

Autonomous Landing of MAVs on an Arbitrarily Textured Landing Site using Onboard Monocular Vision

Shaowu Yang, Sebastian A. Scherer, Konstantin Schauwecker and Andreas Zell

Abstract—This paper presents a novel solution for micro aerial vehicles (MAVs) to autonomously search for and land on an arbitrary landing site using real-time monocular vision. The autonomous MAV is provided with only one single reference image of the landing site with an unknown size before initiating this task. We extend a well-known monocular visual SLAM algorithm that enables autonomous navigation of the MAV in unknown environments, in order to search for such landing sites. Furthermore, a multi-scale ORB feature based method is implemented and integrated into the SLAM framework for landing site detection. We use a RANSAC-based method to locate the landing site within the map of the SLAM system, taking advantage of those map points associated with the detected landing site. We demonstrate the efficiency of the presented vision system in autonomous flights, both indoor and in challenging outdoor environment.

I. INTRODUCTION

The growing research area of Micro Aerial Vehicles (MAVs) has attracted much attention in the robotics community in recent years. One interesting focus has been on using onboard sensors such as cameras and laser scanners, which do not rely on any external signal, to facilitate their autonomous navigation. These onboard sensors are important replacements for GPS sensors in environments where GPS is unavailable or not reliable, such as indoors or in outdoor urban areas.

2D laser scanners have been successfully used for autonomous navigation of MAVs indoor [12], [28]. However, they are difficult to be extended to object recognition tasks. Compared with other sensors, cameras are passive, and have a superior



Fig. 1: Our MAV navigating autonomously to search for a textured landing site.

potential for environment perception, while still being lightweight, relatively low cost and energy efficient. Moreover, unlike stereo cameras with small baselines, a monocular camera does not lose its functionality even for large working distances, when metric scale is properly tracked. Those advantages make monocular vision very attractive for research on autonomous navigation of MAVs, which in general have very limited payload, both in weight and in computational capability.

Autonomous landing is a basic but also challenging phase for autonomous navigation of MAVs. When the exact position of a desired landing site is unknown, the MAVs should be able to search for and locate it autonomously, and then land on it to finish autonomous flights. Monocular visual simultaneous localization and mapping (SLAM) has brought more flexibility to autonomous navigation of MAVs in unknown environments [1]. In fact, it is especially well-suited for the autonomous landing task of an MAV: The problem of slow scale drift, which is inherent to every purely visual monocular SLAM system caused by the unobservability of the scale factor,

S. Yang, S. A. Scherer and K. Schauwecker are PhD students with the Department of Computer Science, University of Tübingen, Tübingen, Germany {shaowu.yang, sebastian.scherer, konstantin.schauwecker}@uni-tuebingen.de

A. Zell is full professor with the Department of Computer Science, University of Tübingen, Tübingen, Germany andreas.zell@uni-tuebingen.de

can hardly cause much effect in such relatively small areas where the MAV is expected to land.

In this paper we show that the rich information provided by a visual SLAM system can also benefit both the real time detection of a known landing site and its localization. Considering the limited computational power that is typically available onboard MAVs, those processes are normally difficult to be performed in parallel to autonomous navigation using onboard visual systems. We achieve autonomous navigation of our MAV by implementing a constant-time monocular visual SLAM framework, while simultaneously detecting an arbitrarily textured landing site using ORB features [21], and estimating its global pose. The resulted monocular vision system enables the MAV to autonomously search for the landing site in unknown environments (as depicted in Fig. 1), and then land on it once it is found.

II. RELATED WORK

Autonomous navigation of Unmanned Aerial Vehicles (UAVs) relying on pose estimation from GPS sensors has been well studied in early research. Those works are usually aided by fusing inertial navigation system (INS) data. UAVs with such navigation systems work well for high altitude and long range tasks, but are not suitable in GPS-denied environments. In recent years, more effort has been focused on using computer vision to enable autonomous flight of UAVs. Computer vision methods do not depend on external signals. Moreover, they fit especially well to cases in which precise position control relative to other objects is required, e.g. for the landing tasks of UAVs. Thus, they are highly appreciated for research towards full autonomy of UAVs.

In [24], the landing task of a helicopter is solved by using image moments for object recognition, while the estimation of the relative position with respect to the landing pad still relies on precise height information provided by differential GPS. Garcia-Pardo et al. [11] present a strategy to find a safe landing area by searching the image for a circular area in which all the pixels have a level of contrast below a given threshold. The vision system developed in [8] allows a remote user to define target areas as waypoints or a landing

area for a UAV from a high resolution aerial or satellite image. In this work, a Scale Invariant Feature Transform (SIFT) based image-matching algorithm is implemented to find the natural landmarks, and an optical-flow-based method is used for the detection of a safe landing area.

Recently, more vision solutions for autonomous navigation and landing have been presented, due to the fast growing interest in MAVs, and especially quadrotors. Mahony et al. [17] provide a tutorial introduction to modelling, pose estimation and control of such multi-rotor MAVs. Meier et al. [19] present a new self-developed quadrotor system capable of autonomous flight with onboard pose estimation from vision and an Inertial Measurement Unit (IMU), while relying on artificial visual markers. Previous work in [30] features an onboard monocular vision solution for autonomous takeoff, hovering and landing of an MAV based on a circular landing pad. Those works, achieving autonomous flight of MAVs, still depend on pose estimates from artificial landmarks, and are thus not flexible enough for long-term autonomy.

One way to be independent of artificial landmarks is to implement visual odometry or visual SLAM systems on MAVs. Fraundorfer et al. [10] extended the system in [19] with autonomous mapping and exploration, based on stereo cameras. In [2], [29], Parallel Tracking and Mapping (PTAM) [13] is implemented as a monocular visual SLAM framework for autonomous navigation of MAVs in unknown and GPS-denied environments. Achtelik [1] also use PTAM to provide position estimates for an MAV, while fusing data from an air pressure sensor and accelerometers to estimate the unknown metric scale factor of the monocular vision system. In [27], a modified PTAM, which integrates depth information as presented in [26], is used for position control of an MAV based on stereo vision.

Concerning SLAM and object recognition, another related work done by Castle et al. [4], [5] can be found in the field of augmented reality (AR). In [4], monoSLAM [9] and SIFT feature [16] are used to recognize and localize objects within a 3D map built by a wearable camera. Those objects were located from a single view using their known sizes and the location fed back to the

EKF. In [5], a multiple map and multiple camera extension to the PTAM algorithm is used to replace monoSLAM. Here, the location of an object is determined by triangulation from the locations of SIFT features matched across different keyframes, and it no longer needs to define the size of objects [5].

In our work, we also implement our visual SLAM framework based on PTAM, to enable autonomous navigation of an MAV, because of its robustness and its ability to generate an accurate map with a large number of map points from the environment. To land an MAV on an arbitrary landing site, we implement an ORB-feature-based method for landing site detection, running in parallel with the visual SLAM. Furthermore, based on the existing map points, we show that it is possible to robustly estimate the 3D pose of the detected landing site even if the size of it is unknown, and without re-triangulation from the landing site features as did in [5]. It is also different from those methods that only consider the relative pose estimation of an MAV with respect to a landing site, based on observations from the landing site itself. An example of such methods is the work in [18], which estimates the 3D pose of a camera for the control of UAVs by tracking a planar object with a known size. Since our pose estimation for MAV position control is provided by a SLAM system, high frequency landing site tracking and pose estimation become unnecessary, while still maintaining the final landing performance.

This paper is an improved and extended version of the work previously presented in a conference [31]. We extend the presentation and demonstrate the robustness of our method when working in outdoor environments by outdoor experiments in a challenging scenario. Moreover, besides using the setpoint method for trajectory control in autonomous navigations, we implemented a more efficient trajectory-following method, which benefits the navigation of MAVs in long trajectories. We further compare the performance of these two typical trajectory control methods in our autonomous navigation scheme.

III. VISUAL SLAM FOR AUTONOMOUS NAVIGATION

The visual SLAM framework we use for autonomous navigation of our MAV is based on

PTAM. In order to overcome the lack of a scale factor, we implemented an automatic initialization method for PTAM, which can cope with cluttered environments and provide a high accuracy. Additionally, we modify the mapping thread of PTAM to achieve a nearly constant processing time during navigation.

A. Basic Functionality of PTAM

The original PTAM implementation can produce detailed environmental maps with a large number of landmarks, which can be used for accurately tracking the pose of a monocular camera at a high frequency. In order to achieve real-time operation, a main idea proposed in PTAM is to split tracking and mapping into two separate threads, which can be processed in parallel on a dual-core computer. One thread is responsible for tracking the camera motion relative to the current map. The other thread extends this map, which consists of 3D point features that are organized in keyframes, and refines it using bundle adjustment.

In the thread responsible for tracking the camera pose, the FAST corner detector [23] is applied to each image at four pyramid levels, and all map points are projected to the current image coordinate frame, based on a prior pose estimate. The map points located inside the image after this projection are then used for tracking: To locate those points in the current camera image, a fixed-range image search around their predicted positions is performed. During this search, only the FAST corner locations are evaluated for finding the best matches. In our work, those FAST corners will further be used for feature-based object detection, without increasing the computation time in this thread.

The mapping thread integrates new keyframes into the map when requested by the tracking thread, and creates new map points by triangulating FAST corner matches between the new keyframe and its closest neighbours. Local bundle adjustment and global bundle adjustment are continuously performed to refine the map for the rest of the time. Since the map points are actually landmarks of the real-world scene, we will take advantage of their known 3D position for our landing site pose estimation.



Fig. 2: A scene when PTAM is initialized. *Top left*, original image. *Top right*, detected circular pattern, labelled with a orange cross. *Bottom left*, vision features in different levels. *Bottom right*, chosen map points.

B. Automatic Initialization of PTAM

Since there exists a common scale ambiguity inherent to monocular camera systems, PTAM naturally requires additional metric scale information. Since PTAM was originally intended for augmented reality applications [13], an accurate metric scale was not necessary, thus only a coarse scale estimate is applied to the triangulation of the initialization phase. We deal with this initialization issue by implementing the monocular solution presented in [30], which can robustly estimate the camera pose based on the image projection of a helicopter landing pad pattern, which also works in cluttered environments. Using this method, we can achieve accurate automatic initialization of PTAM during the takeoff phase of our MAV, without requiring any additional sensors. Fig. 2 shows an example scene and related results of PTAM, when initialized with this method.

1) Pose Estimation from a Circular Pattern:

In [30], we estimated the 6DOF camera pose based on the perspective projection of a typical helicopter landing pad, which consists of a letter “H” surrounded by a circle with a known diameter.

This pad is detected with a method similar to the one presented in [25]. Using adaptive thresholding, we obtain a binarized image that is used to find connected components with a run-based two-scan

labelling algorithm. The components are then classified using an artificial neural network. Finally a geometric constraint is applied, enforcing that the letter “H” must be surrounded by a circle. This allows us to detect the pad robustly in real-time with a high frequency.

After applying a Canny edge detector to the image pattern associated with the above pad, we can retrieve the ellipse that corresponds to the projection of the circle in the pad. At this point, we can obtain a 5DOF pose of the camera coordinate frame C with respect to the world coordinate frame \mathcal{W} , which is defined by the pad and obtained by using a computational geometry method based on the known quadratic equation of the projected ellipse. During this step, we also integrate IMU data to eliminate the remaining geometric ambiguity. Finally, fitting an ellipse to the projected contour of the letter “H” provides us with the last DOF of the camera pose, i.e. its yaw angle.

2) *Initializing PTAM during Takeoff*: Once we obtain an estimate of the camera pose with a height larger than a threshold h_i , then this pose estimate and the image associated with it are sent to PTAM for initialization. If more than a minimum number of FAST features with non-maximum suppression are detected on all four pyramid levels of this image, then we use them to initialize the map of PTAM. We obtain the 3D position of those feature points by assuming that they all lie on the ground plane and by projecting them from their image coordinates to the $z = 0$ plane in the world coordinate frame \mathcal{W} . In this way, the world coordinate frame defined in PTAM coincides with \mathcal{W} .

C. Using PTAM with Constant Computation Time

Bundle adjustment, which is used for map refinement, is the most computationally intensive task in PTAM. To enable PTAM to achieve a nearly constant computation time, we only retain its local bundle adjustment and abandon the global bundle adjustment, since it is rather computationally intensive and may stop the mapping thread from adding enough keyframes to facilitate successful tracking. However, we still keep the complete map during exploration.

IV. LANDING SITE DETECTION AND POSE ESTIMATION

To search for an arbitrary landing site during autonomous navigation of our MAV, we implemented a feature-based object detection scheme. Using one pre-set reference image of the designated landing site, a set of feature matches between the reference image and the currently visible scene can be established. Then the landing site is detected by using a robust RANSAC-based method to estimate the corresponding homography. Because some of the map points produced by PTAM can be associated with the matched features, we can use the 3D position estimates of those map points to estimate the global 3D pose of the landing site, even though there exists no absolute scale information for the landing site. The above process is integrated in the mapping thread of PTAM, as shown in Fig. 3.

A. Brief Overview of ORB

Rublee et al. [21] proposed the ORB (Oriented FAST and Rotated BRIEF) feature based on the FAST keypoint detector and the BRIEF descriptor [7], both of which are known for their high computational efficiency. BRIEF uses a binary string constructed from a set of binary intensity tests as an efficient point feature descriptor. Because BRIEF was not designed to be aware of the orientation of a feature point, it is notably lacking rotation invariance [21], which is, however, important for feature-matching-based object detection.

To cope with this issue, Rublee et al. proposed to compute an orientation component for

each FAST interest point (oFAST) by using the so-called intensity centroid, which is computed from image moments. BRIEF descriptors for those points are then efficiently rotated according to the orientation component, and thus form the steered BRIEF descriptor. Furthermore, a learning method is developed for choosing a good subset of binary tests, in order to increase the feature variance and reduce correlation among the binary tests, both of which are important for a discriminative feature. The resulting descriptor is named rBRIEF.

B. Applying Multi-Scale ORB to the PTAM Framework

We chose ORB as the feature descriptor for our landing site detection because of its low time cost and high discrimination capability for feature matching. ORB achieves scale invariance by applying the FAST detector to a scale-space pyramid of the original image. Since in the tracking thread of PTAM, FAST points have been detected in four-level pyramid images of the current scene, it is straightforward for us to use those points for further feature description. We chose such a multi-scale method in order to avoid the computation of further pyramid levels, as a compromise between matching performance and time cost. In the mapping thread, we compute orientation components of the FAST points to obtain oFAST features, and use rBRIEF for feature description. We perform both of these operations individually at pyramid level 0 and 1, resulting in two sets of descriptors $\{D_i^c | i = 0, 1\}$, each with a size n_i . We discard higher pyramid levels, since at higher levels, a landing site appears too small for us to obtain useful features for matching. For the reference image of the landing site, the number of pyramid levels and the scale factor for producing the pyramid images can vary according to the requirements of scale invariance and available computation time. In this paper we apply a three level pyramid with a scale factor of 1.2 to the reference image, obtaining the reference descriptor sets $\{D_i^r | i = 0, 1, 2\}$. A Gaussian blur is applied to each pyramid level before feature detection.

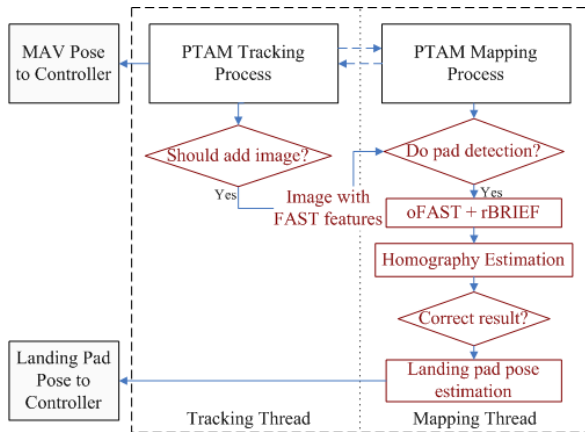


Fig. 3: Landing site detection and pose estimation integrated in the PTAM framework.

C. Landing Site Detection by Feature Matching

1) *Feature Matching*: We use a standard feature matching scheme to obtain a set of good feature

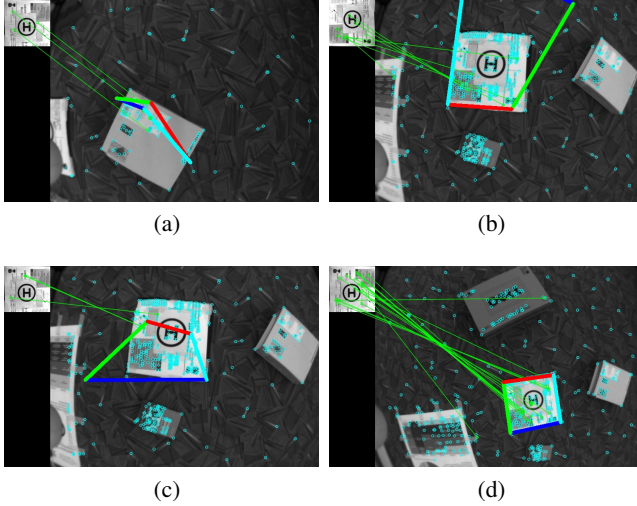


Fig. 4: Examples of homography estimation results shown in one pyramid level. After eliminating false estimates, only the one in (d) will be regarded as a correct homography estimate.

matches from $\{D_i^r | i = 0, 1, 2\}$ to $\{D_i^c | i = 0, 1\}$, for estimating the homography \mathbf{H}_{rc} between the reference image of the landing site and the current image frame. For finding all possible matches, we employ a brute-force matcher without cross checking, implemented in OpenCV [6]. It finds the k descriptors with the closest normalized Hamming distances in $\{D_i^c | i = 0, 1\}$ for each descriptor in $\{D_i^r | i = 0, 1, 2\}$. Similar to [14], [16], we consider a match between a reference descriptor and the corresponding descriptor with the closest distance to be valid, if the ratio of the closest to the second closest distance is smaller than a threshold T_r .

2) *Homography Estimation*: As the ORB feature is applied at different individual pyramid levels of the current camera image, we project all matched feature points to the source pyramid level to calculate the homography. The homography \mathbf{H}_{rc} is estimated by using RANSAC, and then further refined by using the Levenberg-Marquardt method to minimize the image projection error. We can limit the iterations in RANSAC to a relatively small number, in order to make this process more efficient. Since our consecutive landing site pose estimation can cope well with a lower true positive detection rate, we opt for a higher processing performance.

3) *Eliminating False Estimates*: The reference image forms a quadrilateral Q_r when transformed with \mathbf{H}_{rc} to the current image frame. Some examples of the homography estimates we received can be seen in Fig. 4. False Homography estimates may occur due to false matches or too few correct matches of features. We dramatically eliminate those false ones by evaluating some basic properties of this quadrilateral: First, it is required to be a convex polygon. Second, all four vertexes of it should have a reasonable relative distances to their centroid and to each other. This will eliminate estimates like the ones shown in Fig. 4b and Fig. 4c: Although the reference image can be found in the current image frame, we reject this frame since we will not achieve a correct pose estimate of the landing site according to this homography estimate. Third, the number of matched features that are inside of this quadrilateral should be larger than a threshold n_q . We determine whether a point is located inside a polygon using a crossing-number-based method.

D. Locating the Landing Site within the Map

After the landing site has been detected in the current camera image by using the above method, we locate its 3D pose in the world coordinate frame. For this task we take advantage of the environment map produced by PTAM, which can consist of a large number of map points. Doing this provides us with much more tolerance to false negative detections: Even if the landing site is not tracked at camera frame rate, its final pose estimate will be hardly affected, as the landing site should retain a static position in respect to the environment map. Thus, our method is very flexible in respect to the time intervals at which the mapping thread decides to add a new frame for landing site detection. Furthermore, using the map points ensures that only discriminative features are used for locating the landing site.

We first project all map points to a rectified image frame based on their 3D positions and the calibrated camera model [3]. Again, we use a crossing number method to check whether a projected map points is located within the quadrilateral Q_r (see Sect. IV-C.3). Those points inside form the map points subset $\{p_l\}$.

If the size of $\{p_l\}$ is larger than a threshold n_{lmin} , a RANSAC-based method is applied to the points in $\{p_l\}$ to estimate the dominant plane P_d of the landing site. We perform this step in a similar fashion as in [13]: Many sets of three points are randomly selected to form a plane hypothesis, while the remaining points are tested for consensus. The winning hypothesis is further refined by using the consensus set, resulting in the detected plane normal \mathbf{n}_p . Together with the mean 3D coordinate value of all consensus set points \mathbf{x}_m , this normal defines the plane P_d . Once an estimate for P_d is achieved, we use its corresponding measurements \mathbf{n}_p and \mathbf{x}_m as the initial guess for the RANSAC procedure when evaluating the next image frame. Thus, a much smaller threshold for the number of RANSAC iterations can be applied, which further reduce time costs.

The pose of the landing site can be calculated by projecting the quadrilateral Q_r to the plane P_d . We define $\mathbf{x}_i^p, i = 0, 1, 2, 3$, as the four vertices of Q_r , which are the image projections of the four corners P_i of the landing site, with the corresponding world coordinate positions \mathbf{x}_i^w . After projecting \mathbf{x}_i^p to a normalized image frame with rectified lens distortions, we obtain the normalized coordinates $\mathbf{x}_i^n = (x_i^n, y_i^n, 1)^T$. In the camera coordinate frame, we then have $\mathbf{x}_i^c = s \cdot (x_i^n, y_i^n, 1)^T$, with s being an undetermined scale factor. Thus, in the world coordinate frame we have

$$\mathbf{x}_i^w = s \cdot \mathbf{R}_{wc} \cdot \mathbf{x}_i^n + \mathbf{t}_{wc}, \quad (1)$$

with $\{\mathbf{R}_{wc}, \mathbf{t}_{wc}\}$ being the camera pose in the world coordinate frame, obtained by the tracking thread. Since P_i is located on the plane P_d , we have

$$\mathbf{n}_p \cdot (\mathbf{x}_i^w - \mathbf{x}_m) = 0. \quad (2)$$

From (1) and (2), we can calculate \mathbf{x}_i^w . The landing site pose is then obtained as $\mathbf{x}_k = \frac{1}{4} \sum_{i=0}^3 \mathbf{x}_i^w$, where k is the current image frame index.

We further refine the landing site pose by integrating m successful estimates of \mathbf{x}_k . Estimates with a large difference to $\mathbf{x}_{Lm} = \frac{1}{m} \sum_{k=0}^{m-1} \mathbf{x}_k$ are assumed to be outliers. The mean value of the remaining inlier is then assumed to be the final landing site pose estimate \mathbf{x}_L .

V. EXPERIMENTS AND RESULTS

A. Experimental Setup

1) *Quadrotor Platform*: Our MAV is based on the open source and open hardware quadrotor platform developed by the PIXHAWK project from ETH Zürich described in [19], which is depicted in Fig. 1. The onboard computer is a Kontron microETXexpress computer-on-module (COM) featuring an Intel Core 2 Duo 1.86GHz CPU, 2 GB DDR3 RAM and a 32Gb SSD. The pxIMU inertial measurement unit and autopilot board that we use mainly consists of a MicroController Unit (MCU), and sensors including a tri-axis accelerometer and a tri-axis gyroscope. The MCU is a 60 MHz ARM7 microcontroller for sensor readout and fusion, as well as position and attitude control. A PointGrey Firefly MV monochrome camera of only 37 g weight is mounted on the MAV in a downward-facing pose. This camera has an image resolution of 640×480, a maximum frame rate of 60 fps, and is equipped with a lens featuring a 90 degrees viewing angle.

2) *External Tracking System*: To measure ground truth data of the 6 DOF quadrotor pose and landing site poses, we use an external Optitrack tracking system manufactured by Naturalpoint¹, which comprises 9 infrared cameras in our case. After attaching several highly reflective markers to the quadrotor, the tracking system can provide 6 DOF pose estimates of the quadrotor with a frequency of up to 100 Hz. According to our tests, the deviation of the position estimates for a static quadrotor is in the order of only few millimeters.

3) *Coordinate Systems*: The coordinate systems are defined as right-hand systems and calibrated in the same way as we did in [30]. The world coordinates are indicated as the RGB axes (corresponding to $x-y-z$ axes) lie on the ground grids in Fig. 7. Its origin is attached to the center of the circular pattern from which our SLAM system is initialized as described in Sect. III-B.

4) *Software*: We implemented our software system in several modules using the open source Robot Operating System (ROS) [22] on Ubuntu Linux 12.04, as it provides the infrastructure for

¹<http://www.naturalpoint.com/optitrack/products/tracking-tools-bundles>

efficient communication among different modules and for logging all interested onboard data.

B. Navigation and Flight Control Algorithm

1) *Nested PID Pose Control:* Mellinger et al. [20] describe a nested PID controller that consists of a separate attitude and a position controller. Using a dynamic model of a quadrotor and an accurate 6 DoF pose estimate from an external tracking system, it can achieve precise hovering control of a quadrotor MAV. To evaluate our vision system, we control the pose of our quadrotor using a very similar controller, which is implemented in the original pxIMU code from the PIXHAWK project. In our case, we set the desired yaw angle to be a constant value of $\psi^{des} = 0$. The 3D position estimates from the onboard vision system are used as feedback to the position controller, and a basic Kalman Filter is applied to smooth pose estimation for low level control. The attitude controller runs at a frequency of 200 Hz, using the roll and pitch estimates by the IMU, and only the yaw angle is provided by the onboard vision system.

2) *Setpoint Method for Trajectory Control:* In order to search for the landing site, we implemented a setpoint-based method to autonomously navigate our MAV. Thus the MAV can follow a predefined searching path. We assume that the MAV has reached a setpoint, if its distance to this point is smaller than a threshold d_s and the yaw angle difference is smaller than ψ_s , for a period of time t_s . In this case, we advance to the next set point on the searching path. Once an initial pose of the landing site $\mathbf{x}_{l_{ini}}$ is estimated, we change the setpoint to be above this area, keeping our searching height h_s . After the final refined pose of the landing site $\mathbf{x}_l = (x_l, y_l, z_l)^T$ is determined, we define the end of the searching path to be $(x_l, y_l, h_s)^T$. Finally, the desired height of the setpoint is decreased until the MAV reaches a predefined landing height h_l where it can steadily shut down its motors to finish the landing process.

3) *More Efficient Trajectory-following Method:* We implemented a trajectory-following method to enable the MAV to follow the predefined searching path more efficiently. It is based on the path following controller presented in [20]. We simplify it by dropping the feed-forward term of the desired acceleration, and set a constant forward speed, v_s ,

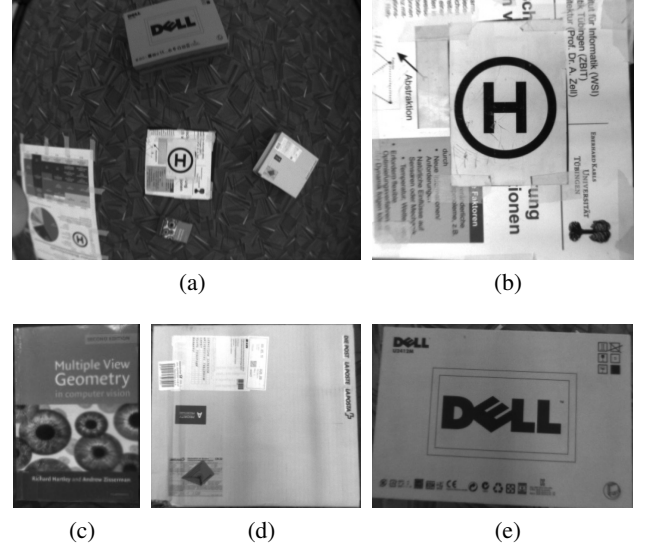


Fig. 5: (a) A scene from the MAV, (b), (c), (d) and (e) are the reference images of the poster landing pad (size 500×500 , height 4.5), the book (size 246×175 , height 33), the mail package (size 380×335 , height 140) and the computer package (size 650×435 , height 235), respectively. All size and height are measured in *mm*.

along the tangent direction of the searching path. The reason for this simplification is that we do not expect our MAV to fly in an aggressive way, and a constant expected speed is sufficient for our searching task.

The general idea of this method is to only consider position errors in the normal direction of the predefined path, while ignoring the errors in the tangent direction, which makes a constant expected forward speed possible. Then the final commanded acceleration is calculated by a PD controller. For more details we refer to [20].

C. Landing Site Pose Estimation Results

We evaluate the landing site detection and pose estimation results by processing a video logfile from a manual flight of our MAV above different objects with planar surfaces, which are used to represent different landing sites: a poster pad, a book, a mail package, and a computer package. Each of them has different texture features. Moreover, they are different in size and height. We control the MAV to take off from another pad nearby those

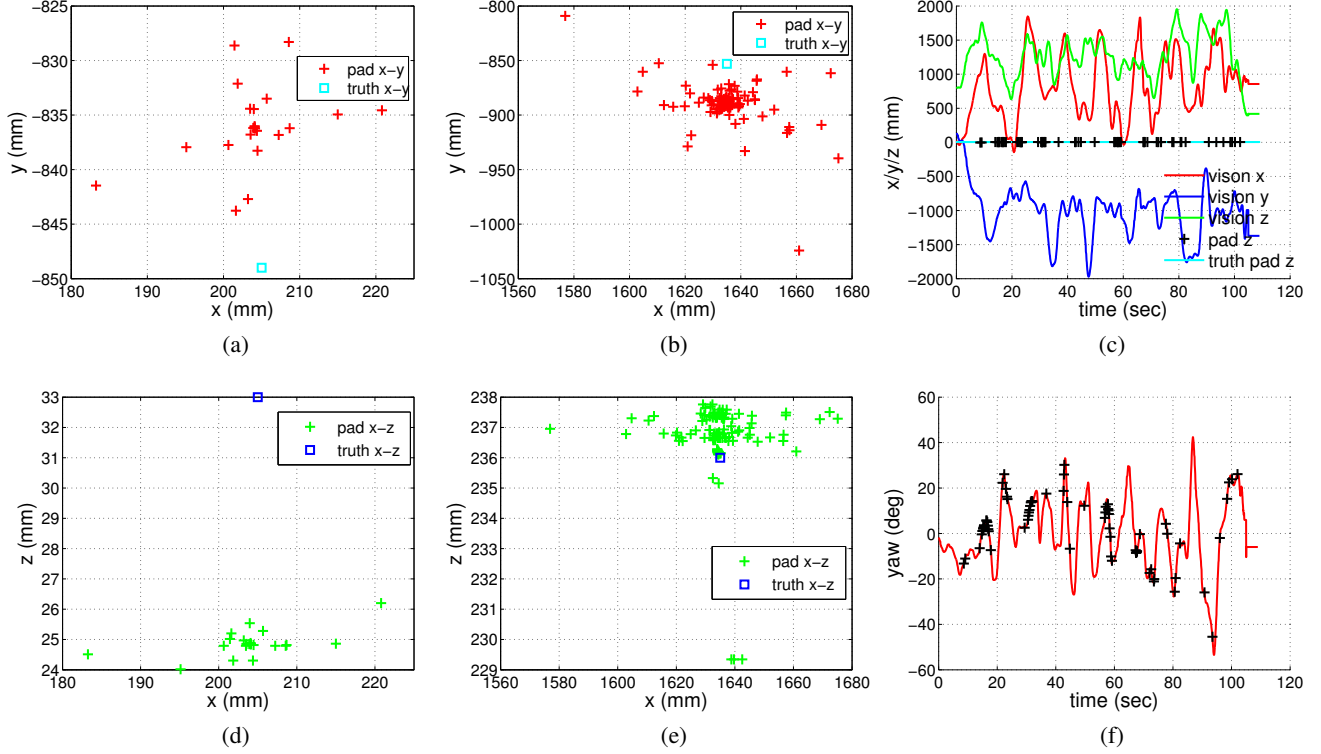


Fig. 6: (a), (d) Position estimates for the book, on $x-y$ and $x-z$ plane respectively, and (b), (e) for the PC package. (c) Trajectory of the MAV, and (f) the corresponding yaw angle estimates (a cross is marked if the landing site is detected at the corresponding time).

objects, such that our SLAM algorithm can be initialized by this pad as described in Sect. III-B. Fig. 5a shows a scene during this flight. Reference images of those landing sites are captured by manually holding the MAV above them in different illumination conditions, as shown in Fig. 5. Note that to confirm that our object detection method is invariant to the rotation of a reference image, they are rotated by 180 degrees for this experiment.

We process the same recorded video sequence four times, selecting a different reference image for landing site detection each time. The identical MAV trajectory estimated by the visual SLAM algorithm is shown in Fig. 6c and 6f. The poster pad provides a total number of 481 ORB features on all three levels, the book 69, the mail package 153 and the computer package 97 features. Despite their differences we mentioned above, they can be correctly detected and located. The Root Mean Square Errors (RMSEs) of their 3D position estimation are listed in Tab. I.

Since pose distributions of the detected landing sites are similar, we only present the results for the book and the computer package, which have the overall smallest and largest RMSEs. Fig. 6a and 6d show the distribution of the pose estimation results for the book, and 6b and 6e show that of the computer package. The pose estimates are projected to the $x-y$ and $x-z$ planes of the world coordinates. The few estimates with relatively large errors do not affect the autonomous navigation since they can be excluded after a pose refine process as described in Sect. IV-D. In Fig. 6c, we mark the height of the detected poster pad with black crosses, if it is detected at the corresponding time. Similarly, in Fig. 6f, we mark the MAV yaw angle estimates when the pad is detected. They show that the poster pad is detected when the MAV is at different positions and yaw angles.

D. Autonomous Navigation and Landing Results

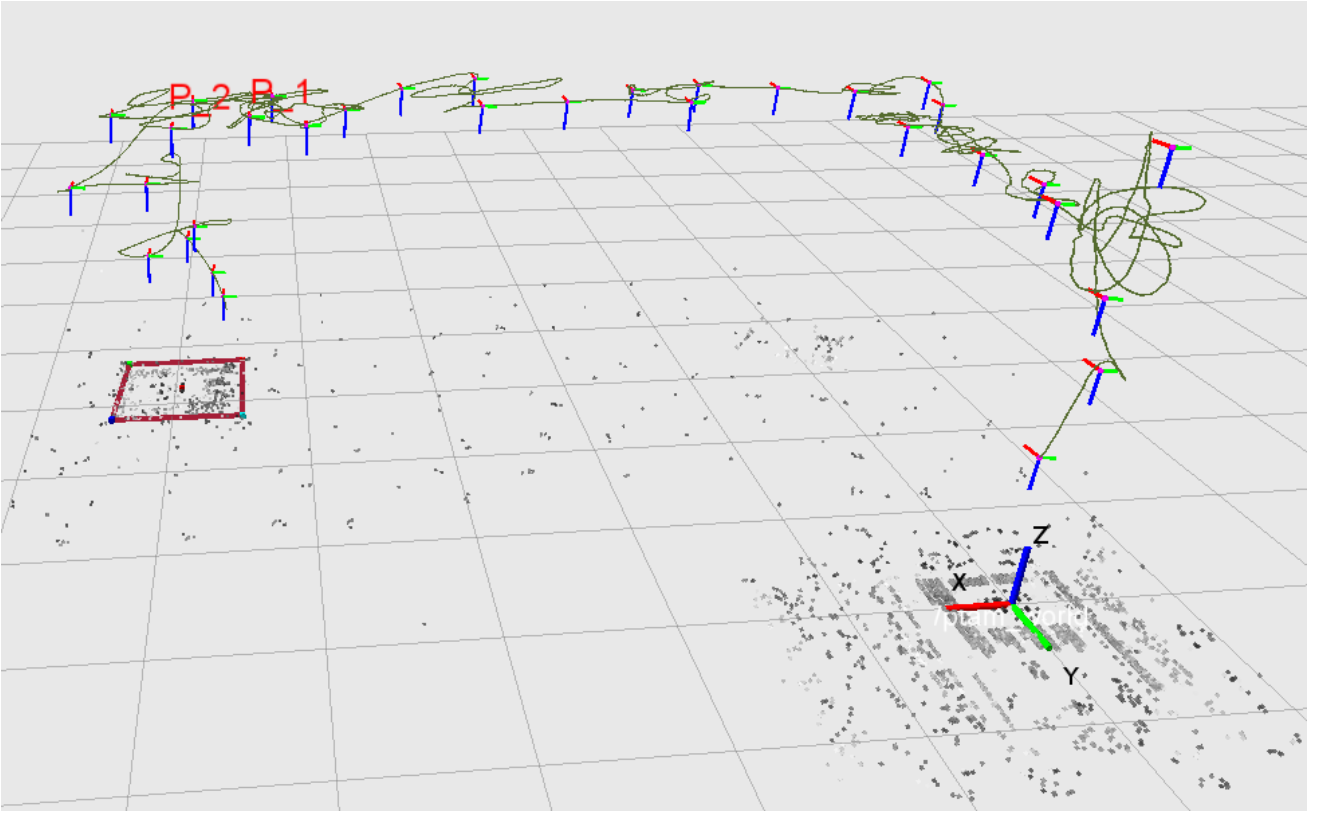


Fig. 7: The built map and MAV trajectory during a searching and landing task.

TABLE I: RMSEs (mm) of position estimates for different landing sites during a manual flight.

RMSE	poster	book	mail Pack.	PC Pack.
x-y	22	15	34	43
z	1	8	6	2
3D	22	17	35	43

1) *Using the setpoint method:* In this experiment, we use the poster pad shown in Fig. 5b as the target landing site. Our MAV autonomously navigates using the setpoint method to search for the landing site and finally lands on it. The trajectory of this searching and landing task, as estimated by our onboard SLAM system, is shown in Fig. 7. The map points built by the SLAM system are triangulated and refined if new keyframes are added. The pose of each keyframe has been depicted as small RGB axes in Fig. 7.

The searching path should depend on the expected complexity of the landing area. Here, a

simple rectangular searching path is defined. We choose four setpoints evenly distributed on each edge of the rectangle. The MAV navigates along this searching path after takeoff and initialization of the visual SLAM system. At each corner of the rectangle, the MAV is commanded to hover for 3 seconds to maintain better stability. The parameters described in Sect. V-B.2 are chosen as shown in Table II. The landing site is detected when the MAV is at the position $\mathbf{P}_1 = (2.001, -1.556, 1.197)^T(m)$, relative to the starting position. When it is at $\mathbf{P}_2 = (2.337, -1.616, 1.201)^T(m)$, landing site detection stops with the computation of the refined landing site pose, which is visualized as a bold colored quadrilateral in Fig. 7. Fig. 8 is the resulting trajectory on different axes of the world coordinates. It shows that the MAV needs to hover for a certain period of time to satisfy the constraint of each setpoint, and then moves to the next one.

Figure 8 also shows that the above MAV trajectory fits well with the ground truth data, which proves the accuracy of both the SLAM algorithm

TABLE II: Parameter setup for the trajectory control using the setpoint method.

Parameter	h_s	h_l	d_s	ψ_s	t_s
Value	1.2m	0.3m	0.2m	15deg	0.2s

TABLE III: MAV pose estimation RMSEs of the whole trajectory, with position errors in *mm* and attitude errors in *degrees*

	x	y	z	3D	roll	pitch	yaw
RMSE	8.6	13.6	14.3	21.6	1.04	0.85	1.49

and its initialization module. Table III lists the RMSE of the on-board MAV-pose estimation when compared to the ground truth data. Position \mathbf{P}_1 is marked with a blue cross in Fig. 8b, \mathbf{P}_2 with a green cross. The initial position estimate of the landing site on the $x-y$ plane is marked with a blue square, and the final refined estimate with a green circle. Both position estimates are close to the ground truth data, which is marked with a black square. The blue and green crosses in Fig. 8a show the initial and final height estimate, comparing to the ground truth height marked with squares. With the landing site size being 500×500 (mm), the initial and final position estimation error is $(-19, -26, -6)^T$ (mm) and $(-11, -27, -5)^T$ (mm), respectively.

2) *Using trajectory-following method:* In this experiment, we compare the efficiency of the trajectory-following method and the setpoint method in autonomous navigation of our MAV. The poster pad is again used.

The trajectories of the MAV when using the two methods are shown in Fig. 9. The PID parameters of low level position and attitude controllers are the same in both cases. Trajectory controller parameters different from those in Sect. V-D.1 are shown in Table IV. Only when we use the method of trajectory following, we can explicitly set the MAV forward speed. However, we should note that it does not directly control the MAV speed. Instead, the forward speed will be fed to the position controller of the MAV.

When we calculate the RMSEs of the actual

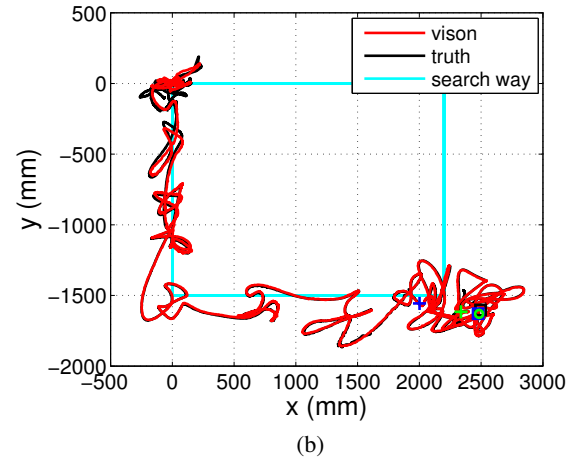
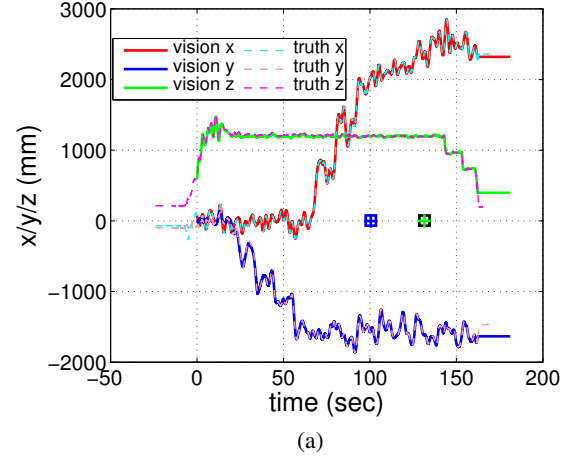


Fig. 8: (a) MAV position estimates from our vision system on x, y, z axis. (b) MAV trajectory projected to $x-y$ plane. The initial and final position estimates of the landing site and the associated MAV poses are also marked.

flight trajectory with respect to the predefined searching path on $x-y$ plane, position errors along the tangent direction of the searching path are ignored. The resulted RMSEs during the period of searching for the landing pad are listed in Table V. It shows the two methods result in similar precision for trajectory control under the parameter setup in Table IV, with the setpoint method performing a bit better. However, the trajectory-following method performs extremely well on yaw angle control. More importantly, this method is much more efficient regarding the spent flight time, which spends 32 seconds to obtain the initial pose of the landing pad. On the contrary, the setpoint

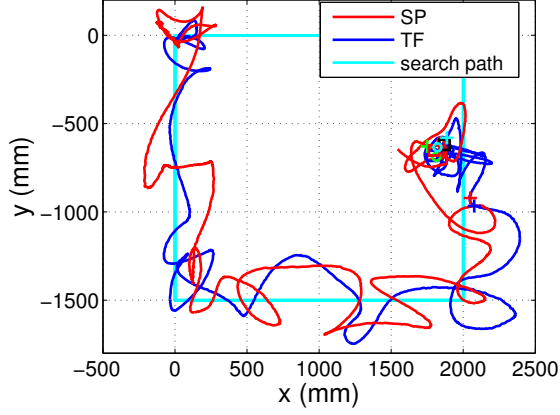


Fig. 9: MAV trajectories projected to $x-y$ plane, when using two different methods for navigation.

TABLE IV: Parameter setup for trajectory control using the two different methods.

Parameter	h_l	t_s	v_s
Value	0.4m	0.1s	0.35m/s

method spends 52.8 seconds to achieve this, even though we have set t_s smaller than that in Sect. V-D.1 to speed up the searching process. To improve the trajectory control precision, which will result in less efficiency in forward speed on the other hand, we suggest to distribute more setpoints along the predefined path and a relatively larger t_s in the setpoint method, while setting a smaller forward speed in the trajectory-following method.

When using the setpoint method, the detected landing pad positions and the corresponding MAV positions are marked in Fig. 9 in the same way as in Fig. 8b. For the result of using the trajectory-following method, the initially and finally detected landing pad positions are marked in red square and cyan circle, and the corresponding MAV positions are marked in red and cyan crosses, respectively.

E. Outdoor Experiment

In this experiment, we want to examine the robustness and extensibility of our vision system when working in outdoor environments. The experiment was carried out on a sunny midday with moderate wind, and the scenario is challenging in three aspects: Our MAV flies above a meadow, as

TABLE V: MAV trajectory control RMSEs, with RMSE1 for the setpoint method, and RMSE2 for the trajectory-following method.

	$x-y$	z	yaw
RMSE1	134.8mm	18.3mm	3.6deg
RMSE2	163.5mm	15.3mm	0.03deg



Fig. 10: The scenario for the outdoor experiment.

shown in Fig. 10. Thus, the first challenging aspect is that the environment is filled with visual features from grass with high self-similarity, which will be discussed in Sect. VI. The second one is that there was wind strong enough to bring disturbance to the pose controller of our MAV, which makes it a good scenario to test the robustness of the controller. The last one is introduced by the grass and other plants below the MAV: When the MAV flies above them in low altitude, their physical structure will be obviously changed by the wind produced by MAV propellers. This introduces more noise to our visual SLAM system.



Fig. 11: The reference image of the landing site in outdoor environment.

TABLE VI: MAV trajectory control RMSEs for the outdoor flight.

	$x - y$	z	yaw
RMSE	192.0mm	62.4mm	0.04deg

The landing site we defined is a piece of package paper with a size of $0.95m \times 0.95m$. Its reference image as shown in Fig. 11 was taken on another day by the same camera on our MAV, which provides a total number of 59 ORB features on the three pyramid levels. The trajectory-following method is used for trajectory control. Parameters for trajectory control are set in the same way as in Sect. V-D.2, except that we define d_s mentioned in Sect. V-B.2 to be $0.3m$. However, we define a longer searching path as $5m \times 5m$ with a height of $3.5m$.

Fig. 12 and Fig. 13 show the built map and MAV trajectory of an outdoor flight. Table. VI shows the RMSEs of the actual flight trajectory with respect to the predefined searching path. We can find that the flight performance is not much worse than in the indoor experiments despite the challenges we mentioned earlier. One reason for such results is that the MAV now flies much higher than in the indoor experiments, which leads to much less ground effect on the MAV.

Fig. 14 shows a scene when the SLAM system is initialized in this experiment. Although this time the assumption of all features lying on the same plane is no longer strictly true, it will not bring obvious effect to the pose tracking process, and those inaccurate map points will be later adjusted by local bundle adjustment.

VI. DISCUSSIONS AND FUTURE WORK

We mentioned three challenging issues in Sect. V-E, which make it difficult for our MAV to autonomously navigate in outdoor complex environments. Those challenges may finally cause pose tracking failure. Among them, the self-similarities of visual features are related to general vision systems. Fig. 15a shows a view of intermediate results from our visual SLAM system, which gives a self-similarity example. Such self-similarities

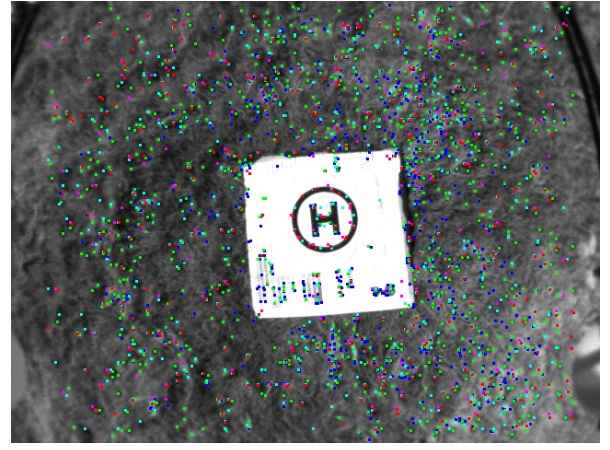


Fig. 14: Visual features in different image levels when the SLAM system is initialized in the outdoor environment.

may make the feature matching difficult, especially when features are in high density, resulting false map points to be triangulated. As can be found in Fig. 15b, most of those map points with relatively large altitude are false ones. A more robust feature matching strategies need to be investigated to tackle this issue. Another way to achieve better tracking robustness could be to augment the monocular SLAM system with more cameras looking at other directions. Thus more diverse vision features can be matched for localization. Illumination change in outdoor environments is another challenge for vision systems, which may easily cause cameras to be overexposed or underexposed, and result in tracking failure. We did our outdoor experiment in sunny daytime, since the illumination does not change much in this case. Thus, to finally go outdoor, we have to develop vision systems employing techniques which can efficiently handle illumination changes, e.g. the method presented in [15].

For an autonomous landing phase at the end of a long-term mission of a UAV, we propose to fuse IMU data to get its accurate short-term relative pose estimates, which can provide a metric scale constraint to initialize the SLAM system. Thus, autonomous searching for and landing on an arbitrary landing site could be achieved with a similar strategy as proposed in this paper. Although we have achieved promising results within relatively small areas, future work could be fusing IMU

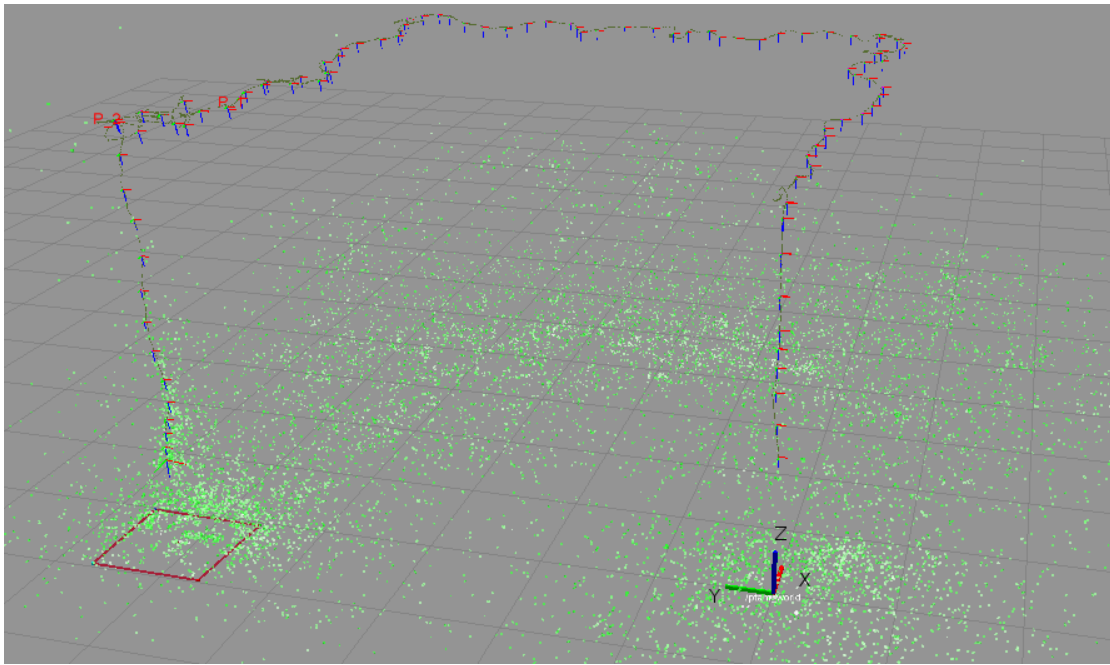


Fig. 12: Visualization of the built map and MAV trajectory during the outdoor searching and landing task. Map points are marked in green.

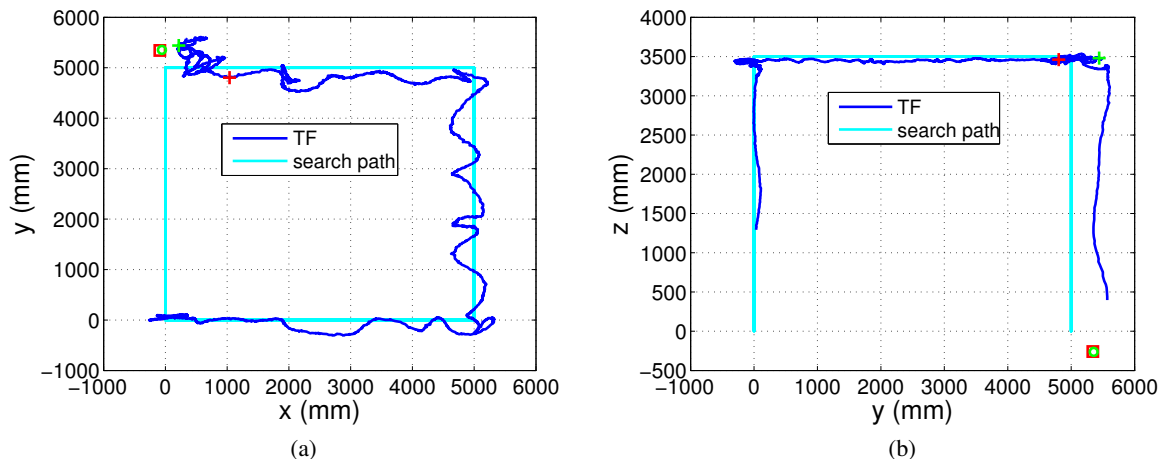


Fig. 13: (a) MAV trajectory projected to $x-y$ plane, and (b) to the $y-z$ plane. The initial and final position estimates of the landing site and the associated MAV poses are also marked.

data to extend the current monocular visual SLAM system to fulfill large scale tasks. This could not only be used to correct the pose estimates resulting from a SLAM system, but also to improve the localization and mapping accuracy of the SLAM system itself.

VII. CONCLUSIONS

In this paper we have presented a monocular vision system which enables an MAV to navigate autonomously in unknown environments, and search for the landing site on which it is designated to land. Our visual SLAM system can provide accurate pose estimates for the control of the MAV. We solve the landing site detection by integrating a

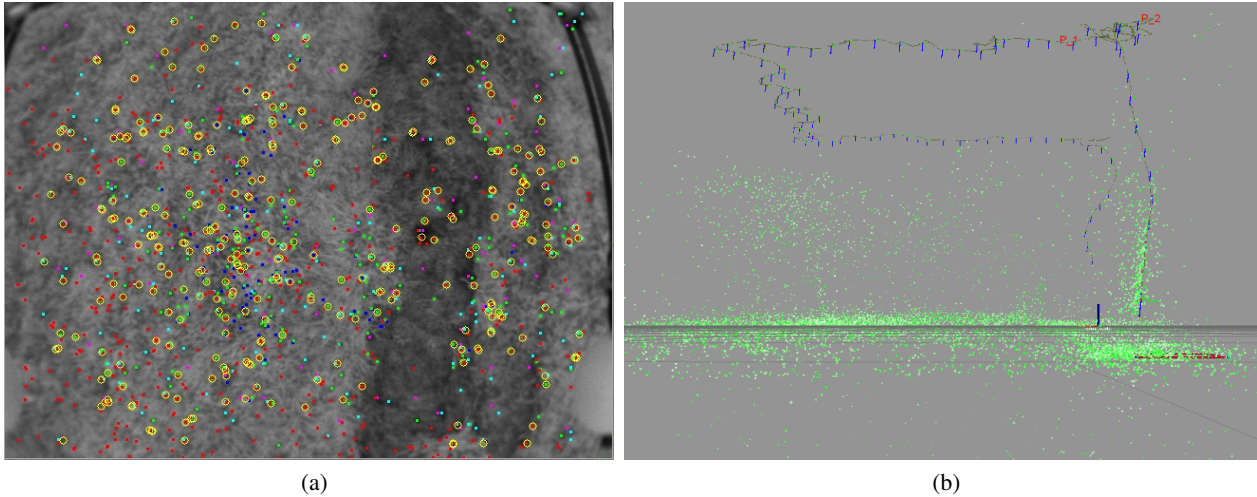


Fig. 15: (a) An intermediate result of the SLAM system when the MAV fly above a meadow. Those blue, green and cyan points are features in different image levels, and the red points are built map points, while those successfully matched map points are marked with yellow circles. (b) A side view of the final built map in the outdoor scenario.

multi-scale ORB feature matching scheme into the mapping thread of the SLAM framework. We further utilize the map points produced by the SLAM system to accurately estimate the 3D pose of the landing site, using a RANSAC-based method. No absolute scale information of the landing site is needed for its pose estimation.

By evaluating the pose estimation results of different landing sites, we show that our method is flexible and accurate enough for the proposed task of searching for and landing on an arbitrary landing site. Finally, we demonstrate our claims by the autonomous navigation and landing flights of our MAV indoor and in outdoor environment. The successful outdoor flight in the challenging scenario proves that our visual system can be extended to complex outdoor environments and enable MAVs to land autonomously. Video demonstration for the work presented in this paper can be found online at http://www.youtube.com/channel/UCQd6_G6qyvGHUmz7NUe1DZQ/videos.

REFERENCES

- [1] M. Achtelik, M. Achtelik, S. Weiss, and R. Siegwart, "On-board imu and monocular vision based control for mavs in unknown in- and outdoor environments", *Proceedings 2011 the IEEE International Conference on Robotics and Automation*, 2011.
- [2] M. Bloesch, S. Weiss, D. Scaramuzza, R. Siegwart, "Vision Based MAV Navigation in Unknown and Unstructured Environments", *Proceedings 2010 IEEE International Conference on Robotics and Automation*, Anchorage, AK, pp. 21–28, May 2010.
- [3] J.Y. Bouguet, "Camera Calibration Toolbox for Matlab", http://www.vision.caltech.edu/bouguetj/calib_doc, 2004.
- [4] R. O. Castle, G. Klein, D. W. Murray, "Combining monoSLAM with object recognition for scene augmentation using a wearable camera". *Image and Vision Computing*, 28(11), 1548–1556.
- [5] R. O. Castle, D. W. Murray, "Keyframe-based recognition and localization during video-rate parallel tracking and mapping". *Image and Vision Computing*, 29(8), 524–532. doi:10.1016/j.imavis.2011.05.002
- [6] G. Bradski, "The OpenCV Library", *Dr. Dobb's Journal of Software Tools*, 2000.
- [7] M. Calonder, V. Lepetit, C. Strecha, and P. Fua, "BRIEF: Binary robust independent elementary features", *In European Conference on Computer Vision*, 2010.
- [8] A. Cesetti, E. Frontoni, A. Mancini, P. Zingaretti, S. Longhi, "A Vision-Based Guidance System for UAV Navigation and Safe Landing using Natural Landmarks", *Journal of Intelligent & Robotic Systems*, 57(1–4), pp. 233–257, 2010.
- [9] A.J. Davison, "Real-time simultaneous localisation and mapping with a single camera", *Proc 9th IEEE Int Conf on Computer Vision*, pp. 1403–1410, 2003.
- [10] F. Fraundorfer, L. Heng, D. Honegger, G. Lee, P. Taniskanen, M. Pollefeys, "Vision-Based Autonomous Mapping and Exploration Using a Quadrotor MAV", *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Oct. 2012.
- [11] P. J. Garcia-Pardo, G. S. Sukhatme, J. F. Montgomery, "Towards vision-based safe landing for an autonomous helicopter", *Robotics and Autonomous Systems*, 38(1), pp. 19–29, 2002.

- [12] S. Grzonka, G. Grisetti, W. Burgard, "Towards a navigation system for autonomous indoor flying". *2009 IEEE International Conference on Robotics and Automation*, 2878–2883, 2009.
- [13] G. Klein and D. Murray, "Parallel tracking and mapping for small AR workspaces", in *International Symposium on Mixed and Augmented Reality*, pp. 225–234, Nov. 2007.
- [14] H. Lu, Z. Zheng, "Two novel real-time local visual features for omnidirectional vision", *Pattern Recognition*, 43(12), pp. 3938–3949, December 2010.
- [15] Huimin Lu, Hui Zhang, Shaowu Yang, and Zhiqiang Zheng, "Camera parameters auto-adjusting technique for robust robot vision," *Robotics and Automation (ICRA)*, 2010 IEEE International Conference on , pp.1518–1523, May 2010.
- [16] D.G. Lowe, "Distinctive image features from scale-invariant keypoints", *International Journal of Computer Vision*, 60(2), pp. 91–110, 2004.
- [17] R. Mahony, V. Kumar, P. Corke, "Multirotor Aerial Vehicles: Modeling, Estimation, and Control of Quadrotor", *Robotics & Automation Magazine, IEEE*, 19(3), pp. 20–32, Sept. 2012.
- [18] I. F. Mondragón, P. Campoy, C. Martínez, and M. A. Olivares-Méndez, "3D pose estimation based on planar object tracking for UAVs control", *2010 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 35–41, May, 2010.
- [19] L. Meier, P. Tanskanen, L. Heng, G.H. Lee, F. Fraundorfer, M. Pollefeys, "PIXHAWK: a micro aerial vehicle design for autonomous flight using onboard computer vision", *Autonomous Robots*, 33(1–2), pp. 21–39, 2012.
- [20] D. Mellinger, N. Michael, V. Kumar, "Trajectory generation and control for precise aggressive maneuvers with quadrotors", *The International Journal of Robotics Research*, 31(5), pp. 664–674, 2012.
- [21] E. Rublee, V. Rabaud, K. Konolige, G. Bradski, "ORB: an efficient alternative to SIFT or SURF", *2011 IEEE International Conference on Computer Vision (ICCV)*, Nov. 2011.
- [22] M. Quigley, B. Gerkey, K. Conley, J. Faust, T. Foote, J. Leibs, E. Berger, R. Wheeler, and A. Ng, "ROS: an open-source robot operating system", in *Open-source software workshop of the Int. Conf. on Robotics and Automation*, Kobe, Japan, 2009.
- [23] E. Rosten and T. Drummond, "Machine learning for high-speed corner detection", In *Proc. 9th European Conference on Computer Vision (ECCV'06)*, Graz, May 2006.
- [24] S. Saripalli, J.F. Montgomery, G.S. Sukhatme, "Visually guided landing of an unmanned aerial vehicle ", *IEEE Transactions on Robotics and Automation*, Vol. 19(3), 2003, pp. 371–380.
- [25] S.A. Scherer, D. Dube, P. Koma, A. Masselli, A. Zell, "Robust real-time number sign detection on a mobile outdoor robot", In *Proceedings of the 6th European Conference on Mobile Robots (ECMR 2011)*, Orebro, Sweden, September 2011.
- [26] S.A. Scherer, D. Dube, A. Zell, "Using Depth in Visual Simultaneous Localisation and Mapping", in *IEEE International Conference on Robotics and Automation (ICRA)*, pp. 5216–5221, 2012.
- [27] K. Schauwecker, N. R. Ke, S. A. Scherer, and A. Zell, "Markerless Visual Control of a Quad-Rotor Micro Aerial Vehicle by Means of On-Board Stereo Processing", in *Autonomous Mobile System Conference (AMS)*, Springer, pp. 11–20, 2012.
- [28] S. Shen, N. Michael, V. Kumar, "Autonomous multi-floor indoor navigation with a computationally constrained MAV". *2011 IEEE International Conference on Robotics and Automation*, 20–25, 2011.
- [29] S. Weiss, D. Scaramuzza, and R. Siegwart, "Monocular-SLAM-Based Navigation for Autonomous Micro Helicopters in GPS-Denied Environments", *Field Robotics*, 28(6), pp. 854–874, 2011.
- [30] S. Yang, S.A. Scherer, A. Zell, "An onboard monocular vision system for autonomous takeoff, hovering and landing of a micro aerial vehicle", *Journal of Intelligent & Robotic Systems*, 69(1–4), pp. 499–515, 2013.
- [31] S. Yang, S.A. Scherer, K. Schauwecker, A. Zell, "Onboard Monocular Vision for Landing of an MAV on a Landing Site Specified by a Single Reference Image", in *2013 International Conference on Unmanned Aircraft Systems (ICUAS'13)*, pp. 317–324, May, 2013.