

Camera-Based Navigation of a Low-Cost Quadrocopter

Jakob Engel, Jürgen Sturm, Daniel Cremers

Abstract—In this paper, we describe a system that enables a low-cost quadrocopter coupled with a ground-based laptop to navigate autonomously in previously unknown and GPS-denied environments. Our system consists of three components: a monocular SLAM system, an extended Kalman filter for data fusion and state estimation and a PID controller to generate steering commands. Next to a working system, the main contribution of this paper is a novel, closed-form solution to estimate the absolute scale of the generated visual map from inertial and altitude measurements. In an extensive set of experiments, we demonstrate that our system is able to navigate in previously unknown environments at absolute scale without requiring artificial markers or external sensors. Furthermore, we show (1) its robustness to temporary loss of visual tracking and significant delays in the communication process, (2) the elimination of odometry drift as a result of the visual SLAM system and (3) accurate, scale-aware pose estimation and navigation.

I. INTRODUCTION

In recent years, research interest in autonomous micro-aerial vehicles (MAVs) has grown rapidly. Significant progress has been made, recent examples include aggressive flight maneuvers [1, 2], ping-pong [3] and collaborative construction tasks [4]. However, all of these systems require external motion capture systems. Flying in unknown, GPS-denied environments is still an open research problem. The key challenges here are to localize the robot purely from its own sensor data and to robustly navigate it even under potential sensor loss. This requires both a solution to the so-called simultaneous localization and mapping (SLAM) problem as well as robust state estimation and control methods. These challenges are even more expressed on low-cost hardware with inaccurate actuators, noisy sensors, significant delays and limited onboard computation resources.

For solving the SLAM problem on MAVs, different types of sensors such as laser range scanners [5], monocular cameras [6, 7], stereo cameras [8] and RGB-D sensors [9] have been explored in the past. In our point of view, monocular cameras provide two major advantages above other modalities: (1) the amount of information that can be acquired is immense compared to their low weight, power consumption, size and cost, which are unmatched by any other type of sensor and (2) in contrast to depth measuring devices, the range of a monocular camera is virtually unlimited – allowing a monocular SLAM system to operate both in small, confined and large, open environments. The drawback however is, that the scale of the environment cannot be determined from monocular vision alone, such that additional sensors (such as an IMU) are required.

J. Engel, J. Sturm and D. Cremers are with the Department of Computer Science, Technical University of Munich, Germany {engelj, sturmju, cremers}@in.tum.de



Fig. 1. A low-cost quadrocopter navigates in unstructured environments using the front camera as its main sensor. The quadrocopter is able to hold a position, fly figures with absolute scale, and recover from temporary tracking loss. Picture taken at the TUM open day.

The motivation behind our work is to showcase that robust, scale-aware visual navigation is feasible and safe on low-cost robotic hardware. As a platform, we use the Parrot AR.Drone which is available for \$300 and, with a weight of only 420 g and a protective hull, safe to be used in public places (see Fig. 1). As the onboard computational resources are utterly limited, all computations are performed externally.

The contribution of this paper is two-fold: first, we derive a maximum-likelihood estimator to determine the map scale in closed-form from metric distance measurements. Second, we provide a robust technique to deal with large delays in the controlled system, which facilitates the use of a ground station in the control loop. Two videos demonstrating the robustness of our approach, its ability to eliminate drift effectively and to estimate the absolute scale of the map are available online:

<http://youtu.be/tZx1Dly7lno>

<http://youtu.be/eznMokFQmpc>

II. RELATED WORK

Previous work on autonomous flight with quadrocopters can be categorized into different research areas. One part of the community focuses on accurate quadrocopter control and a number of impressive results have been published [10, 1, 3]. These works however rely on advanced external tracking systems, restricting their use to a lab environment. A similar approach is to distribute artificial markers in the environment, simplifying pose estimation [11]. Other approaches learn a map offline from a previously recorded, manual flight and thereby enable a quadrocopter to again fly the same trajectory [12]. For outdoor flights where GPS-based pose estimation is possible, complete solutions are available as commercial products [13].

In this work we focus on autonomous flight without previous knowledge about the environment nor GPS signals, while using only onboard sensors. First results towards this goal have been presented using a lightweight laser scanner [5], a Kinect [9] or a stereo rig [8] mounted on a quadcopter as primary sensor. While these sensors provide absolute scale of the environment, their drawback is a limited range and large weight, size and power consumption when compared to a monocular setup [14, 7].

In our work we therefore focus on a monocular camera for pose estimation. Stabilizing controllers based on optical flow were presented in [15], and similar methods are integrated in commercially available hardware [16]. Such systems however are subject to drift over time, and hence not suited for long-term navigation.

To eliminate drift, various monocular SLAM methods have been investigated on quadcopters, both with off-board [14, 5] and on-board processing [7]. A particular challenge for monocular SLAM is, that the scale of the map needs to be estimated from additional metric sensors such as IMU or GPS, as it cannot be recovered from vision alone. This problem has been addressed in recent publications such as [17, 18]. The current state of the art is to estimate the scale using an extended Kalman filter (EKF), which contains scale and offset in its state. In contrast to this, we propose a novel approach which is based on direct computation: Using a statistical formulation, we derive a closed-form, consistent estimator for the scale of the visual map. Our method yields accurate results both in simulation and practice, and requires less computational resources than filtering. It can be used with any monocular SLAM algorithm and sensors providing metric position or velocity measurements, such as an ultrasonic or pressure altimeter or occasional GPS measurements.

In contrast to the systems presented in [14, 7], we deliberately refrain from using expensive, customized hardware: the only hardware required is the AR.Drone, which comes at a costs of merely \$300 – a fraction of the cost of quadcopters used in previous work. Released in 2010 and marketed as high-tech toy, it has been used and discussed in several research projects [19, 20, 21]. To our knowledge, we are the first to present a complete implementation of autonomous, camera-based flight in unknown, unstructured environments using the AR.Drone.

III. HARDWARE PLATFORM

As platform we use the Parrot AR.Drone, a commercially available quadcopter. Compared to other modern MAV's such as Ascending Technology's Pelican or Hummingbird quadcopters, its main advantage is the very low price, its robustness to crashes and the fact that it can safely be used indoor and close to people. This however comes at the price of flexibility: Neither the hardware itself nor the software running onboard can easily be modified, and communication with the quadcopter is only possible over wireless LAN. With battery and hull, the AR.Drone measures 53 cm \times 52 cm and weights 420 g.

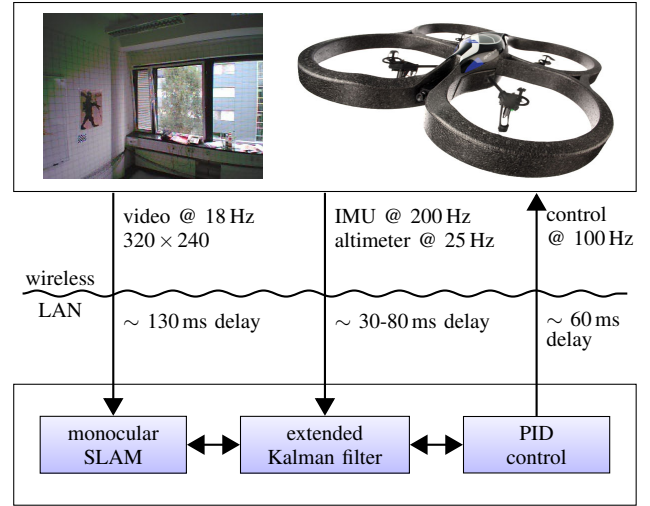


Fig. 2. Approach Outline: Our navigation system consists of three major components: a monocular SLAM implementation for visual tracking, an EKF for data fusion and prediction, and PID control for pose stabilization and navigation. All computations are performed offboard, which leads to significant, varying delays which our approach has to compensate.

A. Sensors

The AR.Drone is equipped with a 3-axis gyroscope and accelerometer, an ultrasound altimeter and two cameras. The first camera is aimed forward, covers a field of view of $73.5^\circ \times 58.5^\circ$, has a resolution of 320×240 and a rolling shutter with a delay of 40 ms between the first and the last line captured. The video of the first camera is streamed to a laptop at 18 fps, using lossy compression. The second camera aims downward, covers a field of view of $47.5^\circ \times 36.5^\circ$ and has a resolution of 176×144 at 60 fps. The onboard software uses the down-looking camera to estimate the horizontal velocity. The quadcopter sends gyroscope measurements and the estimated horizontal velocity at 200 Hz, the ultrasound measurements at 25 Hz to the laptop. The raw accelerometer data cannot be accessed directly.

B. Control

The onboard software uses these sensors to control the roll Φ and pitch Θ , the yaw rotational speed $\dot{\Psi}$ and the vertical velocity \dot{z} of the quadcopter according to an external reference value. This reference is set by sending a new control command $\mathbf{u} = (\bar{\Phi}, \bar{\Theta}, \bar{\dot{z}}, \bar{\dot{\Psi}}) \in [-1, 1]^4$ every 10 ms.

IV. APPROACH

Our approach consists of three major components running on a laptop connected to the quadcopter via wireless LAN, an overview is given in Fig. 2.

1) **Monocular SLAM:** For monocular SLAM, our solution is based on Parallel Tracking and Mapping (PTAM) [22]. After map initialization, we rotate the visual map such that the xy -plane corresponds to the horizontal plane according to the accelerometer data, and scale it such that the average keypoint depth is 1. Throughout tracking, the scale of the map $\lambda \in \mathbb{R}$ is estimated using a novel method described in Section IV-A. Furthermore, we use the pose estimates from the EKF to identify and reject falsely tracked frames.

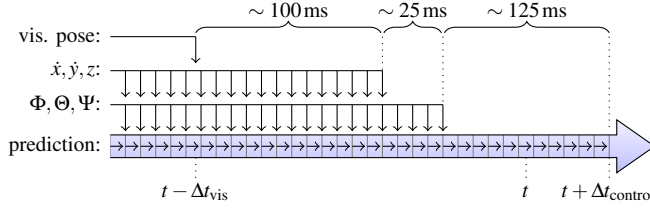


Fig. 3. Pose Prediction: Measurements and control commands arrive with significant delays. To compensate for these delays, we keep a history of observations and sent control commands between $t - \Delta t_{\text{vis}}$ and $t + \Delta t_{\text{control}}$ and re-calculate the EKF state when required. Note the large timespan with no or only partial odometry observations.

2) **Extended Kalman Filter:** In order to fuse all available data, we employ an extended Kalman filter (EKF). We derived and calibrated a full motion model of the quadcopter's flight dynamics and reaction to control commands, which we will describe in more detail in Section IV-B. This EKF is also used to compensate for the different time delays in the system, arising from wireless LAN communication and computationally complex visual tracking.

We found that height and horizontal velocity measurements arrive with the same delay, which is slightly larger than the delay of attitude measurements. The delay of visual pose estimates Δt_{vis} is by far the largest. Furthermore we account for the time required by a new control command to reach the drone $\Delta t_{\text{control}}$. All timing values given subsequently are typical values for a good connection, the exact values depend on the wireless connection quality and are determined by a combination of regular ICMP echo requests sent to the quadcopter and calibration experiments.

Our approach works as follows: first, we time-stamp all incoming data and store it in an observation buffer. Control commands are then calculated using a prediction for the quadcopter's pose at $t + \Delta t_{\text{control}}$. For this prediction, we start with the saved state of the EKF at $t - \Delta t_{\text{vis}}$ (i.e., after the last visual observation/unsuccessfully tracked frame). Subsequently, we predict ahead up to $t + \Delta t_{\text{control}}$, using previously issued control commands and integrating stored sensor measurements as observations. This is illustrated in Fig. 3. With this approach, we are able to compensate for delayed and missing observations at the expense of recalculating the last cycles of the EKF.

3) **PID Control:** Based on the position and velocity estimates from the EKF at $t + \Delta t_{\text{control}}$, we apply PID control to steer the quadcopter towards the desired goal location $\mathbf{p} = (\hat{x}, \hat{y}, \hat{z}, \hat{\Psi})^T \in \mathbb{R}^4$ in a global coordinate system. According to the state estimate, we rotate the generated control commands to the robot-centric coordinate system and send them to the quadcopter. For each of the four degrees-of-freedom, we employ a separate PID controller for which we experimentally determined suitable controller gains.

A. Scale Estimation

One of the key contributions of this paper is a closed-form solution for estimating the scale $\lambda \in \mathbb{R}^+$ of a monocular SLAM system. For this, we assume that the robot is able to make noisy measurements of absolute distances or veloci-

ties from additional, metric sensors such as an ultrasound altimeter.

As a first step, the quadcopter measures in regular intervals the d -dimensional distance traveled both using only the visual SLAM system (subtracting start and end position) and using only the metric sensors available (subtracting start and end position, or integrating over estimated speeds). Each interval gives a pair of samples $\mathbf{x}_i, \mathbf{y}_i \in \mathbb{R}^d$, where \mathbf{x}_i is scaled according to the visual map and \mathbf{y}_i is in metric units. As both \mathbf{x}_i and \mathbf{y}_i measure the motion of the quadcopter, they are related according to $\mathbf{x}_i \approx \lambda \mathbf{y}_i$.

More specifically, if we assume Gaussian noise in the sensor measurements with constant variance¹, we obtain

$$\mathbf{x}_i \sim \mathcal{N}(\lambda \boldsymbol{\mu}_i, \sigma_x^2 \mathbf{I}_{3 \times 3}) \quad (1)$$

$$\mathbf{y}_i \sim \mathcal{N}(\boldsymbol{\mu}_i, \sigma_y^2 \mathbf{I}_{3 \times 3}) \quad (2)$$

where the $\boldsymbol{\mu}_i \in \mathbb{R}^d$ denote the true (unknown) distances covered and $\sigma_x^2, \sigma_y^2 \in \mathbb{R}^+$ the variances of the measurement errors. Note that the individual $\boldsymbol{\mu}_i$ are not constant but depend on the actual distances traveled by the quadcopter in the measurement intervals.

One possibility to estimate λ is to minimize the sum of squared differences (SSD) between the re-scaled measurements, i.e., to compute one of the following:

$$\lambda_y^* := \arg \min_{\lambda} \sum_i \|\mathbf{x}_i - \lambda \mathbf{y}_i\|^2 = \frac{\sum_i \mathbf{x}_i^T \mathbf{y}_i}{\sum_i \mathbf{y}_i^T \mathbf{y}_i} \quad (3)$$

$$\lambda_x^* := \left(\arg \min_{\lambda} \sum_i \|\lambda \mathbf{x}_i - \mathbf{y}_i\|^2 \right)^{-1} = \frac{\sum_i \mathbf{x}_i^T \mathbf{x}_i}{\sum_i \mathbf{x}_i^T \mathbf{y}_i}. \quad (4)$$

The difference between these two lines is whether one aims at scaling the \mathbf{x}_i to the \mathbf{y}_i or vice versa. However, both approaches lead to different results, none of which converges to the true scale λ when adding more samples. To resolve this, we propose a maximum likelihood (ML) approach, that is estimating λ by minimizing the negative log-likelihood

$$\mathcal{L}(\boldsymbol{\mu}_1 \dots \boldsymbol{\mu}_n, \lambda) \propto \frac{1}{2} \sum_{i=1}^n \left(\frac{\|\mathbf{x}_i - \lambda \boldsymbol{\mu}_i\|^2}{\sigma_x^2} + \frac{\|\mathbf{y}_i - \boldsymbol{\mu}_i\|^2}{\sigma_y^2} \right) \quad (5)$$

By first minimizing over the $\boldsymbol{\mu}_i$ and then over λ , it can be shown analytically that (5) has a unique, global minimum at

$$\boldsymbol{\mu}_i^* = \frac{\lambda^* \sigma_y^2 \mathbf{x}_i + \sigma_x^2 \mathbf{y}_i}{\lambda^{*2} \sigma_y^2 + \sigma_x^2} \quad (6)$$

$$\lambda^* = \frac{s_{xx} - s_{yy} + \text{sign}(s_{xy}) \sqrt{(s_{xx} - s_{yy})^2 + 4s_{xy}^2}}{2\sigma_x^{-1} \sigma_y s_{xy}} \quad (7)$$

with $s_{xx} := \sigma_y^2 \sum_{i=1}^n \mathbf{x}_i^T \mathbf{x}_i$, $s_{yy} := \sigma_x^2 \sum_{i=1}^n \mathbf{y}_i^T \mathbf{y}_i$ and $s_{xy} := \sigma_y \sigma_x \sum_{i=1}^n \mathbf{x}_i^T \mathbf{y}_i$. Together, these equations give a closed-form solution for the ML estimator of λ , assuming the measurement error variances σ_x^2 and σ_y^2 are known. By analyzing this result, it can be concluded that

1) λ^* always lies in between λ_x^* and λ_y^* , and

¹The noise in \mathbf{x}_i does not depend on λ as it is proportional to the average keypoint depth, which is normalized to 1 for the first keyframe.

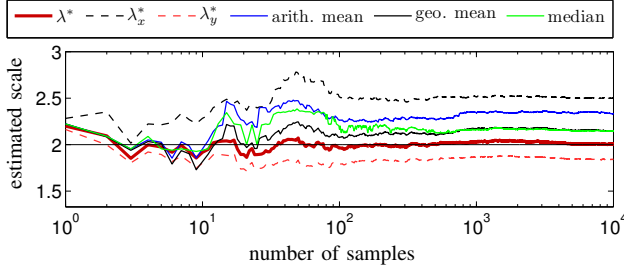


Fig. 4. Comparison of λ^* with Other Estimators: The plot shows the estimated scale as more samples are added. It can be seen that the proposed estimator λ^* is the only consistent estimator, i.e., the only one converging to the correct value. For this plot we used $\lambda = 2$, $\sigma_x = 1$, $\sigma_y = 0.3$ and $\mu_i \sim \mathcal{N}(\mathbf{0}_3, \mathbf{I}_{3 \times 3})$.

- 2) $\lambda^* \rightarrow \lambda_x^*$ for $\sigma_x^2 \rightarrow 0$, and $\lambda^* \rightarrow \lambda_y^*$ for $\sigma_y^2 \rightarrow 0$, i.e., these naïve estimators correspond to the case when one of the measurement sources is noise-free.

We extensively tested our approach on artificially generated data according to (2) and compared it to other, simple estimators, that is the arithmetic mean, geometric mean and the median of the set of quotients $\frac{\|\mathbf{x}_i\|}{\|\mathbf{y}_i\|}$. It can be observed that out of all presented possibilities, our approach is the only consistent estimator, i.e., the only one converging to the true scale for all dimensions d , values for σ_x^2 , σ_y^2 and values for μ_i . An example is shown in Fig. 4. Furthermore, λ^* can be computed efficiently, as each new sample pair only requires one update of the three sums, and the re-evaluation (7). Note that in practice approximations for σ_x^2 and σ_y^2 are sufficient, as their influence on λ^* decreases rapidly the more accurate the measured distances are. More results on the accuracy of this method will be presented in Section V-A.

B. State Prediction and Observation

In this section, we describe the state space, the observation models and the motion model used in the EKF. The state space consists of a total of ten state variables

$$\mathbf{x}_t := (x_t, y_t, z_t, \dot{x}_t, \dot{y}_t, \dot{z}_t, \Phi_t, \Theta_t, \Psi_t, \dot{\Psi}_t)^T \in \mathbb{R}^{10}, \quad (8)$$

where (x_t, y_t, z_t) denotes the position of the quadcopter in m and $(\dot{x}_t, \dot{y}_t, \dot{z}_t)$ the velocity in m/s, both in world coordinates. Further, the state contains the roll Φ_t , pitch Θ_t and yaw Ψ_t angle of the drone in deg, as well as the yaw-rotational speed $\dot{\Psi}_t$ in deg/s. In the following, we define for each sensor an observation function $h(\mathbf{x}_t)$ and describe how the respective observation vector \mathbf{z}_t is composed from the sensor readings.

1) **Odometry Observation Model:** The quadcopter measures its horizontal speed $\hat{v}_{x,t}$ and $\hat{v}_{y,t}$ in its local coordinate system, which we transform into the global frame \dot{x}_t and \dot{y}_t . The roll and pitch angles $\hat{\Phi}_t$ and $\hat{\Theta}_t$ measured by the accelerometer are direct observations of Φ_t and Θ_t . To account for yaw-drift and uneven ground, we differentiate the height measurements \hat{h}_t and yaw measurements $\hat{\Psi}_t$ and treat them as observations of the respective velocities. The resulting observation function $h_I(\mathbf{x}_t)$ and measurement vector

$\mathbf{z}_{I,t}$ is hence given by

$$h_I(\mathbf{x}_t) := \begin{pmatrix} \dot{x}_t \cos \Psi_t - \dot{y}_t \sin \Psi_t \\ \dot{x}_t \sin \Psi_t + \dot{y}_t \cos \Psi_t \\ \dot{z}_t \\ \Phi_t \\ \Theta_t \\ \dot{\Psi}_t \end{pmatrix} \quad (9)$$

$$\mathbf{z}_{I,t} := (\hat{v}_{x,t}, \hat{v}_{y,t}, (\hat{h}_t - \hat{h}_{t-1}), \hat{\Phi}_t, \hat{\Theta}_t, (\hat{\Psi}_t - \hat{\Psi}_{t-1}))^T \quad (10)$$

2) **Visual Observation Model:** When PTAM successfully tracks a video frame, we scale the pose estimate by the current estimate for the scaling factor λ^* and transform it from the coordinate system of the front camera to the coordinate system of the quadcopter, leading to a direct observation of the quadcopter's pose given by

$$h_P(\mathbf{x}_t) := (x_t, y_t, z_t, \Phi_t, \Theta_t, \Psi_t)^T \quad (11)$$

$$\mathbf{z}_{P,t} := f(\mathbf{E}_{DC} \mathbf{E}_{C,t}) \quad (12)$$

where $\mathbf{E}_{C,t} \in \text{SE}(3)$ is the estimated camera pose (scaled with λ), $\mathbf{E}_{DC} \in \text{SE}(3)$ the constant transformation from the camera to the quadcopter coordinate system, and $f: \text{SE}(3) \rightarrow \mathbb{R}^6$ the transformation from an element of $\text{SE}(3)$ to our roll-pitch-yaw representation.

3) **Prediction Model:** The prediction model describes how the state vector \mathbf{x}_t evolves from one time step to the next. In particular, we approximate the quadcopter's horizontal acceleration \ddot{x}, \ddot{y} based on its current state \mathbf{x}_t , and estimate its vertical acceleration \ddot{z} , yaw-rotational acceleration $\ddot{\Psi}$ and roll/pitch rotational speed $\dot{\Phi}, \dot{\Theta}$ based on the state \mathbf{x}_t and the active control command \mathbf{u}_t .

The horizontal acceleration is proportional to the horizontal force acting upon the quadcopter, which is given by

$$\begin{pmatrix} \ddot{x} \\ \ddot{y} \end{pmatrix} \propto \mathbf{f}_{\text{acc}} - \mathbf{f}_{\text{drag}} \quad (13)$$

where \mathbf{f}_{drag} denotes the drag and \mathbf{f}_{acc} denotes the accelerating force. The drag is approximately proportional to the horizontal velocity of the quadcopter, while \mathbf{f}_{acc} depends on the tilt angle. We approximate it by projecting the quadcopter's z -axis onto the horizontal plane, which leads to

$$\ddot{x}(\mathbf{x}_t) = c_1 (\cos \Psi_t \sin \Phi_t \cos \Theta_t - \sin \Psi_t \sin \Theta_t) - c_2 \dot{x}_t \quad (14)$$

$$\ddot{y}(\mathbf{x}_t) = c_1 (-\sin \Psi_t \sin \Phi_t \cos \Theta_t - \cos \Psi_t \sin \Theta_t) - c_2 \dot{y}_t \quad (15)$$

We estimated the proportionality coefficients c_1 and c_2 from data collected in a series of test flights. Note that this model assumes that the overall thrust generated by the four rotors is constant. Furthermore, we describe the influence of sent control commands $\mathbf{u}_t = (\bar{\Phi}_t, \bar{\Theta}_t, \bar{z}_t, \bar{\Psi}_t)$ by a linear model:

$$\dot{\Phi}(\mathbf{x}_t, \mathbf{u}_t) = c_3 \bar{\Phi}_t - c_4 \Phi_t \quad (16)$$

$$\dot{\Theta}(\mathbf{x}_t, \mathbf{u}_t) = c_3 \bar{\Theta}_t - c_4 \Theta_t \quad (17)$$

$$\ddot{\Psi}(\mathbf{x}_t, \mathbf{u}_t) = c_5 \bar{\Psi}_t - c_6 \dot{\Psi}_t \quad (18)$$

$$\ddot{z}(\mathbf{x}_t, \mathbf{u}_t) = c_7 \bar{z}_t - c_8 \dot{z}_t \quad (19)$$

Again, we estimated the coefficients c_3, \dots, c_8 from test flight data. The overall state transition function is now given by

$$\begin{pmatrix} x_{t+1} \\ y_{t+1} \\ z_{t+1} \\ \dot{x}_{t+1} \\ \dot{y}_{t+1} \\ \dot{z}_{t+1} \\ \Phi_{t+1} \\ \Theta_{t+1} \\ \Psi_{t+1} \\ \dot{\Psi}_{t+1} \end{pmatrix} \leftarrow \begin{pmatrix} x_t \\ y_t \\ z_t \\ \dot{x}_t \\ \dot{y}_t \\ \dot{z}_t \\ \Phi_t \\ \Theta_t \\ \Psi_t \\ \dot{\Psi}_t \end{pmatrix} + \delta_t \begin{pmatrix} \dot{x}_t \\ \dot{y}_t \\ \dot{z}_t \\ \ddot{x}(\mathbf{x}_t) \\ \ddot{y}(\mathbf{x}_t) \\ \ddot{z}(\mathbf{x}_t, \mathbf{u}_t) \\ \dot{\Phi}(\mathbf{x}_t, \mathbf{u}_t) \\ \dot{\Theta}(\mathbf{x}_t, \mathbf{u}_t) \\ \dot{\Psi}_t \\ \ddot{\Psi}(\mathbf{x}_t, \mathbf{u}_t) \end{pmatrix} \quad (20)$$

using the model specified in (14) to (19). Note that, due to the many assumptions made, we do not claim the physical correctness of this model. It however performs very well in practice, which is mainly due to its completeness: the behavior of all state parameters and the effect of all control commands is approximated, allowing “blind” prediction, i.e., prediction without observations for a brief period of time (~ 125 ms in practice, see Fig. 3).

V. EXPERIMENTS AND RESULTS

We conducted a series of real-world experiments to analyze the properties of the resulting system. The experiments were conducted in different environments, i.e., both indoor in rooms of varying size and visual appearance as well as outdoor under the influence of sunlight and wind. A selection of these environments is depicted in Fig. 5.

In the following, we present our results on the convergence behavior and accuracy of scale estimation in Section IV-A, the accuracy of the motion model in Section V-B, the responsiveness and accuracy of pose control in Section V-C, and the long-term stability and drift elimination in Section V-D.

As ground truth at time t we use the state of the EKF after all odometry and visual pose information up to t have been received and integrated. It can only be calculated at $t + \Delta t_{\text{vis}}$, and therefore is not used for drone control – in practice it is available ~ 250 ms after a control command for t has been computed and sent to the quadcopter.

A. Scale Estimation Accuracy

To analyze the accuracy of the scale estimation method derived in IV-A, we instructed the quadcopter to fly a fixed figure, while every second a new sample is taken and the scale re-estimated. In the first set of flights, the quadcopter was commanded to move only vertically, such that the samples mostly consist of altitude measurements. In the second set, the quadcopter was commanded to fly a horizontal rectangle, such that primarily the IMU-based velocity information is used. After each flight, we measured the ground truth $\hat{\lambda}$ by manually placing the quadcopter at two measurement points, and comparing the known distance between these points with the distance measured by the visual SLAM system. Note that due to the initial scale normalization, the values for $\hat{\lambda}$ roughly correspond to the

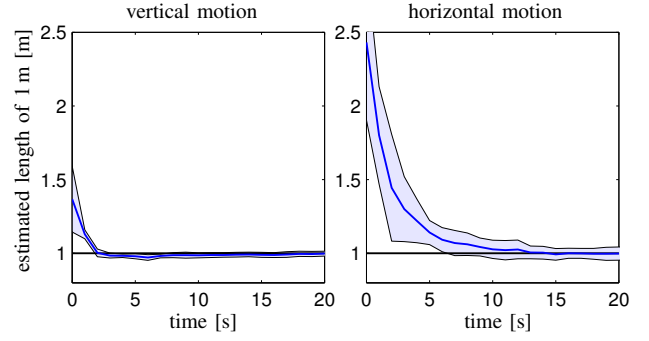


Fig. 6. Scale Estimation Accuracy: The plots show the mean and standard deviation of the the estimation error e , corresponding to the estimated length of 1 m, from horizontal and vertical motion. It can be seen that the scale can be estimated accurately in both cases, it is however more accurate and converges faster if the quadcopter moves vertically.

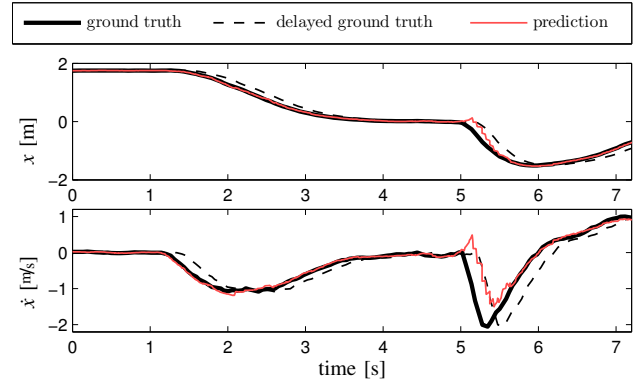


Fig. 7. Comparison of Predicted and Real State. The black curve shows the ground truth, it can only be computed with a delay of ~ 250 ms (dashed curve). At $t = 5$ s, the quadcopter is manually pushed away which cannot be predicted – hence the brief deviation. This plot shows (1) that the prediction approximates the ground truth well and in particular without notable delay and (2) that using visual information, the EKF rapidly recovers from large external disturbances – however with a small delay.

mean feature depth in meters of the first keyframe, which in our experiments ranges from 2 m to 10 m. To provide better comparability, we analyze and visualize the estimation error $e := \frac{\hat{\lambda}^*}{\hat{\lambda}}$, corresponding to the estimated length of 1 m.

Fig. 6 gives the mean error as well as the standard deviation spread over 10 flights. As can be seen, our method quickly and accurately estimates the scale from both types of motion. Due to the superior accuracy of the altimeter compared to the horizontal velocity estimates, the estimate converges faster and is more accurate if the quadcopter moves vertically, i.e., convergence after 2 s versus 15 s, and to a final accuracy $\pm 1.7\%$ versus $\pm 5\%$. Note that in practice, we allow for (and recommend) arbitrary motions during scale estimation so that information from both sensor modalities can be used to improve convergence. Large, sudden changes in measured relative height can be attributed to uneven ground, and removed automatically from the data set.

B. State Prediction Accuracy

In this section we give a qualitative evaluation of the accuracy of the predicted state of the quadcopter, used for control. Fig. 7 shows both the predicted state for time t as well as the ground truth, i.e., the state computed after

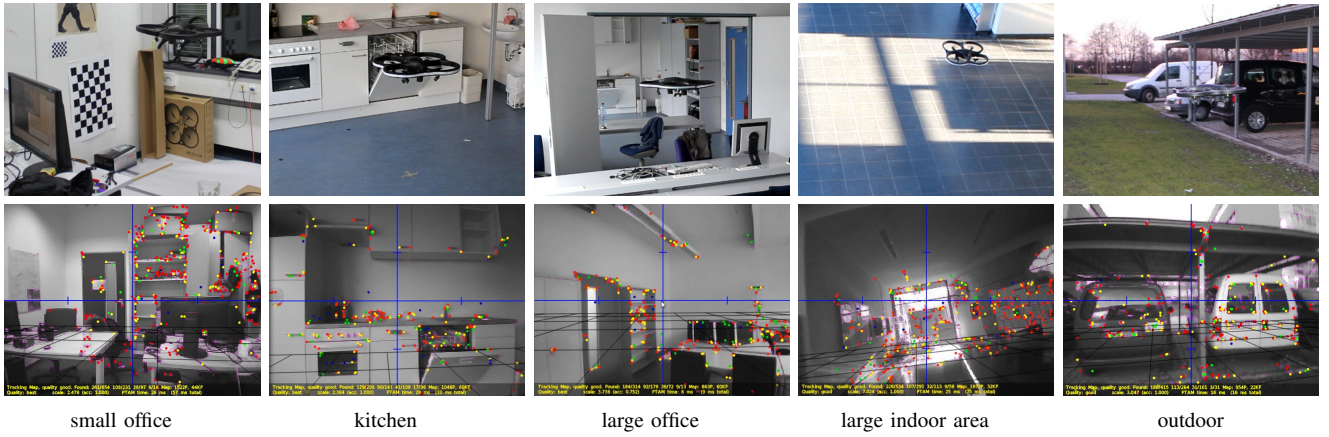


Fig. 5. Testing Environments: The top row shows an image of the quadcopter flying, the bottom row the corresponding image from the quadcopter's frontal camera. This shows that our system can operate robustly in different, real-world environments.

TABLE I
CONVERGENCE SPEED IN POSITION CONTROL

relative motion (x,y,z) [m]	(1,0,0)	(4,0,0)	(0,0,1)	(1,1,1)
t_{conv} [s]	3.1 ± 1.3	4.5 ± 0.8	3.1 ± 0.1	3.9 ± 0.5

all sensor measurements have been evaluated. This is only possible $\sim 250\text{ms}$ after the respective control command has been issued. It can be observed that the prediction approximates the ground truth very well and without notable delay, which is crucial for oscillation-free control.

C. Positioning Accuracy and Convergence Speed

In this Section, we evaluated the performance of the complete system in terms of position control. In particular, we measured the average time to approach a given goal location and the average positioning error while holding this position. Considering the large delay in our system, the pose stability of the quadcopter heavily depends on an accurate prediction from the EKF: the more accurate the pose estimates and in particular the velocity estimates are, the higher the gains can be set without leading to oscillations.

To determine the stability, we instructed the quadcopter to hold a target position over 60 s in different environments and measure the root mean square error (RMSE). The results are given in Fig. 10: the measured RMSE lies between 4.9 cm (indoor) and 18.0 cm (outdoor).

To evaluate the convergence speed, we repeatedly let the quadcopter fly a given distance and measure the convergence time t_{conv} , corresponding to the time required to reach the target position and hold it for at least 5 s. We consider the quadcopter to be at the target position if the Euclidean distance is less than 10 cm. An example of flying a long distance in x -direction is shown in Fig. 8: the plot clearly shows that the quadcopter accelerates initially with maximum pitch, and de-accelerates before reaching the target location at $t = 3.5\text{s}$. Fig. 9 shows an example trajectory in all three dimensions. We repeated this experiment ten times and summarized the results in Tab. I. Reaching a target location at a distance of 4 m took on average 4.5 s.

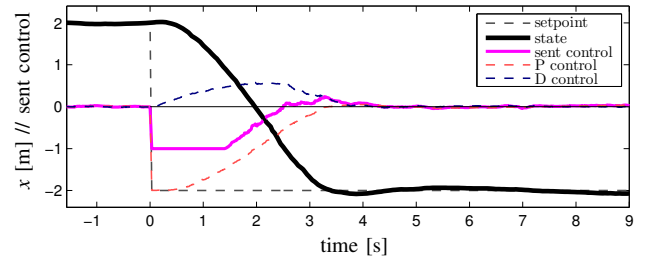


Fig. 8. Flying a Large Distance: The plot shows the behavior of the controller for a large distance. As can be seen, the quadcopter accelerates with maximum pitch for the first second and decelerates before converging on the setpoint.

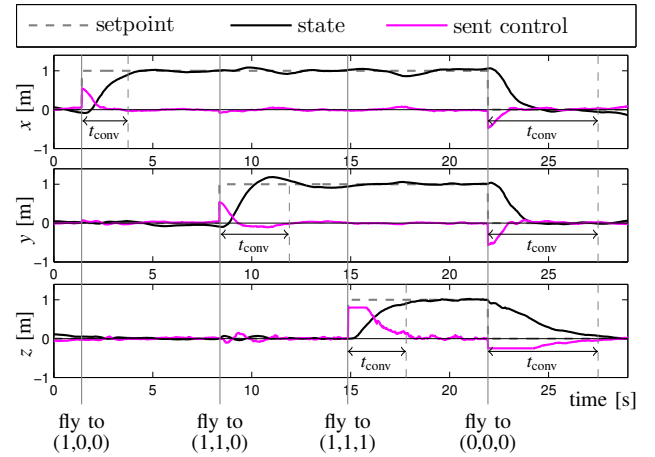


Fig. 9. Example Flight: Flying a simple figure consisting of four waypoints. This plot illustrates the typical behavior of the quadcopter when holding and approaching waypoints (t_{conv} is indicated, see also Tab. I).

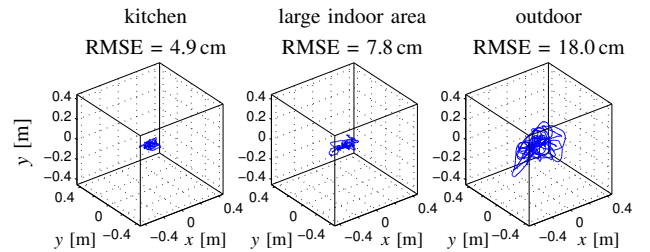


Fig. 10. Flight Stability: Path taken and RMSE of the quadcopter when instructed to hold a target position for 60 s, in three of the environments depicted in Fig. 5. It can be seen that the quadcopter can hold a position very accurately, even when perturbed by wind (right).

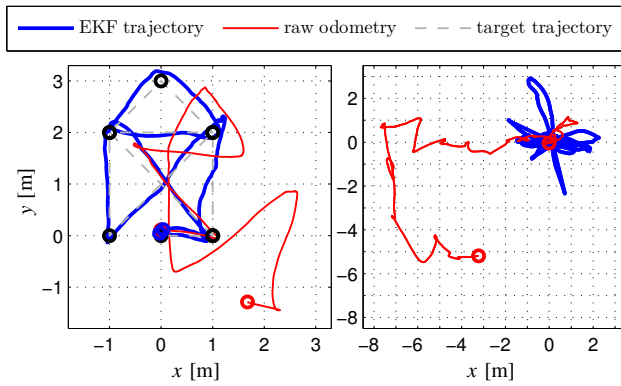


Fig. 11. Elimination of Odometry Drift: Horizontal path taken by the quadcopter as estimated by the EKF compared to the raw odometry (i.e., the integrated velocity estimates). Left: when flying a figure; right: when being pushed away repeatedly from its target position. The odometry drift is clearly visible, in particular when the quadcopter is being pushed away. When incorporating visual pose estimates, it is eliminated completely.

D. Drift Elimination

To verify that the incorporation of a visual SLAM system eliminates odometry drift, we compare the estimated trajectory with and without the visual SLAM system. Fig. 11 shows the resulting paths, both for flying a fixed figure (left) and for holding a target position while the quadcopter is being pushed away (right). Both flights took approximately 35 s, and the quadcopter landed no more than 15 cm away from its takeoff position. In contrast, the raw odometry accumulated an error of 2.1 m for the fixed figure and 6 m when being pushed away. This experiment demonstrates that the visual SLAM system efficiently eliminates pose drift during maneuvering.

E. Robustness to Temporary Loss of Visual Tracking

The system as a whole is robust to temporary loss of visual tracking, e.g. due to occlusions or large rotations, as it continues to navigate based only on odometry measurements. As soon as visual tracking recovers, the EKF state is updated with the absolute pose estimate, eliminating accumulated estimation error. This is demonstrated in the attached video.

VI. CONCLUSION

In this paper, we presented a visual navigation system for a low-cost quadcopter using offboard processing. Our system enables the quadcopter to visually navigate in unstructured, GPS-denied environments and does not require artificial landmarks nor prior knowledge about the environment. The contribution of this paper is two-fold: first, we presented a robust solution for visual navigation with a low-cost quadcopter. Second, we derived a maximum-likelihood estimator in closed-form to recover the absolute scale of the visual map, providing an efficient and consistent alternative to predominant filtering-based methods. Our system was able to estimate the map scale up to $\pm 1.7\%$ of its true value, with which we achieved an average positioning accuracy of 4.9 cm (indoor) to 18.0 cm (outdoor). Furthermore, our approach is able to robustly deal with communication delays of up to 400 ms. We tested our system in a set of extensive

experiments in different real-world indoor and outdoor environments. With these experiments, we demonstrated that accurate, robust and drift-free visual navigation is feasible even with low-cost robotic hardware.

REFERENCES

- [1] D. Mellinger and V. Kumar, "Minimum snap trajectory generation and control for quadrotors," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2011.
- [2] S. Lupashin, A. Schöllig, M. Sherback, and R. D'Andrea, "A simple learning strategy for high-speed quadcopter multi-flips," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2010.
- [3] M. Müller, S. Lupashin, and R. D'Andrea, "Quadcopter ball juggling," in *Proc. IEEE Intl. Conf. on Intelligent Robots and Systems (IROS)*, 2011.
- [4] Q. Lindsey, D. Mellinger, and V. Kumar, "Construction of cubic structures with quadrotor teams," in *Proceedings of Robotics: Science and Systems (RSS)*, Los Angeles, CA, USA, 2011.
- [5] S. Grzonka, G. Grisetti, and W. Burgard, "Towards a navigation system for autonomous indoor flying," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2009.
- [6] M. Blösch, S. Weiss, D. Scaramuzza, and R. Siegwart, "Vision based MAV navigation in unknown and unstructured environments," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2010.
- [7] M. Achtelik, M. Achtelik, S. Weiss, and R. Siegwart, "Onboard IMU and monocular vision based control for MAVs in unknown in- and outdoor environments," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2011.
- [8] M. Achtelik, A. Bachrach, R. He, S. Prentice, and N. Roy, "Stereo vision and laser odometry for autonomous helicopters in GPS-denied indoor environments," in *Proc. SPIE Unmanned Systems Technology XI*, 2009.
- [9] A. S. Huang, A. Bachrach, P. Henry, M. Krainin, D. Maturana, D. Fox, and N. Roy, "Visual odometry and mapping for autonomous flight using an RGB-D camera," in *Proc. IEEE International Symposium of Robotics Research (ISRR)*, 2011.
- [10] D. Mellinger, N. Michael, and V. Kumar, "Trajectory generation and control for precise aggressive maneuvers with quadrotors," in *Proceedings of the Intl. Symposium on Experimental Robotics*, Dec 2010.
- [11] D. Eberli, D. Scaramuzza, S. Weiss, and R. Siegwart, "Vision based position control for MAVs using one single circular landmark," *Journal of Intelligent and Robotic Systems*, vol. 61, pp. 495–512, 2011.
- [12] T. Krajník, V. Vonásek, D. Fišer, and J. Faigl, "AR-drone as a platform for robotic research and education," in *Proc. Research and Education in Robotics: EUROBOT 2011*, 2011.
- [13] "Ascending technologies," 2012. [Online]: <http://www.ascotec.de/>
- [14] M. Blösch, S. Weiss, D. Scaramuzza, and R. Siegwart, "Vision based MAV navigation in unknown and unstructured environments," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2010.
- [15] S. Zingg, D. Scaramuzza, S. Weiss, and R. Siegwart, "MAV navigation through indoor corridors using optical flow," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2010.
- [16] "Parrot AR.Drone," 2012. [Online]: <http://ardrone.parrot.com/>
- [17] S. Weiss, M. Achtelik, M. Chli, and R. Siegwart, "Versatile distributed pose estimation and sensor self-calibration for an autonomous mav," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2012.
- [18] G. Nützi, S. Weiss, D. Scaramuzza, and R. Siegwart, "Fusion of IMU and vision for absolute scale estimation in monocular SLAM," *Journal of Intelligent Robotic Systems*, vol. 61, pp. 287 – 299, 2010.
- [19] C. Bills, J. Chen, and A. Saxena, "Autonomous MAV flight in indoor environments using single image perspective cues," in *Proc. IEEE Intl. Conf. on Robotics and Automation (ICRA)*, 2011.
- [20] T. Krajník, V. Vonásek, D. Fišer, and J. Faigl, "AR-drone as a platform for robotic research and education," in *Proc. Communications in Computer and Information Science (CCIS)*, 2011.
- [21] W. S. Ng and E. Sharlin, "Collocated interaction with flying robots," in *Proc. IEEE Intl. Symposium on Robot and Human Interactive Communication*, 2011.
- [22] G. Klein and D. Murray, "Parallel tracking and mapping for small AR workspaces," in *Proc. IEEE Intl. Symposium on Mixed and Augmented Reality (ISMAR)*, 2007.