and environment and should not be over-interpreted.

### 7.1 Convergence Evaluation

The proposed robocentric formulation allows an undelayed initialization of landmarks but strongly relies on the proper convergence of the corresponding distance estimates (which are initialized with a very high uncertainty). While a decreasing uncertainty is desired since it allows for more accurate tracking of the sensor pose, spurious and inconsistent convergence must be avoided. Especially in the monocular case, the uncertainty should only be decreased if the sensor is moving and sufficient baseline is acquired. In order to investigate the consistency of the distance estimates, a dataset was recorded where a horizontal surface was the only visible object (the camera was directed towards the floor). The virtual groundtruth of the sensor height is inferred by averaging over the height of all converged landmarks (see Figure 7). This allows to evaluate the convergence of the landmarks heights (which are strongly coupled to the distance estimates).

Figures 7 and 8 show the estimated height of the tracked landmarks over time for a monocular and a stereo setup respectively. Since the landmarks are initialized at a fixed distance, which in this experiment tends to relate to points below the surface, a significant estimation error can be observed at initialization. Due to the motion of the sensor, however, the height estimates quickly converge. In the top part of both figures, the estimated standard deviation (average of the estimated standard deviations) is compared against the measured standard deviation (empirical standard deviation of the actual height errors). In both cases, the

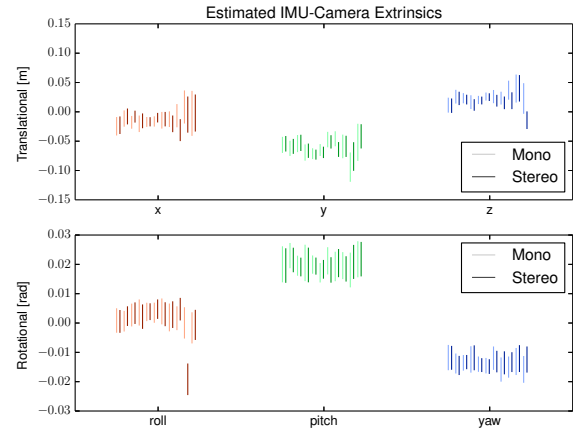


**Figure 8.** Dataset observing a horizontal plane with a stereo setup. Stereo initialization is disabled. Estimated landmark heights (grey) together with their average (red) and empirical standard deviation (red dashes,  $1\text{-}\sigma$ ). The groundtruth (blue) is estimated by averaging the height estimates once converged. In the top plot, the average of the estimated standard deviation over multiple landmark heights is provided (green dashes,  $1\text{-}\sigma$ ). The bottom plot depicts the normalized error together with the 1% confidence threshold.

estimated standard deviation encompasses the measured one. Since we averaged over many landmark tracks, this is not a strict check of consistency but still shows that the estimated covariance must lie within a reasonable range on average. A better analysis is provided in both lower plots which depict the normalized height error (height error divided by estimated standard deviation). In both experiments we can show that the normalized error remains below the 1% confidence threshold, i.e., there are no unreasonably large height estimate errors if compared to the corresponding uncertainty estimates.

The results of the monocular and the stereo setup are similar due to the significant amount of motion present in the recorded data. The final standard deviation of the height errors amounts to 0.0119m for the monocular setup and 0.0073 m for the stereo setup.

Other parameters which have to converge for a proper functioning of the filter are the online calibration parameters, which are composed of the IMU biases and the IMU-camera extrinsics. The later should remain nearly constant for different datasets with the same sensor setup, and we can thus evaluate the extrinsics on multiple datasets and compare the values they have converged to. Figure 9 shows the final estimate of the rotational and translation part of the extrinsics if running the proposed filter on all 11 EuRoC datasets. To make the task more difficult, the initial values were selected as zero translation and closest orthogonal rotation (corresponding to all zero angles in the figure). The resulting estimates, including uncertainties, seem to exhibit a large amount of accordance between the different datasets and between monocular and stereo setup. In comparison to the first half of the datasets, the second half includes datasets with less motion which pose more difficulties for a proper estimation of the extrinsics. Consequently, the estimated uncertainty (length of bar) remains larger as well. In general,



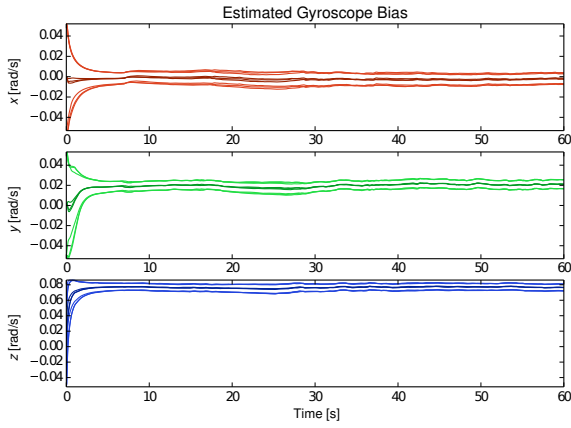
**Figure 9.** Final camera-IMU extrinsics estimates for all 11 EuRoC datasets for monocular and stereo setups. The length of the lines corresponds to the  $3\text{-}\sigma$  bounds of the estimates. The order of the lines corresponds to the datasets V1.01–03, V2.01–03, and MH.01–05.

the stereo setup only brings a marginal reduction of the uncertainties as long as sufficient motion is available. On the contrary, the stereo setup exhibits more difficulties for converging to the proper extrinsic calibration, as can be seen in the lower plot where a single roll angle converges to a biased value. This is due to wrong stereo matches, whereby a single wrong match can bias the estimation, especially if there is not sufficient motion for correcting the wrong convergence.

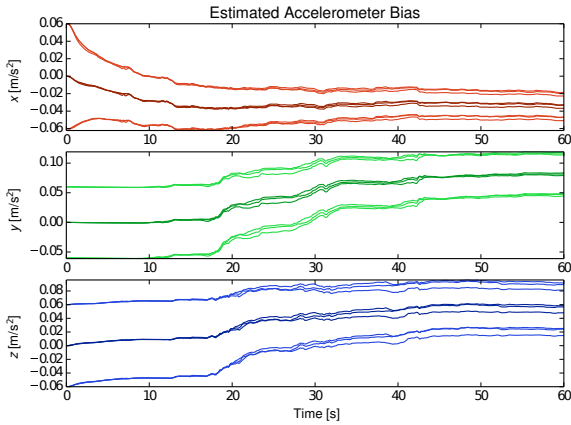
For completeness, we also investigate the convergence of the IMU bias estimation. We only evaluate them on a single dataset (V1.03) since they exhibit intra-dataset variability. Since we have a stereo setup, we can perform two distinct monocular and one stereo evaluation with the same dataset. The results are depicted in Figures 10 and 11, where the estimate over time is plotted together with the  $3\text{-}\sigma$  bounds. In particular, the gyroscope biases seem to converge very rapidly and exhibit only a very small variability. Also the accelerometer biases converge with sufficient excitation of the system. They typically converge faster along the gravity direction which is given by the x-axis at the beginning of the dataset.

## 7.2 Parameter Exploration

Several important framework parameters are evaluated on the EuRoC dataset V1.03. Since the total number of landmarks in the filter state has a major influence on the computational cost of the framework, this is the first parameter that we explore. Figure 12 shows the position error with respect to the traveled distance for different amounts of landmarks. Surprisingly, increasing the number of landmarks does not improve the accuracy of the output once a certain amount of landmarks is reached (roughly 20). While for vision-only systems it has been shown that the number of landmarks is a crucial parameter (Strasdat et al. (2010)), it seems to be different for visual-inertial systems. In visual-inertial systems the IMU provides a good prior on the motion of the systems and merely needs to be stabilized using



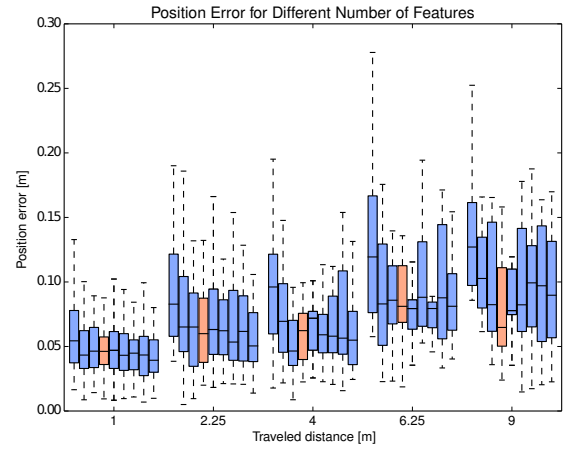
**Figure 10.** Estimated gyroscope biases for dataset V1\_03. Estimates (darker lines) together with the  $3\text{-}\sigma$  bounds (brighter lines). Results for two monocular (left and right camera) and one stereo evaluation are depicted. The estimates converge very quickly and are less motion dependent than the accelerometer biases.



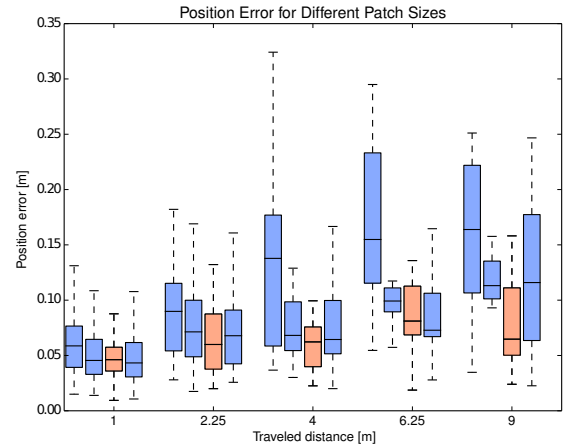
**Figure 11.** Estimated accelerometer biases for dataset V1\_03. Estimates (darker lines) together with the  $3\text{-}\sigma$  bounds (brighter lines). Results for two monocular (left and right camera) and one stereo evaluation are depicted. The accelerometer bias converges quicker along the gravity direction which is mostly along the  $x$ -axis.

recurrent stationary landmarks observations. The amount of required landmarks could be depending on the quality of the employed IMU, whereas an IMU of lower quality would benefit more from higher landmark counts. Within this context we also noticed a relatively strong influence of non-rejected outliers on the output's accuracy, whereas we selected the outlier rejection parameters to be rather strict. We noticed that properly tracking few high-quality landmarks often leads to better results than tracking many landmarks with an increased risk for non-rejected outliers.

In our previous work (Bloesch et al. (2015)) we fixed the patch size to  $8 \times 8$ . Here, we also investigate smaller patch sizes since this reduces the computational load. Results for patch sizes down to  $2 \times 2$  are depicted in Figure 13. It shows that we can reduce the patch size without significantly losing accuracy. Only the case with  $2 \times 2$  patches exhibits



**Figure 12.** Accumulated position error over traveled distance for different landmark counts. The number of landmarks from left to right are 10, 15, 20, 25 (red), 30, 35, 40, 45, 50. The patch size is fixed to 6. The median and the quartiles are depicted.

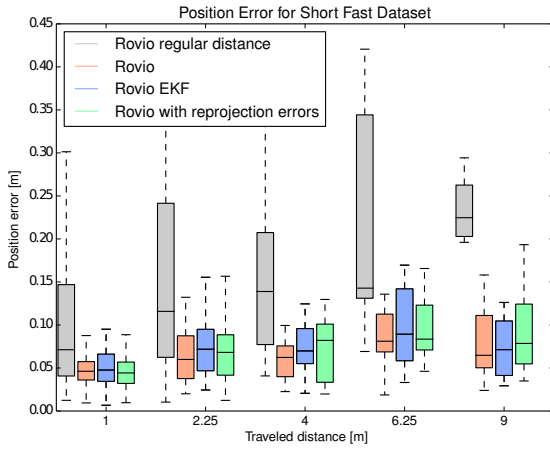


**Figure 13.** Accumulated position error over traveled distance for different patch sizes. The patch sizes from left to right are 2, 4, 6 (red), 8. The landmark count is 25. The median and the quartiles are depicted.

notably increased errors. This could probably be tackled by increasing the amount of pyramid levels which is similar to having larger patches. All in all we propose to employ the combination of patch size  $6 \times 6$  together with 25 filter landmarks which we highlight in red in all box-plots. On a single core of an Intel® Core™ i7 at 2.4 MHz the resulting framework uses 30-50% of the CPU load.

### 7.3 Photometric Feedback Evaluation

The photometric error feedback and its tight integration in the estimation process is a key component of the proposed framework. In order to assess the effect of the photometric feedback we replace it by a traditional reprojection error based feedback combined with an explicit feature tracker. To this end, we implement the KLT tracker mentioned in section 4.3 (including initial guess from the IMU propagation and using patch warping) and try to maintain all settings as similar as possible. For the reprojection error a different measurement covariance is required, which is

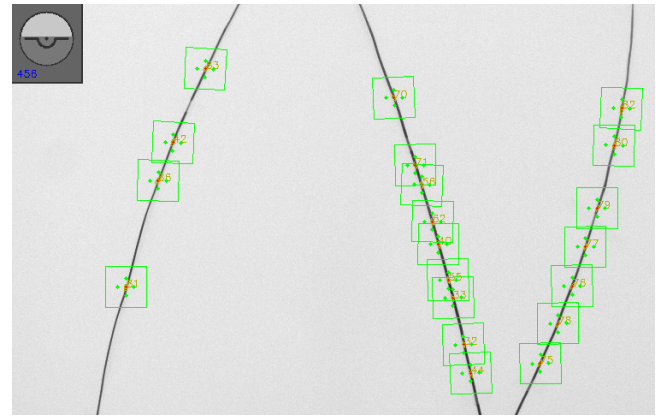


**Figure 14.** Accumulated position error for different Rovio setups on the EuRoC dataset V1.03. The proposed Rovio setup (red) is compared to an implementation with regular distance parametrization (grey), an implementation without update iterations (blue), and an implementation which uses the reprojection error instead of the photometric error feedback (green). While overall the proposed setup exhibits the smallest tracking error, only the setup with the regular distance parametrization shows a significant accuracy loss.

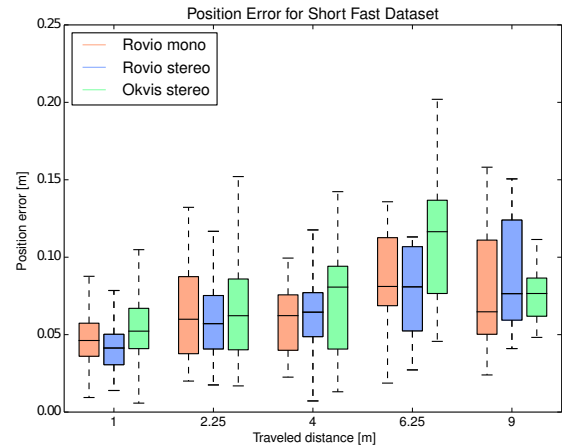
tuned to achieve best performance in terms of accumulated position error on the evaluation dataset (V1.03). The results are depicted in Figure 14. Two related experiments are included in this graph: a comparison with regular distance parametrization and a simple EKF implementation (no update iterations).

An interesting observation is that the reprojection error based implementation does not lead to much larger tracking error. The reason for this are the abundant corner features in the dataset which are relatively well captured by the regular reprojection error. Still, our framework does not rely on any external feature tracker which can be seen as advantage in terms of reduced complexity. Furthermore, in the extreme case where no corner features are available (see Figure 15) the KLT tracker fails and the advantage of the inherent landmark tracking becomes more evident. Failing landmark tracking entails that most sparse state-of-the-art visual-inertial odometry frameworks will be deprived of vision and strongly rely on an IMU-based prediction together with an eventual relocalization. In contrast, Rovio’s implicit landmark tracking is able to keep track of more subtle image regions due to the tight combination with inertial data which can help to overcome visually degenerate sequences.

Finally, we can observe that the use of a regular distance parametrization leads to significantly increased tracking error. This is due to the less accurate stochastic model on the distance when compared to inverse distance parametrization and confirms previous results (Montiel et al. (2006)). The EKF implementation exhibits only a slightly increased error metric. This indicates that a single update is often sufficient for the inherent tracking. This may become more critical if the prediction of the landmark location is less accurate, such as when the initial landmark distance estimate is bad or in cases with high linear velocities.



**Figure 15.** Tracking behavior without strong corner features. Snapshot taken at the end of a 20 s dataset with lines only. The inherent landmark tracking can handle such situations due to the additional prior it receives from the IMU-driven state propagation.

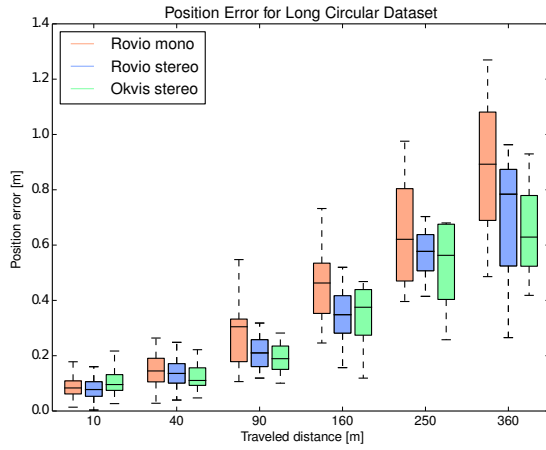


**Figure 16.** Accumulated position error over traveled distance for a monocular and a stereo Rovio setup as well as for Okvis. Red: standard monocular Rovio ( $6 \times 6$  patches, 25 landmarks). Blue: standard stereo Rovio ( $6 \times 6$  patches, 25 landmarks). Green: stereo Okvis with online parameter estimation. All frameworks exhibit similar tracking errors.

## 7.4 Comparison with Related Work

Figure 16 compares both the monocular and stereo Rovio setups against the stereo Okvis framework on the EuRoC dataset V1.03. Okvis is the open-source release version of the work of Leutenegger et al. (2015) and relies on a windowed bundle adjustment approach which includes inertial measurements. Due to the sufficient motion and good texture of the dataset, monocular and stereo setups exhibit very similar tracking errors. For this dataset, the performances of Rovio and Okvis are comparable and show that our approach can compete with state-of-the-art visual-inertial frameworks.

The accuracy of the presented approach was also evaluated on the 1.4 km long circular dataset employed by Leutenegger et al. (2015), Forster et al. (2016), and Usenko et al. (2016). Figures 17 and 18 show the position error and the yaw error over traveled distance for the standard monocular and stereo Rovio setups as well as for Okvis. The performance of Rovio is slightly inferior to Okvis for the 360 m of traveled



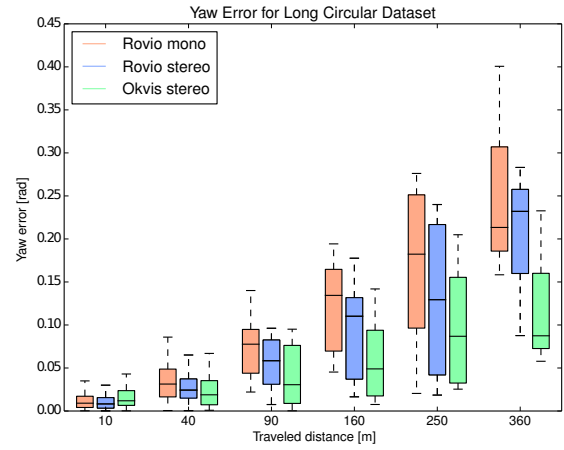
**Figure 17.** Accumulated position error over traveled distance for the long circular dataset. Red: baseline monocular Rovi. Blue: baseline stereo Rovi. Green: stereo Okvis with online parameter estimation.

distance. Again, the performance is strongly depending on the selected tuning parameters which were kept constant for all experiments. Forster et al. (2016) and Usenko et al. (2016) both provide results which show 0.3 m position error after 360 m. Caution should be taken when interpreting these results since only 3 non-overlapping segments of length 360 m are contained in the 1.4 km long dataset (low statistical significance). Rovi seems to perform better for shorter distances and shows similar performance to all other visual-inertial frameworks, especially for a traveled distance of 10 m ( $n=140$ ), where it exhibits a median error of less than 0.1 m. One reason for the decreased long term performance can be found in the increased yaw error which has a strong impact on the position performance for longer distances.

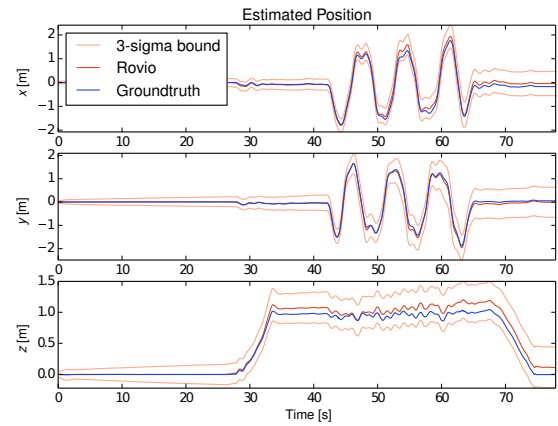
We observed that the above error could be further reduced by choosing stricter outlier rejection parameters and including the first pyramid level into the residual computation (below 0.4 m for 360 m). In the end, however, large scale accuracy should be provided by an enclosing back-end system performing loop closures and re-localization, rather than over-tuning the front-end visual-inertial odometry at the cost of increased computational costs and lower overall robustness.

### 7.5 Robust MAV Control

In this final evaluation section we investigate the applicability of Rovi for feedback control on a MAV for fast aggressive flights under bad lighting conditions and motion blur. The system is initialized on the ground and remains stationary for 30 s. After take off, it performs three fast circular loops before landing at the same location. The trajectory's position and attitude are depicted in Figures 19 and 20, respectively. The  $3\text{-}\sigma$  bounds for the estimates are plotted as well. The observable roll and pitch angles very quickly converge from their initially large uncertainties and accurately track the MAV's inclination angles after take off. The global yaw angle and positions, on the other hand, accumulate uncertainty over time, what confirms the inherent unobservability of those states. All in all, the experiment



**Figure 18.** Accumulated yaw error over traveled distance for the long circular dataset. Red: baseline monocular Rovi. Blue: baseline stereo Rovi. Green: stereo Okvis with online parameter estimation.



**Figure 19.** Estimated MAV position for aggressive flight. Red: estimated position. Light-red:  $3\text{-}\sigma$  bounds. Blue: Vicon groundtruth. The unobservable position accumulates uncertainty over time.

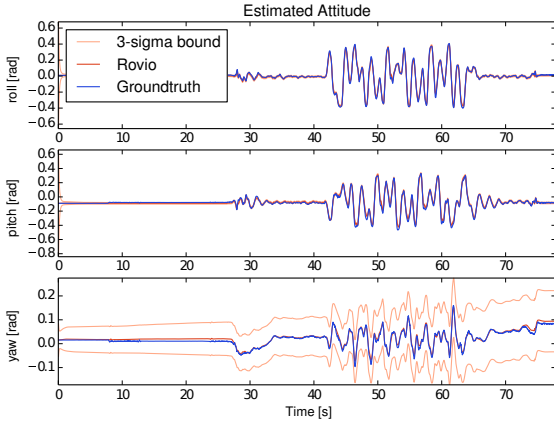
shows that Rovi can handle fast motions and difficult scenes while providing a reliable state estimation for feedback control of an autonomous MAV.

Figure 21 shows the robocentric velocity of the MAV. Due to the robocentric formulation of our filter, the observable states are entirely decoupled from the unobservable states. Hence, the uncertainties are bounded and the estimation error remains minimal which is essential for feedback control. Additionally, velocity estimates are provided where Rovi was reset every 5 s. The estimates very quickly converge to the true velocities for all resets. This highlights the very simple and robust initialization of our robocentric filtering approach where ego-motion estimates are immediately available.

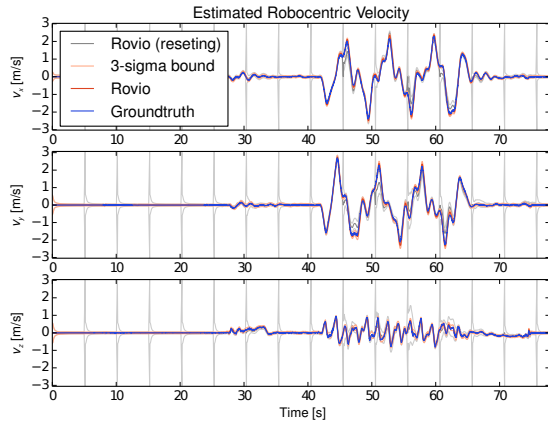
## 8 Conclusion

This paper presented an IEKF-based framework which tightly fuses inertial measurements and image data. The





**Figure 20.** Estimated MAV attitude for aggressive flight. The yaw-pitch-roll decomposition is employed to separate the unobservable yaw from the two inclination angles (only visualization). Red: estimated attitude. Light-red: 3- $\sigma$  bounds. Blue: Vicon groundtruth.



**Figure 21.** Estimated MAV velocity for aggressive flight. Red: estimated velocity. Light-red: 3- $\sigma$  bounds. Blue: Vicon groundtruth. Grey: estimated velocity when resetting Rovio every 5 s. Since the velocity is expressed in the robocentric IMU coordinate frame it is fully observable and the uncertainty remains bounded.

originality and strength of the proposed approach lie in its fully robocentric formulation combined with the direct feedback of photometric error as the Kalman innovation term. This leads to a more robust implementation, since the observable states are not influenced by the growing global covariance. We introduce an iterative update scheme which inherently takes care of landmark tracking. While simplifying the structure of the overall framework, data association is robustified by the tight coupling with the IMU-driven process model. The employed minimal representations of rotations and bearing vectors improve the numerical consistency of the approach and reduce the computational effort. The extensive experimental evaluation shows that the presented approach can compete with state-of-the-art visual-inertial fusion techniques. Interestingly, our approach achieves comparable ego-motion estimation accuracy with a significantly lower landmark count.

Robustness with respect to fast motions and bad lightning conditions as well as the instantaneous initialization procedure were demonstrated in a real autonomous MAV flight experiment. Additional features, such as optional GPS measurements, are included in the updated version of the publicly available open-source software package <sup>\*</sup>.

## Funding

This work was supported by the Swiss National Science Foundation (grant number 200021\_149427/1); the National Centre of Competence in Research Robotics; and the European Community's Seventh Framework Programme (grant number n.608849).

## Appendix A Bearing Vector Calculus

Assuming a stationary 3D point landmark  $\mathcal{F}$  with bearing vector  $\boldsymbol{\mu}$  and distance parameter  $\rho$ , the corresponding differential equations can be obtained by totally differentiating the kinematics:

$$\frac{d}{dt} \left\{ {}^I \mathbf{r}_{IF} = {}^I \mathbf{r}_{IC} + \mathbf{q}_{CI}^{-1}(\boldsymbol{\mu} d(\rho)) \right\}. \quad (77)$$

For this we require the following partial differentials:

$$\frac{d}{dt} ({}^I \mathbf{r}_{IC}) = {}^I \mathbf{v}_C, \quad (78)$$

$$\frac{\partial}{\partial \mathbf{q}_{CI}} (\mathbf{q}_{CI}^{-1}(\boldsymbol{\mu} d(\rho))) = -\mathbf{q}_{CI}^{-1}(\boldsymbol{\mu} d(\rho)) \times \mathbf{C}(\mathbf{q}_{CI}^{-1}), \quad (79)$$

$$= -\mathbf{C}(\mathbf{q}_{CI}^{-1}) \boldsymbol{\mu} \times d(\rho), \quad (80)$$

$$\frac{d}{dt} \mathbf{q}_{CI} = \boldsymbol{\omega}_C, \quad (81)$$

$$\frac{\partial}{\partial \boldsymbol{\mu}} (\mathbf{q}_{CI}^{-1}(\boldsymbol{\mu} d(\rho))) = \mathbf{C}(\mathbf{q}_{CI}^{-1}) \boldsymbol{\mu} \times \mathbf{N}(\boldsymbol{\mu}) d(\rho), \quad (82)$$

$$\frac{\partial}{\partial \rho} (\mathbf{q}_{CI}^{-1}(\boldsymbol{\mu} d(\rho))) = \mathbf{C}(\mathbf{q}_{CI}^{-1}) \boldsymbol{\mu} d'(\rho). \quad (83)$$

In eq. (79) we use the chain rule together with eqs. (14) and (15), in eq. (81) we directly employed eq. (13), and eq. (82) relies on eq. (32). The total differential of eq. (77) can be evaluated and simplified to (left multiplication with  $\mathbf{C}(\mathbf{q}_{CI})$ ):

$$0 = {}^I \mathbf{v}_C - \mathbf{C}(\mathbf{q}_{CI}^{-1}) \boldsymbol{\mu} \times \boldsymbol{\omega}_C d(\rho) + \mathbf{C}(\mathbf{q}_{CI}^{-1}) (\boldsymbol{\mu} \times \mathbf{N}(\boldsymbol{\mu}) \dot{\boldsymbol{\mu}} d(\rho) + \boldsymbol{\mu} d'(\rho) \dot{\rho}), \quad (84)$$

$$0 = \mathbf{v}_C - \boldsymbol{\mu} \times \boldsymbol{\omega}_C d(\rho) + \boldsymbol{\mu} \times \mathbf{N}(\boldsymbol{\mu}) \dot{\boldsymbol{\mu}} d(\rho) + \boldsymbol{\mu} d'(\rho) \dot{\rho}. \quad (85)$$

From this the dynamics for the bearing vector and distance parameter can be obtained by pre-multiplication with  $1/d(\rho) \mathbf{N}(\boldsymbol{\mu})^T \boldsymbol{\mu} \times$  and  $1/d'(\rho) \boldsymbol{\mu}^T$  respectively:

$$\dot{\boldsymbol{\mu}} = \mathbf{N}(\boldsymbol{\mu})^T \left( \hat{\boldsymbol{\omega}}_C + \mathbf{n}(\boldsymbol{\mu}) \times \frac{\hat{\mathbf{v}}_C}{d(\rho)} \right) + \mathbf{w}_\mu, \quad (86)$$

$$\dot{\rho} = -\mathbf{n}(\boldsymbol{\mu})^T \hat{\mathbf{v}}_C / d'(\rho) + w_\rho. \quad (87)$$

Here we used the identities  $\mathbf{N}(\boldsymbol{\mu})^T \boldsymbol{\mu} \times \boldsymbol{\mu} \times = -\mathbf{N}(\boldsymbol{\mu})^T$  and  $\mathbf{N}(\boldsymbol{\mu})^T \mathbf{N}(\boldsymbol{\mu}) = \mathbf{I}$ . Also, some additive process noise has been added.

<sup>\*</sup> <https://github.com/ethz-asl/rovio>

Applying the Euler-forward discretization scheme on the continuous time differential equation of the bearing vectors (67) yields:

$$\begin{aligned}\boldsymbol{\mu}_{k+1} &= \boldsymbol{\mu}_k \boxplus \Delta t \left( \mathbf{N}(\boldsymbol{\mu}_k)^T \left( \hat{\boldsymbol{\omega}}_c + \mathbf{n}(\boldsymbol{\mu}_k) \times \frac{\hat{\mathbf{v}}_c}{d(\rho_k)} \right) + \mathbf{w}_\mu \right), \\ \boldsymbol{\mu}_{k+1} &= \exp \left( \Delta t \mathbf{N}(\boldsymbol{\mu}_k) \left( \mathbf{N}(\boldsymbol{\mu}_k)^T \left( \hat{\boldsymbol{\omega}}_c \right. \right. \right. \\ &\quad \left. \left. \left. + \mathbf{n}(\boldsymbol{\mu}_k) \times \frac{\hat{\mathbf{v}}_c}{d(\rho_k)} \right) + \mathbf{w}_\mu \right) \right) \otimes \boldsymbol{\mu}_k, \\ \boldsymbol{\mu}_{k+1} &= \exp \left( \Delta t \left( (\mathbf{I} - \mathbf{n}(\boldsymbol{\mu}_k) \mathbf{n}(\boldsymbol{\mu}_k)^T) \hat{\boldsymbol{\omega}}_c \right. \right. \\ &\quad \left. \left. + \mathbf{n}(\boldsymbol{\mu}_k) \times \frac{\hat{\mathbf{v}}_c}{d(\rho_k)} \right. \right. \\ &\quad \left. \left. + \mathbf{N}(\boldsymbol{\mu}_k) \mathbf{w}_\mu \right) \right) \otimes \boldsymbol{\mu}_k.\end{aligned}\quad (88)$$

Here we applied the definition of boxplus (25) and used the identity  $\mathbf{N}(\boldsymbol{\mu})\mathbf{N}(\boldsymbol{\mu})^T = \mathbf{I} - \mathbf{n}(\boldsymbol{\mu})\mathbf{n}(\boldsymbol{\mu})^T$ . The three components influencing the bearing vector prediction can be observed here: the perpendicular part of the rotational rate, the linear velocity weighted by the inverse distance, and the additive noise.

## References

- Bell B and Cathey F (1993) The iterated Kalman filter update as a Gauss-Newton method. *IEEE Transactions on Automatic Control* 38(2): 1991–1994.
- Bloesch M, Gehring C, Fankhauser P, Hutter M, Hoepflinger MA and Siegwart R (2013) State Estimation for Legged Robots on Unstable and Slippery Terrain. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*. pp. 6058–6064. DOI:10.1109/IROS.2013.6697236.
- Bloesch M, Omari S, Hutter M and Siegwart R (2015) Robust Visual Inertial Odometry Using a Direct EKF-Based Approach. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*. pp. 298–304.
- Bloesch M, Sommer H, Laidlow T, Burri M, Nützi G, Fankhauser P, Bellicoso D, Gehring C, Leutenegger S, Hutter M and Siegwart R (2016) A Primer on the Differential Calculus of 3D Orientations. *CoRR* abs/1606.0.
- Burri M, Nikolic J, Gohl P, Schneider T, Rehder J, Omari S, Achtelik MW and Siegwart R (2016) The EuRoC micro aerial vehicle datasets. *The International Journal of Robotics Research* DOI:10.1177/0278364915620033.
- Castellanos JA, Neira J and Tardos JD (2004) Limits to the consistency of EKF-based SLAM. In: *IFAC Symposium on Intelligent Autonomous Vehicles*. DOI:10.1109/TAC.2000.880989.
- Civera J, Grasa OG, Davison AJ and Montiel JMM (2009) 1-point RANSAC for EKF-based Structure from Motion. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*. pp. 3498–3504. DOI:10.1109/IROS.2009.5354410.
- Davison AJ (2003) Real-Time Simultaneous Localisation and Mapping with a Single Camera. In: *IEEE International Conference on Computer Vision*. pp. 1403–1410 vol.2. DOI: 10.1109/ICCV.2003.1238654.
- Engel J, Schöps T and Cremers D (2014) LSD-SLAM: Large-Scale Direct Monocular SLAM. In: *European Conference on Computer Vision*. pp. 834–849.
- Forster C, Carlone L, Dellaert F and Scaramuzza D (2016) On-Manifold Preintegration Theory for Fast and Accurate Visual-Inertial Navigation .
- Forster C, Pizzoli M and Scaramuzza D (2014) SVO : Fast Semi-Direct Monocular Visual Odometry. In: *IEEE International Conference on Robotics and Automation*. DOI:10.1109/ICRA.2014.6906584.
- Geiger A, Lenz P and Urtasun R (2012) Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite. In: *Computer Vision and Pattern Recognition*.
- Hertzberg C, Wagner R, Frese U and Schröder L (2011) Integrating generic sensor fusion algorithms with sound state representations through encapsulation of manifolds. *Information Fusion* 14(1): 57–77.
- Huang GP, Mourikis AI and Roumeliotis SI (2008) Analysis and improvement of the consistency of extended Kalman filter based SLAM. In: *IEEE International Conference on Robotics and Automation*. pp. 473–479. DOI:10.1109/ROBOT.2008.4543252.
- Jin H, Favaro P and Soatto S (2003) A semi-direct approach to structure from motion. In: *The Visual Computer*, volume 19. pp. 377–394. DOI:10.1007/s00371-003-0202-6.
- Jones ES and Soatto S (2011) Visual-inertial navigation, mapping and localization: A scalable real-time causal approach. *The International Journal of Robotics Research* 30(4): 407–430. DOI:10.1177/0278364910388963.
- Julier SJ and Uhlmann JK (2001) A counter example to the theory of simultaneous localization and map building. In: *IEEE International Conference on Robotics and Automation*. pp. 4238–4243. DOI:10.1109/ROBOT.2001.933280.
- Kelly J and Sukhatme GS (2011) Visual-Inertial Sensor Fusion: Localization, Mapping and Sensor-to-Sensor Self-calibration. *The International Journal of Robotics Research* 30(1): 56–79. DOI:10.1177/0278364910382802.
- Leutenegger S, Lynen S, Bosse M, Siegwart R and Furgale P (2015) Keyframe-based visualinertial odometry using nonlinear optimization. *The International Journal of Robotics Research* 34(3): 314–334. DOI:10.1177/0278364914554813.
- Li M and Mourikis AI (2013) High-precision, consistent EKF-based visual-inertial odometry. *The International Journal of Robotics Research* 32(6): 690–711. DOI:10.1177/0278364913481251.
- Lucas BD and Kanade T (1981) An Iterative Image Registration Technique with an Application to Stereo Vision. In: *International Joint Conference on Artificial Intelligence*. pp. 674–679.
- Lynen S, Sattler T, Bosse M, Hesch J, Pollefeys M and Siegwart R (2015) Get Out of My Lab: Large-scale, Real-Time Visual-Inertial Localization. In: *Proceedings of Robotics: Science and Systems*. DOI:10.15607/RSS.2015.XI.037.
- Ma J, Bajracharya M, Susca S, Matthies L and Malchano M (2015) Real-time pose estimation of a dynamic quadruped in GPS-denied environments for 24-hour operation. *The International Journal of Robotics Research* DOI:10.1177/0278364915587333.
- Molton N, Davison A and Reid I (2004) Locally Planar Patch Features for Real-Time Structure from Motion. In: *British Machine Vision Conference*. pp. 90.1–90.10.
- Montiel J, Civera J and Davison A (2006) Unified Inverse Depth Parametrization for Monocular SLAM. In: *Proceedings of*

*Robotics: Science and Systems.*

- Mourikis AI and Roumeliotis SI (2007) A Multi-State Constraint Kalman Filter for Vision-aided Inertial Navigation. In: *IEEE International Conference on Robotics and Automation*. pp. 3565–3572. DOI:10.1109/ROBOT.2007.364024.
- Rosten E and Drummond T (2006) Machine Learning for High Speed Corner Detection. *Computer Vision – ECCV 2006* 1: 430–443. DOI:10.1007/11744023\_34.
- Shen S, Mulgaonkar Y, Michael N and Kumar V (2014) Initialization-free monocular visual-inertial estimation with application to autonomous MAVs. In: *International Symposium on Experimental Robotics*.
- Shi J and Tomasi C (1994) Good features to track. In: *Computer Vision and Pattern Recognition*. pp. 593–600. DOI:10.1109/CVPR.1994.323794.
- Silveira G, Malis E and Rives P (2008) An Efficient Direct Approach to Visual SLAM. *IEEE Transactions on Robotics* 24(5): 969–979. DOI:10.1109/TRO.2008.2004829.
- Solà J, Vidal-Calleja T, Civera J and Montiel JMM (2012) Impact of landmark parametrization on monocular EKF-SLAM with points and lines. *International Journal of Computer Vision* 97: 339–368. DOI:10.1007/s11263-011-0492-5.
- Stelzer A, Hirschmuller H and Gornier M (2012) Stereo-vision-based navigation of a six-legged walking robot in unknown rough terrain. *The International Journal of Robotics Research* 31(4): 381–402. DOI:10.1177/0278364911435161.
- Strasdat H, Montiel JMM and Davison AJ (2010) Real-time monocular SLAM: Why filter? In: *IEEE International Conference on Robotics and Automation*. pp. 2657–2664. DOI: 10.1109/ROBOT.2010.5509636.
- Tanskanen P, Naegeli T, Pollefeys M and Hilliges O (2015) Semi-Direct EKF-based Monocular Visual-Inertial Odometry. In: *IEEE/RSJ International Conference on Intelligent Robots and Systems*.
- Usenko V, Engel J, Stückler J and Cremers D (2016) Direct Visual-Inertial Odometry with Stereo Cameras. In: *Proceedings - IEEE International Conference on Robotics and Automation*.
- Weiss S, Achtelik M, Lynen S, Kneip L, Chli M and Siegwart R (2013) Monocular Vision for Long-term Micro Aerial Vehicle State Estimation: A Compendium. *Journal of Field Robotics* 30(5): 803–831.